

Machine Learning Methods for Discovering Metabolite Structures from Mass Spectra

by

Samuel Lucas Goldman

A.B. in Computer Science
Harvard College, 2019

Submitted to the Program of Computational and Systems Biology in Partial Fulfillment
of the Requirements for the Degree of
Doctor of Philosophy in Computational and Systems Biology

at the

Massachusetts Institute of Technology

February 2024

© Samuel Lucas Goldman. All rights reserved.

The author hereby grants to MIT a nonexclusive, worldwide, irrevocable, royalty-free license to exercise any and all rights under copyright, including to reproduce, preserve, distribute and publicly display copies of the thesis, or release the thesis under an open-access license.

Author.....
Computational and Systems Biology Graduate Program
January 12, 2024

Certified by.....
Connor W. Coley
Assistant Professor of Chemical Engineering
Assistant Professor of Electrical Engineering and Computer Science
Thesis Supervisor

Accepted by.....
Christopher Burge
Professor of Biology
Co-Director, Computational and Systems Biology Graduate Program

Machine Learning Methods for Discovering Metabolite Structures from Mass Spectra

by

Samuel Lucas Goldman

Submitted to the Program of Computational and Systems Biology
on January 12, 2024 in partial fulfillment of the requirements for the degree of

DOCTOR OF PHILOSOPHY

ABSTRACT

Small molecule metabolites mediate myriad biological and environmental phenomena across host-microbiome interactions, plant chemistry, cancer biology, and various other processes. Mass spectrometry is often used as an analytical technique to investigate the small molecules present in a sample, measuring both their masses and fragmentation spectra. However, the complexity and high dimensionality of spectral data makes it difficult to identify unknown metabolites and their roles, with a large majority of detected metabolites remaining unidentified in public data.

This thesis proposes a suite of new computational methodologies for higher accuracy annotation of small molecule metabolites from mass spectrometry data that integrate chemistry-informed priors with modern deep learning advancements. I begin by decomposing and framing the metabolite annotation pipeline into four key tasks well-fit for supervised deep learning including (A) molecular formula prediction, (B) spectrum-to-molecule property prediction, (C) molecule-to-spectrum prediction, and (D) *de novo* generation of molecular candidates. To address these various tasks, I first introduce the Molecular Formula Transformer to predict molecular property fingerprints from spectra by changing the tandem mass spectrum input basis from scalar mass values to plausible molecular formula annotations. This method is then extended to an energy-based-model formulation to predict the molecular formula of an unknown molecule from its tandem mass spectrum. Following these initial efforts to learn better representations of fragmentation spectra, I develop new neural networks capable of *generating* fragmentation spectra from small molecules through two-step autoregressive modeling. I show how this can be accomplished by generating either molecular formula peaks or molecular fragment peaks.

Downstream of metabolite prediction, a separate key question is to identify the function of discovered small molecules. To this end, I study and probe the ability to model enzyme-substrate compatibility from high throughput screens within a single enzyme family. In a final collaborative work, I further demonstrate how a new method for epistemic uncertainty quantification, evidential deep learning, can be applied to molecular property prediction. Altogether, this work outlines a path forward to a fully neuralized pipeline for the high throughput identification of small molecule metabolites and their functions.

Thesis supervisor: Connor W. Coley

Title: Assistant Professor of Chemical Engineering

Assistant Professor of Electrical Engineering and Computer Science

Acknowledgments

I am indebted to far more people than I can fully acknowledge in this space, both for supporting me throughout my graduate work and for preparing me to take on this endeavor.

To my advisor Connor, what a fun ride it has been. Thank you for your continual mentorship, friendship, and for taking a chance on me to help build out your then-nascent lab. You exemplify hard work, patience, uncompromising rigor, and kindness. I have full confidence your lab will do many great things in the years ahead.

Jing-Ke Weng, Matt Vander Heiden, Rafa Gómez-Bombarelli, and Regina Barzilay were continuously inspiring as members of my thesis committee. Jing-Ke and Matt were great thought partners in our adventure to consider various metabolomics applications. Rafa was a great collaborator in thinking through uncertainty quantification, and he brought enthusiasm and positivity to all our interactions that I hope to emulate. Regina served as an inspiration for how to build computational methods with real world impact; her words often propelled me to think more critically about my work and impact. I owe Jacquie Carota and Chris Burge additional gratitude for their support, especially in organizing the CSB PhD program. Several early scientific mentors were instrumental in giving me the confidence to push down this path including Gloria Coruzzi, Manpreet Katari, Philippe Cluzel, and Debora Marks.

Many colleagues have more directly shaped and influenced this thesis. Ava and Alexander Amini were early collaborators who trained me to be effective and organized in my work. Jeremy Wohlwend, Martin Strazar, Guy Haroush, and John Bradshaw have been wonderful collaborators in our forays into metabolomics. My time in the Coley Lab would not have been the same without upbeat and positive labmates including Priyanka Raghavan, Itai Levin, Rocío Mercado, Nicholas Casetti, Jenna Fromer, Thijs Stuyver, Wenhao Gao, Guangqi Wu, Keir Adams, Runzhong Wang, and David Graff. Several outstanding undergraduates have also contributed in various ways including Ria Das, Janet Li, and Jiayi Xin. Mingxun Wang was also unreasonably generous with his time and support for our metabolomics efforts.

Outside of research, the MIT Biotechnology Group, Nucleate, and MPM BioImpact, deserve special recognition. The friends I've made and experiences I've had in each of these organizations have made me a more thoughtful scientist and pushed me to consider the translational impact of research.

I am fortunate for my friends and housemates who have also inexplicably continued to tolerate me. Theo Diamandis, Alex Lenail, Axel Feldmann, Mark Tebbe, Tessa Buchan, and Jenny Walker joined me for many runs that always restored my sanity. Meilin Zhu and Constantine Tzouanas are at the center of some of my fondest memories in graduate school. Peter Bearse, Ben Milliken, and Tim O'Brien are less helpful for my sanity, but never fail to make me laugh.

These four years would of course not have been remotely the same if not for Louisa: thank you for adding levity to my life and reminding me to never take myself too seriously.

To my family, particularly Mom, Dad, my sister Sarah, Uncle Michael, Cousin Max, and Grandma Marilyn: your unconditional love and support means everything. This thesis—and all my accomplishments—are possible because of you.

Contents

| | |
|---|-----------|
| Titlepage | 1 |
| Abstract | 3 |
| Acknowledgments | 5 |
| List of Figures | 11 |
| List of Tables | 31 |
| 1 Introduction | 39 |
| 1.1 Metabolites: “Canaries of the Genome” | 39 |
| 1.2 Measuring the Metabolome | 40 |
| 1.3 Toward an AI/ML Analysis Pipeline | 41 |
| 1.4 Thesis Organization | 44 |
| 2 Predicting Molecular Properties and Substructures from Spectra | 47 |
| 2.1 Introduction | 47 |
| 2.2 Results | 49 |
| 2.2.1 The MIST method | 49 |
| 2.2.2 Formula transformers accurately predict fingerprints | 51 |
| 2.2.3 Contrastive fine-tuning improves metabolite retrieval | 54 |
| 2.2.4 Neural embeddings elucidate compound class clusters | 55 |
| 2.2.5 MIST uncovers putative dipeptides in clinical cohort data | 55 |
| 2.2.6 Public data benchmarking for future comparison to MIST | 60 |
| 2.3 Discussion | 60 |
| 2.4 Methods | 62 |
| 2.4.1 MIST architecture | 62 |
| 2.4.2 Fingerprints | 67 |
| 2.4.3 SIRIUS software | 67 |
| 2.4.4 Feed forward baseline | 67 |
| 2.4.5 Contrastive fine-tuning for database retrieval | 68 |
| 2.4.6 Model training | 68 |
| 2.4.7 Retrieval | 69 |
| 2.4.8 Training datasets | 69 |
| 2.4.9 Clinical data reanalysis | 70 |

| | | |
|----------|---|------------|
| 2.4.10 | Chemical classifications | 70 |
| 2.4.11 | Latent distance comparisons | 70 |
| 2.4.12 | Code and data availability | 71 |
| 2.5 | Additional Results | 71 |
| 2.5.1 | Additional model details | 71 |
| 2.5.2 | CSI:FingerID annotations for Mills et al. cohort | 71 |
| 2.5.3 | Dataset splits | 74 |
| 2.5.4 | Unfolding layer ablation | 74 |
| 2.5.5 | MAGMa substructure labels | 76 |
| 2.5.6 | Alternative transformer baselines | 77 |
| 2.5.7 | Spectra noising | 78 |
| 2.5.8 | Hyperparameter optimization | 79 |
| 3 | Inferring the Precursor Mass Molecular Formula from Spectra | 89 |
| 3.1 | Introduction | 89 |
| 3.2 | Methods | 93 |
| 3.2.1 | Preliminaries | 93 |
| 3.2.2 | An energy-based model for formula annotation | 94 |
| 3.2.3 | The MIST-CF architecture | 95 |
| 3.2.4 | Generating formula candidates | 97 |
| 3.2.5 | Downselecting candidate formulae with a learned filter | 97 |
| 3.2.6 | Datasets | 98 |
| 3.2.7 | Baseline models | 98 |
| 3.3 | Results | 99 |
| 3.3.1 | FastFilter constrains candidate formulae with high precision | 99 |
| 3.3.2 | Molecular Formula Transformers provide meaningful representations | 100 |
| 3.3.3 | MIST-CF compares favorably to existing formula annotation tools | 103 |
| 3.3.4 | MIST-CF is competitive on the CASMI2022 challenge | 104 |
| 3.4 | Conclusion | 105 |
| 3.5 | Data and Software Availability | 106 |
| 3.6 | Additional Results | 106 |
| 3.6.1 | SIRIUS configuration | 106 |
| 3.6.2 | Timing experiments | 107 |
| 3.6.3 | Hyperparameter setting | 107 |
| 4 | Generating Mass Spectra as Subformula Sets with Prefix-Trees | 109 |
| 4.1 | Introduction | 109 |
| 4.2 | Background | 111 |
| 4.2.1 | Tandem mass spectrometry | 111 |
| 4.2.2 | Predicting mass spectra from molecules (spectrum predictors) | 111 |
| 4.2.3 | Mass spectrum libraries | 112 |
| 4.3 | Model | 113 |
| 4.3.1 | SCARF-Thread: Generating product formulae using prefix trees | 114 |
| 4.3.2 | SCARF-Weave: Predicting intensities given product formulae | 115 |
| 4.3.3 | Training and inference | 116 |

| | | |
|----------|---|------------|
| 4.4 | Experiments | 116 |
| 4.4.1 | Dataset | 116 |
| 4.4.2 | Spectra prediction | 117 |
| 4.4.3 | Retrieval | 119 |
| 4.5 | Related Work | 120 |
| 4.6 | Conclusion | 121 |
| 4.7 | Additional Results | 122 |
| 4.7.1 | Extended results | 122 |
| 4.7.2 | Dataset preparation | 124 |
| 4.7.3 | Baselines | 125 |
| 4.7.4 | Retrieval subsets vs. PubChem | 130 |
| 4.7.5 | Model details | 130 |
| 4.7.6 | Limitations and future work | 132 |
| 5 | Generating Mass Spectra as Substructure Sets by Breaking Bonds | 137 |
| 5.1 | Introduction | 137 |
| 5.2 | Methods | 140 |
| 5.2.1 | Datasets | 140 |
| 5.2.2 | Canonical DAG construction | 140 |
| 5.2.3 | Model details | 141 |
| 5.2.4 | Baselines | 143 |
| 5.3 | Results | 143 |
| 5.3.1 | ICEBERG is trained as a two-stage generative and scoring model | 143 |
| 5.3.2 | ICEBERG enables highly accurate spectrum prediction | 146 |
| 5.3.3 | Model predictions are interpretable | 148 |
| 5.3.4 | Fragmentation simulations lead to improved structural elucidation | 149 |
| 5.3.5 | Challenging test splits better explain retrieval performance | 150 |
| 5.4 | Discussion | 151 |
| 5.5 | Additional Results | 151 |
| 5.5.1 | Graph neural network details | 151 |
| 5.5.2 | Additional baseline details | 152 |
| 5.5.3 | Hyperparameters | 156 |
| 6 | Predicting Enzyme-Substrate Compatibility | 161 |
| 6.1 | Introduction | 161 |
| 6.2 | Results | 164 |
| 6.2.1 | Data summary | 164 |
| 6.2.2 | Models | 165 |
| 6.2.3 | Enzyme discovery | 166 |
| 6.2.4 | Substrate discovery | 168 |
| 6.2.5 | Reanalysis of kinase inhibitor discovery | 170 |
| 6.2.6 | Improving enzyme discovery models | 172 |
| 6.3 | Discussion | 174 |
| 6.4 | Conclusion | 174 |
| 6.5 | Methods | 174 |

| | | |
|----------|--|------------|
| 6.5.1 | Dataset preparation | 174 |
| 6.5.2 | Hyperparameter optimization | 175 |
| 6.5.3 | Evaluation metrics | 175 |
| 6.5.4 | Filtering imbalanced tasks | 175 |
| 6.5.5 | Kinase inhibitor reanalysis | 176 |
| 6.5.6 | Pooling strategies | 176 |
| 6.5.7 | Active site pooling | 176 |
| 6.6 | Additional results | 177 |
| 6.6.1 | Data | 177 |
| 6.6.2 | Models | 181 |
| 6.6.3 | Extended results | 183 |
| 7 | Evidential Molecular Prediction and Discovery | 195 |
| 7.1 | Introduction | 195 |
| 7.2 | Approach | 197 |
| 7.2.1 | Formulating evidential learning for molecules | 197 |
| 7.3 | Results | 198 |
| 7.3.1 | Uncertainty benchmarking | 198 |
| 7.3.2 | Calibration and tunability | 201 |
| 7.3.3 | Application I: Uncertainty-guided learning | 203 |
| 7.3.4 | Application II: Uncertainty-guided inference for virtual screening | 206 |
| 7.4 | Discussion | 208 |
| 7.4.1 | Contributions | 208 |
| 7.4.2 | Advantages of evidential learning in chemical science | 209 |
| 7.4.3 | Opportunities and applications in molecular property prediction | 209 |
| 7.4.4 | Scope and future work | 210 |
| 7.5 | Conclusion | 210 |
| 7.6 | Methods | 211 |
| 7.6.1 | Evidential deep learning formulations | 211 |
| 7.6.2 | Network architectures | 212 |
| 7.6.3 | Datasets | 213 |
| 7.6.4 | Uncertainty quantification baselines | 214 |
| 7.6.5 | Uncertainty benchmarking experiments | 215 |
| 7.6.6 | Active learning on QM9 dataset | 217 |
| 7.6.7 | Bayesian optimization on docking dataset | 218 |
| 7.6.8 | Uncertainty-guided virtual screening for antibiotic discovery | 219 |
| 7.7 | Additional Results | 220 |
| 8 | Conclusion | 229 |
| 8.1 | Contributions | 229 |
| 8.2 | Future Work | 230 |
| 8.2.1 | Building deployment-ready systems | 230 |
| 8.2.2 | Toward new discovery paradigms | 231 |
| | References | 233 |

List of Figures

| | | |
|-----|--|----|
| 1.1 | Spectrum prediction workflow tasks. Raw data is first processed into a mass spectrum with a precursor MS1 mass. Subsequently, a molecular formula can be predicted. Beyond this, three separate classes of models can be used to help identify the molecule structure, and the predictive direction of each model is shown. | 44 |
| 2.1 | MIST replaces the critical spectrum-to-fingerprint prediction step in the computational workflow. A. A set of fragmentation peaks are extracted from a sample with tandem mass spectrometry. After predicting molecular property fingerprints, several different downstream software tools can be used to annotate molecules from databases, link metabolites to genomes, <i>de novo</i> generate plausible candidates, classify the unknown metabolite into compound classifications, generate statistics for all observed spectra, and calibrate certainty of annotations. B. The first step in the fingerprint prediction pipeline common to MIST and CSI:FingerID is to annotate the spectrum with a molecular formula for the full metabolite, then label each sub-peak with a molecular formula. C. CSI:FingerID arranges annotated sub-peak formulae into a "fragmentation tree", and subsequently uses many independent support vector machine models to predict fingerprint property bits. D. MIST does not require a fragmentation tree, but rather directly utilizes molecular formulae as inputs to a Formula Transformer in order to predict molecular fingerprints. | 48 |
| 2.2 | MIST accurately predicts compound fingerprints from mass spectra. A. Overview of MIST model architecture. An input spectrum of (<i>mass-to-charge, intensity</i>) pairs is transformed into molecular formulae vectors, encoded, and fed into a Molecular Formula Transformer along with pairwise neutral losses between formulae. An unfolding module predicts full molecular fingerprints and a separate auxiliary module predicts substructure fingerprints as a secondary training signal. B,C. Molecular fingerprints are predicted by MIST and CSI:FingerID for every spectrum in the test set. The performance for each spectrum by log likelihood to the true fingerprint (B) or cosine similarity (C) is evaluated and plotted. Points below the line represent instances where MIST is more performant. D. Equivalent evaluation showing the likelihood of predicting each fingerprint bit correctly across all spectra. | 50 |

| | | |
|-----|--|----|
| 2.2 | <p>E. The performance of CSI:FingerID [53], MIST, and FFN, a baseline inspired by MetFID [74], are shown. "Tanimoto," "Cosine," and "Log likelihood (spectra)" indicate performance across spectra and "Log likelihood (bits)" indicates performance across bits (higher is better). Median lines are shown; boxes show the interquartile range, with whiskers indicating 1.5x interquartile ranges. Statistics are pooled across $n = 18,700$ test spectra (Tanimoto, Cosine, and Spectra log likelihood) and $n = 5,496$ fingerprint bits (Bit log likelihood). F. Performance differences by compound class as assigned by NPCClassifier. All classes with > 40 molecules are shown. G-I. Example molecules from stilbenoid, aromatic polyketide, and sphingolipid classes. All results indicate predictions aggregated across 3 separate independent splits and model re-trainings. MIST fingerprints are created by averaging predictions from 5 separately trained models per split.</p> | 51 |
| 2.3 | <p>Contrastive fine-tuning improves compound annotation by database retrieval. A. Overview of the contrastive fine-tuning workflow for MIST. Model weights in the formula embedding module and attention layers are initialized from the fingerprint prediction tasks ("Pre-trained base model") and fine-tuned to project fingerprints into the a latent hidden representation space shared with spectra. The model minimizes distance between the spectra representations and their true structure’s representation and maximizes distance to decoy small molecules. B,C. MIST top-k retrieval accuracy assessed against CSI:FingerID by querying the PubChem reference database using predicted fingerprint cosine distances, contrastive latent space distances, or a custom CSI:FingerID fingerprint distance function ("Bayes"). Top-k accuracy is computed by ranking candidate compounds for each spectra and computing the fraction of spectra for which the true compound appears in the top-k rankings, aggregated across 3 test set data folds.</p> | 53 |
| 2.3 | <p>D. Retrieval accuracy for all folds are segmented into various chemical classes. For each class for which > 40 examples exist, the fraction of examples on which MIST performs better, equivalently, or worse than CSI:FingerID are shown. E. MIST latent space spectra distances more effectively cluster high similarity compounds than MS2DeepScore [82], Spec2Vec [81], or cosine distances. Test set compounds pairs are sorted by similarity according to various metrics and the average structural similarity of pairs is computed using Tanimoto distance at various percentiles for a single test fold. F. Contrastive spectra latent vectors for a single test fold are projected into 2D space with UMAP and colored by their respective compound class. Compounds 2 and 4 are example spectra for which NPCClassifier fails to return an annotation. G. Example compounds from (F) reveal structural similarity of neighboring compounds in the UMAP space derived from MIST’s spectral embeddings. Maximal Tanimoto similarity to any example in the training set is shown in parentheses. 54</p> | 54 |

| | | |
|-----|--|----|
| 2.4 | <p>MIST annotates putative and clinically relevant dipeptides. A. Metabolite classes ranked by their Pearson correlation with IBD disease severity. All putative metabolites are assigned disease classes with NPClassifier [101] and relative abundances are summed across classes. Within the ulcerative colitis (UC) and Crohn’s disease (CD) cohorts, metabolite classes are correlated with disease severity and sorted according to their correlation with UC (top 20 shown; each computed for n=60 patients with UC, n=102 CD patients). Error bars indicate 95% confidence intervals around the mean. B. Dipeptide relative abundances for each patient are plotted against disease severity within the healthy (n=19), UC, and CD cohorts. Dipeptides correlate with disease severity for UC patients ($R^2 = 0.23$, $p = 0.01$, two-sided the Wald test, adjusted with Benjamini-Hochberg) as observed by Mills et al. [102]. In the left boxplot, the data median line is shown; boxes show the interquartile range, with whiskers indicating 1.5x interquartile ranges. Regression lines are shown with 95% confidence intervals computed with bootstrapping. C. Within the Dipeptide class, the distribution of individual metabolite correlations are shown for both UC and CD cohorts.</p> | 57 |
| 2.4 | <p>D. All putative dipeptide metabolites (red) are embedded into the latent space of a single, contrastive MIST model alongside embeddings for reference standard spectra from the original model training set and also labeled as dipeptides (blue). All embeddings are projected in 2D space with UMAP and plotted. Metabolites are colored by their coefficient of determination R^2 with UC disease severity as computed in C. E. Annotations for example metabolites from D are shown across different compound clusters. Metabolites with molecular formulae not in standards library are colored brown. F-G. Example spectra and their MIST annotated compound are shown. Maximal Tanimoto similarity to any example in the training set is shown in parentheses. Explained subpeaks are annotated using CFM-ID [51] and shown as validation of plausibility. Compound annotations are made using an ensemble of 5 MIST models, contrastive distance retrieval, and the HMDB and PubChem reference databases of molecules.</p> | 58 |

2.5 **Putative alkaloids separate healthy and diseased IBD cohorts. A,B.** Putative metabolite classes from the Mills et al. cohort [102] are scattered showing their fold change and respective p value comparing and ulcerative colitis (UC) (**A**) and Crohn’s disease (CD) (**B**) cohorts to the healthy cohort. **C.** Distributions of patients’ relative abundances for piperidine and pyridine alkaloids classes are shown for healthy, UC, and CD groups (n=19 healthy, n=22 UC, n=38 CD). Median lines are shown; boxes show the interquartile range, with whiskers indicating 1.5x interquartile ranges. **D.** The total number of putative metabolites in both significant classes are shown, highlighting the number of metabolites with putative novel structures not observed in the standards library. **E,F.** Individual metabolite abundance fold changes and statistical significance are shown comparing the healthy cohort to UC (**E**) and CD cohort (**F**). **G.** Select putative annotated molecule structures are shown for compounds **15** (isonicotinic acid, HMDB:0060665), **16** (triacetonamine, HMDB:0031779), **17** (piperettine, HMDB:0034371), and **18** (unknown). Spectra **15** and **16** are differentially less abundant for both UC and CD. Novel compounds indicate those without spectra standards and are labeled brown; compounds with respective standards are shown in blue. Maximal Tanimoto similarity to any example in the training set is shown in parentheses. Compound annotations are made using an ensemble of 5 MIST models, contrastive distance retrieval, and the HMDB and PubChem reference database of molecules. The top chemical class annotation is found by running the top putative prediction through NPClassifier [101]. P values were computed using independent two-sided t-tests and adjusted for multiple hypothesis testing with the Benjamini Hochberg method. To have equal cohort sizes for healthy and diseased patients, UC and CD cohorts were subsetted to patients with disease severity score > 0.2. Blue scattered points indicate differentially less abundant compounds or classes; red indicates differentially increased abundance.

| | | |
|-----|---|----|
| 2.6 | Model ablations affirm the value of domain-inspired model components. A,B. MIST performance on fingerprint prediction compared to ablated variants of MIST without pairwise interactions, without MAGMa substructure supervision, using no simulated spectra, without fingerprint unfolding and using only a feed forward network (FFN) using metrics of (A) cosine similarity and (B) log likelihood (n=2,438 test spectra). C,D. MIST fingerprint prediction accuracy improves as a function of dataset size for both (C) cosine similarity and (D) log likelihood (n=819 test spectra). E. MIST retrieval performance using a weighted sum of contrastive and fingerprint distance outperforms contrastive distance or fingerprint distance alone, a contrastive version of the model without finetuning, a FFN contrastive model, or FFN fingerprints. F. Full MIST contrastive + fingerprint retrieval accuracy improves with training set size. All results are computed on the public GNPS subset data released with CANOPUS. Model ablations are conducted with 3 random splits of the data by molecular structure; dataset ablations are conducted for a single split of the data. Log likelihood values are clamped to a minimum of -5 . Unless otherwise stated, error bars show 95% confidence intervals for the mean. | 61 |
| 2.7 | Unfolding fingerprint prediction module. A. The module starts with a hidden representation and subsequently predicts increasing resolutions of the fingerprint. B. At each intermediate fingerprint prediction $\hat{y}^{(i)}$, a corresponding target vector is generated by "folding" the previous, higher resolution fingerprint. An intermediate loss is calculated at each of these resolutions. | 72 |
| 2.8 | Feed forward network binned spectra baseline. An input spectrum is binned, a multilayer perceptron (MLP) is applied, and the output is projected to a fingerprint prediction. | 72 |
| 2.9 | Forward simulation of spectra from molecules. A. An input spectrum, molecule pair is first denoised and cleaned by labeling chemical sub-formulae. This is converted into a binned spectra representation. B. To generate new model training examples, molecules are fingerprinted, binned spectra predictions are predicted by a neural network, and sub-formula are assigned to each of the bins such that the spectra can be used as input to MIST. C. The forward simulation neural architecture utilizes unfolding layers. A multilayer perceptron (MLP) projects the molecular fingerprint into a hidden representational space. Forward-reverse models, G , are used to progressively grow the binned prediction such that resolution increases at each step, similar to unfolding layers. D. Illustration of the forward-reverse module, where intensities (z_f) and neutral losses (z_r) are predicted. Neutral losses are reversed and mapped onto intensity bins using the full mass as input then summed according to a learned gate to get a single prediction of intensities in each bin. | 73 |

- 2.10 **Repeated analysis of clinically relevant dipeptides with the PubChem retrieval library.** **A.** Metabolite classes ranked by their Pearson correlation with IBD disease severity. All putative metabolites are assigned disease classes with NPClassifier [101] and relative abundances are summed across classes. Within the ulcerative colitis (UC) and Crohn’s disease (CD) cohorts, metabolite classes are correlated with disease severity and sorted according to their correlation with UC (top 20 shown; each computed for n=60 patients with UC, n=102 CD patients). Error bars indicate 95% confidence intervals around the mean. **B.** Dipeptide relative abundances for each patient are plotted against disease severity within the healthy (n=19), UC, and CD cohorts. Dipeptides correlate with disease severity for UC patients ($R^2 = 0.20$, $p = 0.02$, two-sided the Wald test, adjusted with Benjamini-Hochberg) as observed by Mills et al. [102]. In the left boxplot, the data median line is shown; boxes show the interquartile range, with whiskers indicating 1.5x interquartile ranges. Regression lines are shown with 95% confidence intervals computed with bootstrapping. **C.** Within the dipeptide class, the distribution of individual metabolite correlations are shown for both UC and CD cohorts. Compound annotations are made using an ensemble of 5 MIST models, contrastive distance retrieval, and PubChem reference databases of molecules. Input spectra are subsetted to be under 500 Da, have $[M+H]^+$ adducts, and utilize molecular formula as output by CSI:FingerID at the end of the compound identification step. 75
- 2.11 **Reanalyzing putative metabolite class differences between healthy and diseased IBD cohorts using the PubChem database.** **A,B.** Putative metabolite classes from the Mills et al. cohort [102] are scattered showing their fold change and respective p value comparing and ulcerative colitis (UC) (**A**) and Crohn’s disease (CD) (**B**) cohorts to the healthy cohort. P values were computed using independent two-sided t-tests and adjusted for multiple hypothesis testing with the Benjamini Hochberg method. To have equal cohort sizes for healthy and diseased patients, UC and CD cohorts were subsetted to patients with disease severity score > 0.2 . Blue scattered points indicate differentially less abundant compounds or classes; red indicates differentially increased abundance. **C.** Distributions of patients’ relative abundances for piperidine and pyridine alkaloid classes are shown for healthy (n=19), UC (n=22), and CD (n=38) groups. Median lines are shown; boxes show the interquartile range, with whiskers indicating 1.5x interquartile ranges. All compound annotations are made using an ensemble of 5 MIST models, contrastive distance retrieval, and the PubChem reference database of molecules. The top chemical class annotation is found by running the top putative prediction through NPClassifier [101]. Input spectra are subsetted to be under 500 Da, have $[M+H]^+$ adducts, and utilize molecular formula as output by CSI:FingerID at the end of the compound identification step. 76

| | | |
|------|---|----|
| 2.12 | Dataset split similarity. Distribution plot of the most similar molecule in the training set according to Tanimoto similarity for molecules in the train and molecules in the test sets. Distributions are computed on the public GNPS CANOPUS dataset (A.) and private proprietary CSI dataset (B.). | 77 |
| 2.13 | Labeling data substructures. A. MIST is trained to build representations of each peak in the spectrum. The model uses the precursor peak representation to predict a fingerprint of the full molecule, and when substructures can be labeled for other peaks, the model learns to predict these as well. B-D. Example substructure peak labels in the training set computed with MAGMa. The full compound is outlined in red. Other substructures for the top 6 highest intensity peaks are outlined in black. | 78 |
| 2.14 | Single MIST model comparison to CSI:FingerID. Molecular fingerprints are predicted by MIST and CSI:FingerID for every spectrum in the test set using a single MIST model as in Figure 2. The performance for each spectrum by cosine similarity to the true fingerprint (A) or log likelihood (B) is evaluated and plotted. Points below the line represent instances where MIST is more performant. C. Equivalent evaluation showing the likelihood of predicting each fingerprint bit correctly across all spectra. D. The performance of CSI:FingerID, MIST, and FFN, a baseline inspired by MetFID, are shown. "Tanimoto," "Cosine," and "Log likelihood (spectra)" indicate performance across spectra and "Log likelihood (bits)" indicates performance across bits (higher is better). Median lines are shown; boxes show the interquartile range, with whiskers indicating 1.5x interquartile ranges. Statistics are pooled across $n = 18,700$ test spectra (tanimoto, cosine, and spectra log likelihood) and $n = 5,496$ fingerprint bits (bit log likelihood). All results are aggregated across 3 splits of the data. | 80 |
| 2.15 | Dataset ablations on the CSI2022 dataset validate that model performance is limited by data set size. MIST fingerprint prediction accuracy improves as a function of dataset size in terms of both (A) cosine similarity and (B) log likelihood. Model performance is computed on the partially proprietary CSI2022 dataset. Dataset ablations are conducted for a single split of the data. Log likelihood values are clamped to a minimum of -5 . Error bars show 95% confidence intervals for the mean ($n=3,116$ test spectra). For this ablation, MIST models are trained without simulated data and without MAGMa auxiliary loss to reduce model training time. | 81 |
| 2.16 | Example molecules and their training precedents for Figure 2.3. Each molecule is shown along side its nearest neighbor in the training set, with its Tanimoto similarity listed. | 86 |
| 2.17 | Example molecules and their training precedents for Figure 2.4. Each molecule is shown along side its nearest neighbor in the training set, with its Tanimoto similarity listed. Even in cases for which the proposed molecular structure appears in the training set, it is predicted on the basis of a new experimental spectrum that does not exactly match the reference spectrum. | 87 |

| | | |
|------|--|----|
| 2.18 | Example molecules and their training precedents for Figure 2.5. Each molecule is shown along side its nearest neighbor in the training set, with its Tanimoto similarity listed. Even in cases for which the proposed molecular structure appears in the training set, it is predicted on the basis of a new experimental spectrum that does not exactly match the reference spectrum. | 88 |
| 3.1 | MIST-CF and SIRIUS both address the MS/MS molecular formula annotation problem. A. Input samples are first processed, recording a tandem mass spectrum. Before assigning full molecular structures, the molecular formula can first be inferred to constrain structure assignment. B. Methodological similarities and differences between MIST-CF and SIRIUS. A candidate precursor mass is first decomposed into plausible molecular formulae and adduct pairs. MIST-CF (left) learns in a data-driven fashion to assign scores and circumvents the need of fragmentation tree construction compared to SIRIUS (right). Both methods rely on assigning subformulae (“subforms”), which are molecular formula subsets of the candidate precursor formula that match individual MS/MS peak masses. | 92 |
| 3.2 | Large numbers of formula candidates can be quickly filtered with a simple feed forward neural network, FastFilter. A. Generated candidate formulae for a spectrum can be prioritized with a learned model. A scalar score is generated for each candidate formula given an MS1 value and mass tolerance via a feed forward neural network. Only the top ranked formulae are selected for consideration with MIST-CF. B. The distribution of the number of candidates for all recorded MS1 spectra in our spectra libraries NPLIB1 and NIST20. All candidate molecular formulae are generated for all 6 adduct types considered with an allowed mass deviation of 10 ppm. C. FastFilter model accuracy at various top k cutoff values. The top formula is nearly always recovered within the top 256 candidates. The 95% confidence interval of the mean for recovery is shown across the 3 different formula splits considered in this work. All results are computed for a single trained FastFilter model on a large database of biologically-relevant molecules. | 99 |

| | | |
|-----|--|-----|
| 3.3 | <p>MIST-CF is a highly-effective architecture for learning to rank plausible MS1 formula annotations A. The MIST-CF architecture uses the candidate molecular formulae to generate subformulae and encode these with a Formula Transformer. B. The baseline feed forward network (FFN) separately encodes the spectra and formula before feeding their concatenation into a multilayer perceptron (MLP). C. For all spectra in the test set, the fraction recovered at various top k values for all methods is computed and shown. D-E. The top 1 accuracy for the methods is grouped by the mass of the MS1 precursor or adduct type. All results are computed for MIST-CF and the following baselines: MS1 only (a feed forward network utilizing only the molecular formula and context vector), a feed forward network, and a Transformer. All results are computed over 3 random formula splits and respective training runs as described in Section 3.2, where the NIST20 is included in the training sets. All spectra include up to 256 candidates to select from as selected with the “COMMON” filter [50] and FastFilter. Error bars and shaded regions show 95% confidence intervals of the mean.</p> | 100 |
| 3.4 | <p>MS/MS subpeaks drive MIST-CF performance. A. The (sub)Formula Transformer in MIST-CF operates on only the N_p highest intensity MS2 peaks. Here, from that set of labeled MS2 peaks, the number of peaks used as input to the spectrum is limited. As the number of included peaks increases, the formula retrieval accuracy does too. B. The top 1 retrieval accuracy is shown for all maximum subpeak numbers. All results are shown for MIST-CF trained on the joint NPLIB1 and NIST20 dataset and tested on a single test split of the data.</p> | 102 |
| 3.5 | <p>MIST-CF more frequently assigns correct molecular formulae than the SIRIUS formula assignment module on the NPLIB1 test set. A. The number of spectra for which each method is able to predict a formula—regardless of accuracy. B. The distribution of molecular masses for the spectra both methods are able to annotate—regardless of accuracy— compared to the distribution of molecular masses for spectra on which only MIST-CF succeeds. C. Top k accuracy for both methods is shown. D. The rank at which each method is able to recover the true molecular formula. This visualization highlights that there are many spectra for which MIST-CF achieves rank 1 that SIRIUS does not predict in its top 3. MIST-CF models are trained on the joint NIST20 and NPLIB1 dataset. SIRIUS is executed with a compound and tree timeout of 300 seconds. All values are shown for 3 random formula splits of the NPLIB1 data. Error bars and confidence intervals show 95% confidence intervals for the standard error of the mean. Accuracy is computed with respect to the <i>total</i> MS1 composition (i.e., summed adduct and formula) to avoid biasing results against SIRIUS.</p> | 103 |
| 4.1 | <p>Tandem mass spectrometers measure fragmentation patterns of molecules, resulting in characteristic peaks that are indicative of their structure. SCARF simulates these fragmentation patterns <i>in silico</i>.</p> | 110 |

| | | |
|-----|--|-----|
| 4.2 | Overview of various approaches to spectrum prediction. A. Fragmentation prediction approaches use heuristics and scoring rules to break down the molecule into fragments and their associated intensities. B. Binned prediction approaches discretize the possible mass-to-charge values and predict intensities for each possible bin. C. Formula prediction approaches predict spectra as sets of molecular formulae and intensities. Our model SCARF utilizes a two stage approach, first by predicting the product formulae present (constrained by the precursor formula), which defines the x-axis locations of the peaks, before secondly assigning intensities to these formulae (defining the peaks’ y-axis values). | 112 |
| 4.3 | Illustration of the SCARF-Thread architecture. A. The formulae of the product fragments can be represented using a prefix tree. SCARF-Thread predicts this tree for new molecules at test time. It does so by expanding each node at a given depth in parallel, treating the counts of subsequent elements as dependent only on the counts of elements predicted so far (i.e., the prefix) and the original molecular structure. B. The SCARF-Thread predictive task at the C ₇ node from the prefix tree diagram shown in A. Here the network takes as input (i) an embedding of the overall molecule; (ii) a vector representing the counts of each element in the prefix so far (counts yet to be predicted are represented using a special token), (iii) the difference of the counts predicted so far from the precursor molecule, and (iv) a one-hot representation of the element for which the counts are currently being predicted. The network predicts which counts are valid next nodes in the prefix tree (where counts that are greater than those in the original precursor molecular formula are automatically masked out as invalid). See also Algorithm 9. | 114 |
| 4.4 | The SCARF-Weave network, which takes in the product formulae (e.g., predicted by SCARF-Thread) and predicts their intensities. We use a Set Transformer architecture [91], such that our model takes in the details of the other product formulae present when predicting intensities. . . | 116 |
| 4.5 | SCARF enables more accurate retrieval of ground truth molecules within the NIST20 dataset. A. Average retrieval accuracy of SCARF at various top k thresholds. Retrieval is conducted on the same test split, and retrieval accuracy is averaged across models trained for three separate random seeds. B-C. Example spectrum predictions made by SCARF (top) compared to the ground truth spectrum (bottom). Up to 5 predicted peaks with the highest intensity are annotated with their molecular formula explanation as predicted by SCARF. The full molecule is shown inset. Further examples are in the Appendix (Figure 4.6). | 120 |
| 4.6 | Example spectra predictions from the NIST20 dataset for 10 randomly selected test molecules. The ground truth spectra are shown underneath in black, with predictions above in teal. Molecules are shown inset. | 123 |

| | | |
|-----|--|-----|
| 4.7 | Cosine similarity of predicted spectra stratified by properties. Results are stratified across molecular weight for both NPLIB1 (A) and NIST20 (B). We further stratify results across putative chemical classes of input molecules using NPClassifier [101] for both NPLIB1 (C) and NIST20 (D). The dotted line indicates the average predictive cosine similarity of SCARF across all examples and averaged over three random splits. | 126 |
| 4.8 | Spectra dataset molecule characterizations. A. Distribution of the molecular weight of compounds across NPLIB1 and NIST20. B-C. Chemical classes contained in NPLIB1 (B) and NIST20 (C) with the top 15 classes shown and all others grouped in ‘Other’. Chemical classes are computed using NPClassifier [101]. | 126 |
| 5.1 | ICEBERG enables the prediction of tandem mass spectra by efficiently navigating the space of possible fragmentation events. A. Example experimental mass spectrum. An input molecule, benzocaine, is depicted entering a mass spectrometer collision cell and fragmenting. The observation of the resulting charged fragments results in a characteristic spectrum. B. A combinatorial mass spectrum simulation. The root molecule, benzocaine, is iteratively fragmented by removing atoms or breaking bonds, resulting in a large fragmentation tree. Heuristic rules score nodes in the tree to predict intensities. C. ICEBERG spectrum simulation. ICEBERG learns to generate only the most relevant substructures. After generating fragments, a neural network module scores the resulting fragments to predict intensities. | 138 |
| 5.2 | Overview of ICEBERG. A. The target fragmentation directed acyclic graph (DAG) for an example molecule \mathcal{M} , benzocaine. Fragments are colored in black with missing substructures in gray. B. Example illustration for the generative process at a single step in the DAG generation predicting subfragments of $\mathcal{S}^{(2)}$. The root molecule \mathcal{M} , fragment of interest $\mathcal{S}^{(2)}$, and context vector C are encoded and used to predict fragment probabilities at each atom of the fragment of interest. A sample disconnection is shown at atom a_2 , resulting in fragment $\mathcal{S}^{(7)}$. C. ICEBERG Score module. Fragments generated from A are encoded alongside the root molecule. A Set Transformer module predicts intensities for each fragment, allowing mass changes corresponding to the loss or gain of hydrogen atoms, resulting in the final predicted mass spectrum. | 144 |
| 5.3 | ICEBERG predictions are highly accurate. A. Cosine similarities to true spectra on NPLIB1 (left) and NIST20 respectively (right) for CFM-ID [73], 3DMolMS [157], FixedVocab [152], NEIMS (FFN) [84], NEIMS (GNN) [150], MassFormer [86], SCARF [130], and ICEBERG. Error bars represent 95% confidence intervals using the standard error of the mean across three random seeds on a single test set split. B. Time required to predict spectra for 100 molecules randomly sampled from NIST20 on a single CPU, including the time to load models into memory. C,D. Comparison of NPLIB1 and NIST20 molecules in terms of synthetic accessibility (SA) score [182] and molecular weight (Mol. weight). | 146 |

| | | | |
|-----|---|---|-----|
| 5.4 | Examples of predicted spectra from ICEBERG. | Predictions are shown as generated by ICEBERG trained on NPLIB1 for select test set examples GNPS: CCMSLIB00003137969 (A), MoNA: 001659 (B), and GNPS: CCMSLIB00000080524 (C). The input molecular structures are shown (left); fragmentation spectra are plotted (right) with predictions (top, blue) and ground truth spectra (bottom, black). Molecular fragments are shown inset. Spectra are plotted with m/z shifted by the mass of the precursor adduct. All examples shown were not included in the model training set. | 147 |
| 5.5 | ICEBERG enables improved spectrum retrieval over other methods on both NPLIB1 (A) and NIST20 (B) compared to other spectrum prediction models. | Top k retrieval accuracy is computed by ranking a list of 49 additional candidates by putative cosine similarity as predicted by the model and determining the fraction of times that the true molecule is within the first k entries. A 95% confidence interval across three random model training seeds is shaded. | 149 |
| 6.1 | Enzyme-substrate interaction modeling strategies. | (A) Current machine learning-directed evolution strategies, which involve design-build-test-model-learn cycles measuring protein variant activity on a single substrate of interest. (B) The “dense screen” setting where homologous enzyme variants from one protein family are profiled against multiple substrates. In this setting, we can aim to generalize to either new enzymes (“enzyme discovery”) or new substrates (“substrate discovery”). (C) Three different styles of models evaluated in this study, where single task models independently build predictive models for rows and columns from panel (B), whereas a CPI model takes both substrates and enzymes as input. (D) An example CPI model architecture where pretrained neural networks extract features from the substrate and enzyme to be fed into a top-level feed forward model for activity prediction. | 163 |

6.2 **Assessing enzyme discovery in family wide screens.** (A) CPI models are compared against the single task setting by holding out enzymes for a given substrate and allowing models to train on either the full expanded data (CPI) or only data specific to that substrate (single-task). (B) AUPRC is compared on five different datasets, arranged from left to right in order of increasing number of enzymes in the dataset. Baseline models are compared against multi-task models, CPI models, and single-task models. K-nearest neighbor (KNN) baselines are calculated using Levenshtein edit distances to compare sequences; multi-task models use a shared feed forward network (FFN) to compute predictions against all substrate targets, CPI models utilize FFN with either concatenation (“{prot repr.}, {sub repr.}”) or dot product interactions (“{prot repr.}•{sub repr.}”), and ridge regression is used for single-task models. ESM-1b features indicate protein features extracted from a masked language model trained on UniRef50 [206]. Halogenase and glycosyltransferase datasets are evaluated using leave-one-out splits, whereas BKACE, phosphatase, and esterase datasets are evaluated with 5 repeats of 10 different cross validation splits. Standard error bars indicate the standard error of the mean of results computed with 3 random seeds. Each method is compared to the single-task L2-regularized logistic regression model (“Ridge: ESM-1b”) using a 2-sided Welch T test, with each additional asterisk representing significance at [0.05, 0.01, 0.001, 0.0001] thresholds respectively after application of a Benjamini-Hochberg correction. (C) Average AUPRC on each individual “substrate task” is compared between compound protein interaction models and single-task models. Points below 1 indicate substrates on which single-task models better predict enzyme activity than CPI models. CPI models used are FFN: [ESM-1b, Morgan] and single-task models are Ridge: ESM-1b. (D) AUPRC values from the ridge regression model are plotted against the average enzyme similarity in a dataset, with higher enzyme similarity revealing better predictive performance. (E) AUPRC values from the ridge regression model broken out by each task are plotted against the fraction of active enzymes in the dataset. Best fit lines are drawn through each dataset to serve as a visual guide.

| | | | |
|-----|--|---|-----|
| 6.3 | Assessing substrate discovery in family wide screens. | CPI models and single task models are compared on the glycosyltransferase, esterase, and phosphatase datasets, all with 5 trials of 10-fold cross validation. Error bars represent the standard error of the mean across 3 random seeds. Each model and featurization is compared to “Ridge: Morgan” using a 2-sided Welch T test, with each additional asterisk representing significance at [0.05, 0.01, 0.001, 0.0001] thresholds respectively, after applying a Benjamini-Hochberg correction. Pretrained substrate featurizations used in “Ridge: JT-VAE” are features extracted from a junction-tree variational auto-encoder (JT-VAE) [236]. Two compound protein interaction architectures are tested, both concatenation and dot-product, indicated with “[{prot repr.}, {sub repr.}]” and “{prot repr.}•{sub repr.}” respectively. In the interaction based architectures, ESM-1b indicates the use of a masked language model trained on UniRef50 as a protein representation [206]. Models are hyperparameter optimized on a held out halogenase dataset. AUCROC results can be found in Figure 6.9 in Section 6.6. | 169 |
| 6.4 | Evaluating single-task models on kinase repurposing and discovery tasks | Kinase data from Davis et al. is extracted, featurized, and split as prepared in Hie et al. Multilayer perceptrons (MLP) and Gaussian process + multilayer perceptron (GP+MLP) models are employed. We add variants of these models without CPI training separate single-task models for each enzyme and substrate in the training set, as well as linear models using both pretrained featurizations (“Ridge: JT-VAE”) and fingerprint based featurizations of small molecules (“Ridge: Morgan”). Spearman correlation is shown for (A) held out kinases not in the training set and (B) held out small molecules not in the training set across 5 random initializations. (C) We repeat the retrospective evaluation of lead prioritization. The top 5 average acquired K_d values are shown for the CPI models in Hie et al. compared against a linear, single-task ridge regression model using the same features. (D) The top 25 average acquired K_d values are shown. | 171 |

| | | |
|-----|---|-----|
| 6.5 | Structure-based pooling improves enzyme activity predictions. (A) | |
| | Different pooling strategies can be used to combine amino acid representations from a pretrained protein language model. Yellow coloring in the schematic indicates residues that will be averaged to derive a representation of the protein of interest. (i) We introduce active site pooling, where only embeddings corresponding to residues within a set radius of the protein active site are averaged. By increasing the angstrom radius from the active site, we increase the number of residues pooled. Crystal structures shown are taken from the BKACE reference structure, PDB: 2Y7F rendered with Chimera [242]. (ii, iii) We also introduce two other alignment based pooling strategies: coverage and conservation pooling average only the top- k alignment columns with the fewest gaps and highest number of conserved residues respectively. (iv) Current protein embeddings often take a mean pooling strategy to indiscriminately average over all sequence positions. (B) Enzyme discovery AUPRC values are computed for various different pooling strategies. Each strategy is tested for different thresholds of residues to pool, comparing against both KNN Levenshtein distance baselines and a mean pooling baseline. The same hyperparameters are used as set in Figure 6.2 for ridge regression models. The kinase repurposing regression task from Hie et al. is shown with Spearman’s ρ instead of AUPRC as interactions are continuous, not binarized. All experiments and are repeated for 3 random seeds. | 173 |
| 6.6 | Dataset substrates 6 exemplar molecule substrates are randomly chosen from each dataset and displayed. | 179 |
| 6.7 | Dataset diversity Distributions of top-5 enzyme similarity (left) and substrate similarity (right) are shown across enzyme datasets collected. Enzyme similarity is calculated as the percent overlap between two sequences in their respective multiple sequence alignment, excluding positions where both sequences contain gaps. Substrate similarity is computed using Tanimoto similarity between 2048-bit chiral Morgan fingerprints. | 180 |
| 6.8 | Enzyme discovery benchmarking with AUCROC On the 5 different datasets tested, K-nearest neighbor baselines with Levenshtein edit distance are compared against feed-forward networks using various featurizations and ridge regression models in terms of AUC ROC performance. ESM-1b features indicate protein features extracted from a masked language model trained on UniRef50 [206]. Concatenation and dot product architectures are indicated with “[{prot repr.}, {sub repr.}]” and “[prot repr.]•{sub repr.}” respectively. Halogenase and glycosyltransferase datasets are evaluated using leave-one-out splits. BKACE, phosphatase, and esterase datasets are evaluated with 5 repeats of 10 different cross validation splits. AUC ROC is calculated using scikit-learn for each substrate task separately before being averaged. Error bars represent the standard error of the mean across 3 random seeds. Each model and featurization is compared to “Ridge: ESM-1b” using a 2-sided Welch T test, with each additional asterisk representing significance at [0.05, 0.01, 0.001, 0.0001] thresholds respectively after applying a Benjamini-Hochberg correction. | 186 |

| | | |
|------|--|-----|
| 6.9 | Full substrate discovery AUC ROC results. CPI models and single task models are compared on the glycosyltransferase, esterase, and phosphatase datasets, all with 5 trials of 10-fold cross validation. Error bars represent the standard error of the mean across 3 random seeds. Each model and featurization is compared to “Ridge: Morgan” using a 2-sided Welch T test, with each additional asterisk representing significance at [0.05, 0.01, 0.001, 0.0001] thresholds respectively after applying a Benjamini-Hochberg correction. Pretrained substrate featurizations used in “Ridge: JT-VAE” are features extracted from a junction-tree variational auto-encoder (JT-VAE) [236]. Concatenation and dot-product architectures are indicated with “[{prot repr.}, {sub repr.}]” and “[{prot repr.}•{sub repr.}]” respectively. In the interaction based architectures, “ESM-1b” indicates the use of a masked language model trained on UniRef50 as a protein representation [206]. Models are hyperparameter optimized on a held out halogenase dataset. | 187 |
| 6.10 | Substrate Discovery Extended Analysis (A) Average AUPRC on each individual “enzyme task” is compared between compound protein interaction models and single-task models. Points below 1 indicate substrates on which single-task models better predict enzyme activity than CPI models. CPI models used are “FFN: [ESM-1b, Morgan]” and single-task models are “Ridge: Morgan”. (B) AUPRC values from the ridge regression model broken out by each task are plotted against the fraction of active enzymes in the dataset. Best fit lines are drawn through each dataset to serve as a visual guide. | 187 |
| 6.11 | MSA and structure based pooling across all datasets tested (A) Active site, coverage, conservation, and mean pooling are plotted for all 5 enzyme discovery datasets tested. Both AUCROC and AUPRC values are shown. These are compared against the Levenshtein distance baseline (dotted). (B) Equivalent analysis is conducted on the filtered kinase dataset extracted from Davis et al. with MAE, RMSE, and Spearman rank correlation shown [223]. The same hyperparameters are used as set in Figure 6.2 for ridge regression models. All experiments are repeated for 3 random seeds following the same split evaluation as in other enzyme discovery model benchmarking. | 188 |
| 6.12 | Enzyme discovery halogenase prediction results Ground truth binary enzyme-substrate activities (left) are compared against a single seed of predictions made through cross validation using a single-task ridge regression model (middle) and a CPI based model, FFN: [ESM-1b, Morgan] (right). | 189 |
| 6.13 | Enzyme discovery glycosyltransferase prediction results Ground truth binary enzyme-substrate activities (left) are compared against a single seed of predictions made through cross validation using a single-task ridge regression model (middle) and a CPI based model, FFN: [ESM-1b, Morgan] (right). | 190 |
| 6.14 | Enzyme discovery BKACE prediction results Ground truth binary enzyme-substrate activities (left) are compared against a single seed of predictions made through cross validation using a single-task ridge regression model (middle) and a CPI based model, FFN: [ESM-1b, Morgan] (right). | 190 |

| | | |
|------|---|-----|
| 6.15 | Enzyme discovery esterase prediction results Ground truth binary enzyme-substrate activities (left) are compared against a single seed of predictions made through cross validation using a single-task ridge regression model (middle) and a CPI based model, FFN: [ESM-1b, Morgan] (right). | 191 |
| 6.16 | Enzyme discovery phosphatase prediction results Ground truth binary enzyme-substrate activities (left) are compared against a single seed of predictions made through cross validation using a single-task ridge regression model (middle) and a CPI based model, FFN: [ESM-1b, Morgan] (right). | 192 |
| 6.17 | Substrate discovery glycosyltransferase prediction results Ground truth binary enzyme-substrate activities (left) are compared against a single seed of predictions made through cross validation using a single-task ridge regression model (middle) and a CPI based model, FFN: [ESM-1b, Morgan] (right). . . | 193 |
| 6.18 | Substrate discovery esterase prediction results Ground truth binary enzyme-substrate activities (left) are compared against a single seed of predictions made through cross validation using a single-task ridge regression model (middle) and a CPI based model, FFN: [ESM-1b, Morgan] (right). . . | 193 |
| 6.19 | Substrate discovery phosphatase prediction results Ground truth binary enzyme-substrate activities (left) are compared against a single seed of predictions made through cross validation using a single-task ridge regression model (middle) and a CPI based model, FFN: [ESM-1b, Morgan] (right). . . | 194 |
| 7.1 | Evidential uncertainty for molecular prediction and discovery. (A) Evidential direct message passing or atomistic neural networks learn molecular representations, predict target properties, and infer the parameters of an underlying evidential distribution that captures the evidence in support of each prediction and enables uncertainty estimation. (B) Uncertainties are applied during learning (I) to guide sample acquisition and during deployment (II) to discover high confidence candidates with high empirical success rates. | 197 |
| 7.2 | Building and training evidential models. (A) Evidential layers can be added to the end of existing molecular feature extractor neural networks. The output of the evidential layer is the parameters (\mathbf{m}) defining the molecule's evidential distribution. (B) For continuous regression learning problems, the evidential distribution $p(\boldsymbol{\theta} \mathbf{m})$ can take the form of a Normal-Inverse-Gamma distribution, where $\mathbf{m} = \{\gamma, \nu, \alpha, \beta\}$, to yield both property prediction and uncertainty estimates. Color represents likelihood density (darker = greater density). (C) The model (feature extractor and evidential layer) is trained end-to-end using backpropagation with a dual-objective loss that jointly maximizes model fit and inflates uncertainty (i.e., minimizes evidence) on errors. | 199 |

| | | |
|-----|---|-----|
| 7.3 | Benchmarking evidential uncertainty for molecular property prediction. (A) Lower-N regression tasks using 2D molecular representations for uncertainty benchmarking. (B, C) Prediction error, measured as root mean squared error (RMSE) or mean average error (MAE), at different confidence percentile cutoffs for the Delaney (B) and QM7 (C) datasets. Mean \pm 95% confidence interval (c.i.), $n = 10$ independent trials. (D) Higher-N regression tasks using 2D or 3D molecular representations. (E, F) Prediction error at different confidence percentile cutoffs for the Docking (E) and QM9 (F) datasets for 2D direct message passing (D-MPNN; E) and 3D atomistic (SchNet; F) neural networks, respectively. Mean \pm 95% c.i., $n = 5$ independent trials. | 201 |
| 7.4 | Tunability of the evidential uncertainty. (A) The evidential regression method can be fine tuned with a single hyperparameter, λ , in order to achieve more calibrated predictions for a given dataset. (B) Estimated confidence (cumulative probability) against the observed proportion correct for an evidential D-MPNN evaluated on the Delaney dataset. Dotted line represents perfect calibration. Mean \pm 95% c.i., $n = 5$ independent trials. (C) Area between the observed calibration curve and the perfect calibration line across several lower-N datasets for evidential D-MPNNs trained with varying λ . Dotted lines represent calibration of an ensemble of models. Mean \pm 95% c.i., $n = 10$ independent trials. | 203 |
| 7.5 | Evidential active learning and Bayesian optimization. (A) Experimental scheme. (B) Active learning with explorative (solid) versus random (dashed) sampling for D-MPNN evaluated on the QM9 dataset. Mean \pm 95% c.i., $n = 10$ independent trials. (C) Change in sample efficiency for explorative acquisition in (B), evaluated as the percent decrease in predictive error relative to a randomly-selected training set. (D) Bayesian optimization performance on Enamine 50k data, measured by the percentage of top-500 scores found as a function of the number of ligands explored. Solid traces represent an upper confidence bound (UCB) acquisition strategy. Mean \pm 95% c.i., $n = 10$ independent trials. (E) Average 10-nearest training set neighbors (10-NN) Tanimoto distance for batch samples after the first round of acquisition in Bayesian optimization experiments. Dots represent median; bars represent interquartile range; lines represent upper and lower adjacent values. $n = 10$ independent trials, two-tailed unpaired t-test, **** $P < 0.0001$ | 204 |
| 7.6 | Uncertainty guided nomination in virtual screens. (A) Experimental framework in which trained uncertainty-aware models are deployed on a discovery dataset, and molecules are ranked based on predicted properties. Uncertainty filtering is used to prioritize candidates among the top ranking molecules. (B) Performance of evidential D-MPNN after training to predict <i>E. coli</i> growth inhibition. (C) t-SNE visualization of training set (orange) and discovery dataset (Broad library), colored by predicted evidential uncertainties (blue). (D) Application of confidence filters to prioritize sets of antibiotic candidates with high experimental hit rates. Mean \pm 95% c.i., $n = 10$ independent trials. | 208 |

| | | |
|------|--|-----|
| 7.7 | Uncertainty benchmarking and calibration for lower-N datasets. (A, B) Prediction error, measured as MAE (A) or RMSE (B), at different confidence percentile cutoffs for D-MPNNs evaluated on each of the 2D lower-N datasets. (C) Estimated confidence (cumulative probability) against the observed proportion correct for an evidential D-MPNN evaluated on each of the 2D lower-N datasets, with regularization parameter $\lambda = 0.2$. Mean \pm 95% c.i., $n = 10$ independent trials. | 221 |
| 7.8 | Uncertainty benchmarking and calibration for additional lower-N Therapeutics Data Commons datasets. (A, B) Prediction error, measured as MAE (A) or RMSE (B), at different confidence percentile cutoffs for D-MPNNs evaluated on each of the 2D lower-N datasets. (C) Estimated confidence (cumulative probability) against the observed proportion correct for an evidential D-MPNN evaluated on each of the 2D lower-N datasets, with regularization parameter $\lambda = 0.2$. Mean \pm 95% c.i., $n = 10$ independent trials. | 222 |
| 7.9 | Spearman rank correlation between error and uncertainty. Spearman rank correlation coefficient between the estimated uncertainty and the absolute error for each point across lower-N (A) , additional lower-N Therapeutics Data Commons (TDC) (B) , higher-N (C) , and atomistic (D) datasets. Mean \pm 95% c.i., $n = 10$ independent trials for lower-M and $n = 5$ independent trials for atomistic and higher-N datasets. For QM9 in the higher-N D-MPNN setting (C), standard error bars are computed across all tasks and independent trials ($n = 60$). | 223 |
| 7.10 | Task-specific cutoffs for QM9 dataset. RMSE at different confidence percentile cutoffs for D-MPNNs evaluated on each of the individual tasks from the QM9 dataset. Mean \pm 95% c.i., $n = 5$ independent trials. | 224 |
| 7.11 | Uncertainty benchmarking and calibration for higher-N 2D and 3D datasets. (A, B) Prediction error, measured as MAE (A) or RMSE (B), at different confidence percentile cutoffs for models evaluated on the higher-N 2D and 3D datasets tested. (C) Estimated confidence (cumulative probability) against the observed proportion correct for the higher-N 2D and 3D datasets tested. Mean \pm 95% c.i., $n = 5$ independent trials. | 225 |
| 7.12 | Effect of λ on uncertainty calibration. Evidential D-MPNNs are trained with different regularization coefficients λ on each of the lower-N datasets. Estimated confidence (cumulative probability) against the observed proportion correct is computed and plotted across different λ . Mean \pm 95% c.i., $n = 10$ independent trials. | 225 |
| 7.13 | Cutoff RMSE is robust to small λ shifts. Evidential D-MPNNs are trained with different regularization coefficients λ on each of the lower-N datasets. RMSE is computed at the 10% confidence cutoff for each model. All computed errors are reported as the fold change of the top 10% RMSE for the evidential method relative to that of the ensemble baseline. Mean \pm 95% c.i., $n = 10$ independent trials. | 226 |

| | | |
|------|---|-----|
| 7.14 | Task-specific calibration for QM9 dataset. Estimated confidence (cumulative probability) against the observed proportion correct is computed for an evidential D-MPNN evaluated on QM9 and then broken down into task-specific plots. Mean \pm 95% c.i., $n = 5$ independent trials. | 227 |
| 7.15 | Antibiotic discovery datasets and uncertainty predictions. (A) Distribution of OD ₆₀₀ values for the training dataset of small molecules and their <i>in vitro</i> growth inhibitory activity against <i>E. coli</i> , as originally measured by Stokes et al [43]. Lower OD ₆₀₀ values indicate less <i>E. coli</i> growth and hence correspond to greater antibiotic activity. (B) Distribution of predicted OD ₆₀₀ values and evidential uncertainties for molecules in the Broad Drug Repurposing Hub discovery dataset. (C) Distribution of empirically determined OD ₆₀₀ for the subset of the discovery set (162 out of 6,111 total molecules) that was experimentally tested for <i>in vitro</i> growth inhibitory activity against <i>E. coli</i> . . | 228 |

List of Tables

| | | |
|-----|---|----|
| 2.1 | Comparison of SIRIUS and MIST for annotating IBD cohort data. | 74 |
| 2.2 | Fingerprint prediction accuracy on CANOPUS GNPS subset dataset for MIST top model variations trained. All models are trained without any auxiliary MAGMa loss and without any simulated spectra added to the training set. “MIST - unfolding (linear)” indicates a model without unfolding layers featuring only a single linear layer to predict fingerprint outputs. “MIST - unfolding (FFN)” indicates a MIST model trained without unfolding layers, but instead using a dense feedforward network of 3 layers. Each experiment was conducted on a single split of the data. Results are shown \pm the standard error of the mean. | 77 |
| 2.3 | Full performance metrics for CSI:FingerID, FFN, MIST, and a MIST ensemble of 5 models are shown averaged across 3 structural splits of the data. Log likelihoods are clamped such that they have a minimum of -5 . Results are shown \pm standard error with an arbitrary number of significant figures. | 79 |
| 2.4 | Full retrieval accuracy performance metrics for CSI:FingerID, FFN, MIST, and a MIST ensemble of 5 models are shown averaged across 3 structural splits of the data. | 80 |
| 2.5 | Full fingerprint prediction accuracy on CANOPUS GNPS subset dataset for various model ablations. All model performances are computed on a merged dataset of fingerprint predictions computed for 3 independent random structural split hold outs of the CANOPUS dataset. FFN: Feed forward neural network. Transformer: A partial reimplementaion of [83]. Results are shown plus or minus the standard error of the mean. | 81 |
| 2.6 | Fingerprint prediction accuracy on primary CSI2022 dataset to compare models with and without pairwise featurizations. Model performance values are computed on a merged dataset of fingerprint predictions computed for 3 independent random structural split test sets. All MIST models described were trained without simulated data or MAGMa auxiliary loss features for quicker model runs. Log likelihoods are clamped such that they have a minimum of -5 . Results are shown \pm standard error with an arbitrary number of significant figures. | 82 |

| | | |
|------|--|-----|
| 2.7 | Fingerprint prediction accuracy on CANOPUS GNPS subset dataset for MIST models with various featurizations of peak intensities. All models are trained without any auxiliary MAGMa loss and without any simulated spectra added to the training set. “MIST + float intensity” indicates a concatenation of a single floating point value of intensity between 0 and 1; “MIST + float intensity & pooling” indicates a concatenation of a floating point intensity and selective intensity-weighted pooling of hidden states in the transformer; “MIST + log intensity” indicates a concatenation of a single log-transformed intensity value; “MIST + cat. intensity” indicates a concatenation of a one-hot vector of intensity transformed into one of 10 evenly spaced categorical bins; “MIST + zero intensity” indicates a concatenation of a constant float of zero rather than the true intensity. Each experiment was conducted on a single split of the data and average results are shown \pm the standard error of the mean. | 82 |
| 2.8 | Full retrieval accuracy on CANOPUS GNPS subset dataset for various model ablations. All model performance values are computed on a merged dataset of retrieval predictions computed for 3 independent random structural splits of 3 random structural splits of the CANOPUS dataset. FFN: Feed forward neural network. | 82 |
| 2.9 | Average molecular similarity between spectra pairs with the 0.1% highest similarity according to the specified distance function. | 83 |
| 2.10 | Forward model ablations demonstrate the utility of unfolding and reverse mass predictions. Three variants of the forward simulation model are trained on a single split of the public GNPS dataset: “FFN” (full model), “FFN - unfolding” (full model a single linear layer to map to the 15,000 dimension output), and “FFN - reverse” (full model without predicting mass differences as in [84]). The average cosine similarity and coverage (i.e., fraction of overlapping peaks between the top 100 predicted peaks and ground truth spectra) between the binned spectra and true spectra are shown. Values are plotted \pm the standard error of the mean. | 83 |
| 2.11 | Model hyperparameters searched and set for CSI fingerprint prediction. | 84 |
| 2.12 | Model hyperparameters searched and set for Morgan fingerprint prediction. | 85 |
| 3.1 | MIST-CF outperforms comparable neural network baselines at molecular formula annotation from MS/MS. Models were trained to predict held out NPLIB1 test examples using training sets consisting of “NPLIB1” or “NPLIB1 + NIST20.” The best value in each column is typeset in bold. Models were evaluated using three independent formula splits. Values are shown \pm standard errors of the mean. | 101 |

| | | |
|-----|--|-----|
| 3.2 | Model accuracy of MIST-CF and SIRIUS evaluated on CASMI2022. “Accuracy (\mathcal{F})” indicates accuracy of predicting the exact formula correctly. “Accuracy (\mathcal{A})” indicates the fraction of spectra for which the first predicted (formula, adduct) pair includes the correct adduct. “Accuracy ($\mathcal{F} + \mathcal{A}$)” indi- cates the fraction of spectra for which the total elemental composition of the top predicted formula and adduct matches the elemental composition of the summed true formula and adduct. “Predicted” indicates the number of spectra for which a formula and adduct prediction was made. “SIRIUS” makes predictions with the default command line tool described above; “SIRIUS (CSI:FingerID)” predicts formula, adduct pairs using CSI:FingerID rankings against PubChem; SIRIUS (Submission) is taken directly from the previous CASMI2022 results. MIST-CF (20 peaks) uses default 20 subformula peaks, whereas MIST-CF (50 peaks) utilizes 50 subformula peaks. Accuracy (\mathcal{F}) and Accuracy (\mathcal{A}) are not reported for SIRIUS because the method does not distinguish between formula, adduct pairs that sum to the same molecular formula, and our metrics default to the worst-case ranking for ties (i.e., some- what unfairly penalizing methods such as SIRIUS that have more ties). All accuracies are reported for the top 1 prediction and divided by the total num- ber of spectra (304), not the total number of predicted spectra. Numbers typeset in bold indicate the best result in the column. | 105 |
| 3.3 | Hyperparameter settings for MIST-CF, FFN, Transformer, and the FastFilter model to select batches of molecular formula. | 108 |
| 4.1 | Model coverage (higher better) of true peak formulae as determined by MAGMa at various max formula cutoffs for the NIST20 and NPLIB1 datasets. Best result for each column is in bold. Results are computed for a single test set; all re-trained models (i.e., Autoregressive and SCARF variants) are averaged across three random seeds. | 118 |
| 4.2 | Spectra prediction in terms of cosine similarity, coverage (propor- tion of ground-truth peaks that are covered by the top 100 non-zero predictions), validity (the fraction of predicted peaks for which a chemically plausible explanation is possible), and time. Best value in each column is typeset in bold (higher is better for all metrics but time). Results are averaged across 3 random seeds on a single data split for all re- trainable models (i.e., not CFM-ID). | 119 |
| 4.3 | NIST20 spectra prediction retrieval top k accuracy for different values of k. All values represent the mean across three separate random seeds \pm the standard error of the mean for a single test set. | 122 |
| 4.4 | NPLIB1 spectra prediction retrieval top k accuracy for different values of k. All values represent the mean across three separate random seeds \pm the standard error of the mean for a single test set. | 122 |

| | | |
|------|--|-----|
| 4.5 | Retrieval accuracy on NIST20 for a single 500 molecule subset of the test set using a library of 49 decoys or all decoys contained in PubChem (“None”). Results were repeated for 3 random training seeds of the model and are shown \pm the standard error of the mean. The top value is shown in bold. | 122 |
| 4.6 | Model coverage of true peak formulae as determined by MAGMa at various max formula cutoffs for the NPLIB1 dataset. Results are calculated for a single held out test split, shown \pm the standard error of the mean across three random seeds for all models that were retrained. The best value (i.e., highest) is typeset in bold for each column. | 124 |
| 4.7 | Model coverage of true peak formulae as determined by MAGMa at various max formula cutoffs for the NIST20 dataset. Results are calculated for a single held out test split, shown \pm the standard error of the mean across three random seeds for all models that were retrained. The best value (i.e., highest) is typeset in bold for each column. | 125 |
| 4.8 | Spectra prediction in terms of cosine similarity, coverage (proportion of ground-truth peaks that are covered by the top 100 non-zero predictions), validity (the fraction of predicted peaks for which a chemically plausible explanation is possible), and time. Best value in each column is typeset in bold (higher is better for all metrics but time). Values are shown \pm the standard error of the mean computed across three random seeds on a single test set for all models that could be retrained (i.e., not CFM-ID). | 127 |
| 4.9 | Spectra prediction accuracy comparing inclusion (Cosine sim.) and exclusion (Cosine sim. (no MS1)) of the precursor mass. For all compounds, the peak at the mass of the input compound is masked in the prediction and ground truth to compute Cosine sim. (no MS1). All results represent an average on a single test set across three random seeds. | 127 |
| 4.10 | NIST20 retrieval accuracy averaged across three random seeds of model training stratified by the weight of the target molecule. | 128 |
| 4.11 | Graph neural network (GNN) atom features. | 131 |
| 4.12 | Model and baseline hyperparameters. | 133 |
| 5.1 | Comparing the accuracy of spectrum prediction on NIST20 using random (easier) or scaffold (harder) split. | 150 |

| | | |
|------|--|-----|
| 5.2 | Spectra prediction accuracy on NIST20, and NIST20 (scaffold split). Results are shown for CFM-ID [73], 3DMolMS [157], FixedVocab [152], NEIMS (FFN) [84], NEIMS (GNN) [150], MassFormer [86], SCARF [130], and ICEBERG. Cosine similarity is calculated at a bin resolution of 0.1 m/z; Coverage indicates the fraction of peaks in the ground truth spectrum explained by the predicted spectrum on average; Valid indicates the fraction of predicted peaks that can be explained as a subset of the precursor molecule’s molecular formula (filtered to formulae with ring-double bond equivalents of greater or equal to zero); Time (s) indicates the number of seconds required to predict 100 random spectra from NIST20 on a single CPU including the time to load the model. Values are shown \pm 95% confidence intervals of the mean across three random seeds for the random split. Scaffold split experiments are only repeated once and shown without confidence intervals. The best value in each column is typeset in bold. | 152 |
| 5.3 | Spectra prediction accuracy on NPLIB1. Results are shown for CFM-ID [73], 3DMolMS [157], FixedVocab [152], NEIMS (FFN) [84], NEIMS (GNN) [150], MassFormer [86], SCARF [130], and ICEBERG. Cosine similarity is calculated at a bin resolution of 0.1 m/z; Coverage indicates the fraction of peaks in the ground truth spectrum explained by the predicted spectrum on average; Valid indicates the fraction of predicted peaks that can be explained as a subset of the precursor molecule’s molecular formula (filtered to formulae with ring-double bond equivalents of greater or equal to zero). Values are shown \pm 95% confidence intervals as measured by three random seeds on a single test split. The best value in each column is typeset in bold. | 153 |
| 5.4 | NIST20 spectra retrieval top k accuracy for different values of k, \pm a 95% confidence interval across three random seeds on the same test set. | 154 |
| 5.5 | NPLIB1 spectra retrieval top k accuracy for different values of k, \pm a 95% confidence interval across three random seeds on the same test set. | 154 |
| 5.6 | Dataset details. | 155 |
| 5.7 | Graph neural network (GNN) atom features. | 155 |
| 5.8 | Hyperparameter descriptions for ICEBERG. | 156 |
| 5.9 | ICEBERG Generate and Score hyperparameter grid and selected values. | 157 |
| 5.10 | Baseline hyperparameters. | 157 |
| 6.1 | Summary of curated datasets with the number of unique enzymes, unique substrates, and unique pairs in each dataset in addition to an exemplar structure for the protein family. | 166 |
| 6.2 | Active site structure references used in pooling. All structure informed pooling strategies require a catalytic center in order to define various angstrom shells of residues to pool over. This table provides the PDB reference crystal structure as well as the reference residues or structural elements used to define the pooling center, from which spherical radii originate. | 178 |

| | | |
|-----|--|-----|
| 6.3 | Summary of valid substrate and sequence “tasks”. In each dataset, only certain substrates and sequences are defined as valid “tasks” based upon the balance between active and inactive examples. Each substrate or sequence used for an enzyme or substrate discovery task respectively requires at least 2 positive examples and at a minimum, 10% of examples in that task must be part of the minority class. This table defines the number of valid substrate and sequence tasks. | 179 |
| 6.4 | Summary and classifications of different models utilized. | 182 |
| 6.6 | Full enzyme discovery area under the precision recall curve (AUPRC) results. On the 6 different datasets tested (thiolase datasets used for hyperparameter optimization), K-nearest neighbor baselines with Levenshtein edit distance are compared against feed-forward networks using various featurizations and ridge regression models. Pretrained features (“ESM-1b”) indicate protein features extracted from a masked language model trained on UniRef50 [206]. Two compound protein interaction architectures are tested, both concatenation and dot products, indicated with “[{prot repr.}, {sub repr.}]” and “[prot repr.]•{sub repr.}” respectively. Halogenase and glycosyltransferase datasets are evaluated using leave-one-out splits, whereas BKACE, phosphatase, and esterase datasets are evaluated with 5 repeats of 10 different cross validation splits. Average precision is calculated using scikit-learn for each substrate task separately before being averaged. Average values are presented across 3 random seeds \pm standard error. ¹ Used for hyperparameter optimization | 183 |
| 6.5 | Summary and classifications of different models utilized in our re-analysis of Hie et al. [226] | 183 |
| 6.7 | Full enzyme discovery area under the receiver operating curve (AUC-ROC) results. On the 6 different datasets tested (thiolase datasets used for hyperparameter optimization), K-nearest neighbor baselines with Levenshtein edit distance are compared against feed-forward networks using various featurizations and ridge regression models. ESM-1b features indicate protein features extracted from a masked language model trained on UniRef50 [206]. Two compound protein interaction architectures are tested, both concatenation and dot products, indicated with “[{prot repr.}, {sub repr.}]” and “[prot repr.]•{sub repr.}” respectively. Halogenase and glycosyltransferase datasets are evaluated using leave-one-out splits, whereas BKACE, phosphatase, and esterase datasets are evaluated with 5 repeats of 10 different cross validation splits. AUC ROC is calculated using scikit-learn for each substrate task separately before being averaged. Average values are presented across 3 random seeds \pm standard error. ¹ Used for hyperparameter optimization | 184 |

| | | |
|-----|---|-----|
| 6.8 | <p>Full substrate discovery area under the precision recall curve (AUPRC) results. CPI models and single task models are compared on the glycosyltransferase, esterase, and phosphatase datasets, all with 5 trials of 10-fold cross validation. Each model and featurization is compared to “Ridge: Morgan” using a 2-sided Welch T test, with each additional asterisk representing significance at [0.05, 0.01, 0.001, 0.0001] thresholds respectively after applying a Benjamini-Hochberg correction. Pretrained substrate featurizations used in “Ridge: JT-VAE” are features extracted from a junction-tree variational auto-encoder (JT-VAE) [236]. Two compound protein interaction architectures are tested, both concatenation and dot-product, indicated with “[{prot repr.}, {sub repr.}]” and “[{prot repr.}•{sub repr.}]” respectively. In the interaction based architectures, ESM-1b indicates the use of a masked language model trained on UniRef50 as a protein representation [206]. Average precision is calculated using scikit-learn for each substrate task separately before being averaged. Models are hyperparameter optimized on a held out halogenase dataset. Values represent mean values across 3 random seeds \pm standard error. ¹Used for hyperparameter optimization</p> | 184 |
| 6.9 | <p>Full substrate discovery area under the receiver operating curve (AUC-ROC) results. CPI models and single task models are compared on the glycosyltransferase, esterase, and phosphatase datasets, all with 5 trials of 10-fold cross validation. Each model and featurization is compared to “Ridge: Morgan” using a 2-sided Welch T test, with each additional asterisk representing significance at [0.05, 0.01, 0.001, 0.0001] thresholds respectively after applying a Benjamini-Hochberg correction. Pretrained substrate featurizations used in “Ridge: JT-VAE” are features extracted from a junction-tree variational auto-encoder (JT-VAE) [236]. Two compound protein interaction architectures are tested, both concatenation and dot-product, indicated with “[{prot repr.}, {sub repr.}]” and “[{prot repr.}•{sub repr.}]” respectively. In the interaction based architectures, “ESM-1b” indicates the use of a masked language model trained on UniRef50 as a protein representation [206]. Models are hyperparameter optimized on a held out halogenase dataset. Values represent mean values across 3 random seeds \pm standard error. ¹Used for hyperparameter optimization</p> | 185 |
| 7.1 | <p>Model error at various confidence percentile cutoffs. For a given confidence percentile cutoff, top performing methods based on prediction standard error of the mean (\pm s.e.m.) are bolded. A cutoff of 0.95 indicates that only the top 5% most confident predictions are considered. Full confidence plots for all datasets are shown in Figure 7.3 and Figures 7.7, 7.10, 7.11. Mean \pm s.e.m. (RMSE for all D-MPNN models, MAE for atomistic); $n = 10$ independent trials for lower-N datasets, $n = 5$ independent trials for higher-N datasets.</p> | 202 |

| | | |
|-----|--|-----|
| 7.2 | Extended model error at various confidence percentile cutoffs including TDC lower-N data. For a given confidence percentile cutoff, top performing methods based on prediction standard error of the mean (\pm s.e.m.) are bolded. A cutoff of 0.95 indicates that only the top 5% most confident predictions are considered. Full confidence plots for all datasets are shown in Figure 7.3 and Figures 7.7, 7.10, 7.11. Mean \pm s.e.m. (RMSE for all D-MPNN models, MAE for atomistic); $n = 10$ independent trials for lower-N datasets, $n = 5$ independent trials for higher-N datasets. | 220 |
| 7.3 | Statistical significance tests for Enamine 50k Bayesian optimization. Pairwise comparison of uncertainty quantification methods for upper confidence bound (UCB) acquisition in Bayesian optimization on Enamine 50k data (Figure 7.5D). Fold changes (FC; mean \pm s.d.) reflect the fold change between the mean percentage of top-500 scores found for the first method listed relative to the second method listed. Significance values reflect the result of two-tailed unpaired t-tests over $n = 10$ independent trials. | 220 |

Chapter 1

Introduction

Nothing was formerly further from my thoughts, than that I should ever trouble the world with anything in chemistry. There are so many books already upon this subject, and many of them wrote so well, that it is hardly possible for me... to offer anything that has not been said before.

Herman Boerhaave
1668 – 1738

In the first half of the 18th century, Herman Boerhaave used careful elemental analysis and techniques to isolate and identify the metabolite urea [1]. Fast forward nearly three full centuries since, high throughput measurements of biochemical systems have now transformed our understanding of biology at a range of scales. An exciting array of technologies have emerged in the past 20 years to measure the DNA [2], RNA [3], and proteins [4] in biological samples. Yet, despite these advances for measuring common macromolecules, the high throughput measurement of small molecules like urea—critical building blocks that mediate numerous processes—has received far less attention. The overarching goal of this thesis is to introduce a framework of new computational tools to identify the structure and potential functional properties of previously unknown metabolites.

1.1 Metabolites: “Canaries of the Genome”

The relative lack of attention toward small molecules can in part be explained by the high structural complexity and diversity of small molecule structures. In contrast, DNA and RNA molecules are long, polymeric strings of four repeating nucleic acid bases—adenine, thymine/uracil, guanine, and cytosine. This highly structured repetitive nature has enabled “sequencing” measurements [5] that read nucleic acids base-by-base. Proteins are similarly composed of a constrained repetition of a 20 amino acid alphabet, which, combined with the well-defined understanding of protein synthesis (i.e., DNA encodes proteins), simplifies high throughput proteomic measurements. Small molecules stand apart with their branched and complex structures.

Nevertheless, biologically produced small molecules, referred to as metabolites, perform a vast array of often underappreciated functions. As David Wishart, a leader in the

field of metabolomics, elegantly describes, metabolites can be seen as the “canaries of the genome” [6]. That is, small biological genetic changes in upstream enzymes can be magnified by nearly 10,000 fold [7], so metabolite quantities may be sensitive responders and signatures for understanding biochemical phenotypes.

Well-characterized metabolites serve as energy stores, basic building blocks for larger molecules, and transport mechanisms. New and unanticipated functions for even common metabolites are continuously being discovered. An endogenously found metabolite 2-hydroxyglutarate was found to be an oncometabolite, produced more abundantly due to mutations in the IDH1/IDH2 enzymes, often associated with cancers such as glioblastoma [8]. More recently, a high throughput screen of proteome stability revealed lactate to regulate protein degradation pathways by active site inhibition of SENP1 [9].

Beyond primary metabolites, secondary metabolites (i.e., those not directly necessary for survival or growth) produced across the kingdom of life perform their own considerable and important functions [10]–[12]. Within the microbiome in particular, recent effort has focused on identifying metabolite products that impact human health. New microbially derived bile acid conjugates have been discovered [13] and separate bile acid synthesis pathways have been correlated with longevity [14]. Another line of inquiry focused on understanding the impact of the microbiome on cancer risk has identified a previously unknown class of microbiome-derived metabolites, indolimines, as having genotoxic properties [15]. In the same way as metabolites may mediate cancer risk, metabolites have also been found to drive gut-brain axis interactions [16]. Collectively, metabolites have emerged as potential targets for therapeutic intervention. Not only do they find use as targets, but natural product metabolites produced by plants, bacteria, and fungi can be repurposed outside of their natural contexts as therapeutic medicines themselves [17].

While this thesis is motivated primarily by applications within human health, metabolites and small molecules mediate key interactions and phenomena within the environment more broadly [18]. Recent efforts have scrutinized the impact of pollutants. In a canonical example, Tian *et al.* identified 6PPD-quinone, a tire rubber antioxidant, as a toxin causing coho salmon mortality [19]. This was accomplished through an extensive analytical chemistry investigation.

1.2 Measuring the Metabolome

Metabolomics cumulatively represents the science of measuring and identifying these such small molecules (i.e., often defined as the chemical species under 1,500 Da) [20]. Tactically, because of the discrete structure of small molecules, we cannot use sequencing-based methods and instead utilize analytical chemistry techniques, most often mass spectrometry. While other techniques such as nuclear magnetic resonance (NMR) exist to structurally characterize molecules, these often lack the sensitivity to characterize low quantities of sample and perform *bulk* measurements of many different peaks [21]. As such, we consider mass spectrometry measurements herein, despite NMR remaining the gold standard for confirming an individually identified structure of interest.

As the metabolomics practitioner’s weapon of choice, a mass spectrometer is an exquisitely sensitive analytical chemistry instrument that can measure the mass of molecules down to

less than single digit parts-per-million error (i.e., mass over charge errors on the order of 10^{-4}) with up to femtogram-level sensitivity [22]. A common pipeline is to utilize a multi-step workflow of liquid chromatography followed by tandem mass spectrometry, denoted as LC-MS/MS. In the first chromatography step, the contents of an input sample of interest are separated out by time according to their affinities to reduce complexity. Then, at each time point, the eluting compounds are ionized (i.e., given a charge) and their masses are measured by the first mass analyzer, known as the MS1 measurement. Subsequently, the highest intensity, most abundant molecules are isolated in the mass spectrometer to go through a collision step, whereby the isolated molecule is fragmented and each fragment is measured. This produces a characteristic spectrum that provides clues as to the identity of the molecule. There are many different incarnations of this technology and excellent reviews of the various technologies [23] and their histories [24]; I further define the relevant components of this technology in each chapter when necessary.

Paradoxically, mass spectrometers are capable of detecting thousands of different molecules, yet our ability to assign the detected masses to small molecule structures remains limited. As reviewed by Alseekh *et al.*, current methods cannot fully define the number of metabolites present [25]. Even from our pool of putative metabolites we expect are present, we can only identify a narrow 8,000 of the 114,100 metabolites predicted to be in humans [26] and 14,000 of the >200,00 metabolites predicted in the plant kingdom [27].

While some of this gap is due to challenges with the physical instrumentation and separation techniques used to isolate metabolites, a large part of this gap surprisingly is due to our inability to identify which compounds occur from their analytical chemistry signature. Anecdotally, even for human serum measurements, which we would expect to be well characterized, only roughly 300 of 2,000 consistently detected metabolites can be matched to compounds. This same trend is true for public databases, with the largest Global Natural Products Social (GNPS) database [28] reportedly only featuring compound annotations for 13% of its identified metabolites [29]. A key question then, is how to build new computational tools that will increase this number of metabolites we can discover? And beyond this, how can we build technologies to not only identify, but also understand the function, of these various metabolites?

1.3 Toward an AI/ML Analysis Pipeline

Parallel to the improving instrumentation and growing interest in metabolite-driven biology, a separate revolution has been brewing over the last decade in the field of artificial intelligence and machine learning (AI/ML). Beginning with superhuman image classification [30], [31], deep learning methodologies have since replaced statistical learning approaches for everything from natural language translation [32], [33] to image generation [34] to protein folding [35]–[37]. At its core, these techniques eschew the traditional statistical learning paradigm of carefully crafted input features in favor of adaptive basis functions; with enough data, the power of auto-differentiation, and a well-framed loss function, large and over-parameterized models are capable of learning to approximate even the most complex functions.

Most recently, especially large “foundation” models have taken center stage [38]. In principle, with enough data, a large model can be trained once and then adapted to perform

many disparate, downstream tasks. The most compelling case of this phenomenon can be seen in natural language, where large language models trained with the objective of completing the sentence subsequently find use in a range of tasks from translation to text summarizing [32].

Unsurprisingly, there has been interest in targeting these powerful classes of deep learning models to advance and solve some of the most pressing problems in science [39]. Some more recent examples include small molecule synthesis planning [40], [41], small molecule property prediction and discovery [42], [43], viral gene therapy vector optimization [44], cancer detection [45], protein folding [35]–[37], and protein design [46]—all examples with profound importance for discovering new therapeutics and medicines.

While foundation models find uses in natural language processing, given the varied number of tasks considered, each of the demonstrations above has required separate innovation for neural network architectures, loss function, and training data. Protein folding is an especially apt example. Over the course of a decade, old statistical learning pipelines and models were slowly converted into deep generative models via a process Mohammed AlQuraishi dubbed the “neuralization of structure prediction pipelines” [47]. Previous statistical methods considered a set of related protein sequences known as a “multiple sequence alignment,” which could help constrain a protein, as well as a “template” structure of similar proteins. To reach groundbreaking performance, the AlphaFold2 system developed custom deep neural network components to gracefully accept, featurize, and convolve upon these varied inputs [35]. Simply inputting a sequence and training it to match a protein structure with existing deep neural network methods was insufficient to rival statistical methods.

If the first lesson from the protein folding analogy is that custom neural network solutions are needed to consider the biases imparted by statistical methods when working in scientific domains, the second is that we need well-defined tasks and objectives. For over two decades, the protein biology community organized public, carefully curated challenges and metrics [48], as well as data [49], to test the ability to predict protein structure. Importantly, these competitions provided the data, guidelines, problem framing, and inspiration that enabled an outside party, DeepMind, to solve a well-posed problem with AlphaFold2.

Returning to the case of mass spectrometry and metabolite identification, the field currently exists in a somewhat pre-neural network regime. At the time I commenced work on this thesis, and even at the time of its completion, the two most widely used computational prediction methods for resolving molecules from spectra, SIRIUS [50] and CFM-ID [51], are either impossible or highly difficult to retrain and do not use deep learning. Prior to this work, we have found vanishingly few public data benchmarks that jointly include training data, evaluation data, and the performance of common methods.

In part, this lack of benchmarks stems more so from lack of clarity on the tasks themselves; metabolomics processing pipelines have many moving parts and different steps with various utility. To add rigidity to this framework, I introduce and articulate four key tasks that I posit can be solved with deep learning (Figure 1.1):

- A. **MS1 formula prediction.** Initial data processing steps provide a mass fragmentation spectrum and its parent MS1 mass. The MS1 mass substantially narrows down the chemical space and set of structures to consider, but we can go one step further by predicting the exact molecular formula from a set of equivalent-mass options. Knowing

the set of elements in the target molecule further constrains the structure elucidation task. Several approaches such as SIRIUS [50] and BUDDY [52] have been proposed already for this task but do not have standardized training data or evaluations.

- B. Inverse spectrum-to-molecule prediction.** Given the input molecular formula and spectrum, we can then use models to predict key structural and functional properties of the molecule, in the form of a molecular fingerprint. A core challenge of this task is *how to represent mass spectra* into a neural network. Once a molecular fingerprint is predicted, these can be compared to a database of molecular candidates to rank the most likely identity for the spectrum of interest. CSI:FingerID [53] within SIRIUS [54] was the leading approach to this task, with others more recently emerging [55], [56]. Such inverse models can also be used to predict the chemical class of the unknown molecule, rather than features of the molecule itself [57], which we consider to be an extension of this same task.
- C. Forward molecule spectrum prediction.** Rather than try to invert the physical process, we can attempt to predict a spectrum from a molecule directly and follow the so-called “forward” data generating process by which molecules are fragmented in a mass spectrometer. This creates questions not about representing spectra as inputs to a model, but rather how to represent spectra as outputs of a model. A forward molecule-to-spectrum model can be used to create large synthetic libraries of putative standards that can then be queried as an orthogonal strategy to identify a molecule from its spectrum of interest. Many methods exist to address this challenge such as CFM-ID [51]; we review these in more depth within Chapters 4 and 5.
- D. Molecule candidate generation.** A core limitation of tasks B and C as we have described them is that they enable molecular annotation via a “retrieval” process. That is, the user provides a reference database of potential molecule structures, like PubChem [58] or the Human Metabolome Database [59], that could explain the candidate spectrum, and the molecules are either ranked by how closely their fingerprints match the predicted molecular fingerprint of the spectrum (task B) or by how closely the predicted spectrum of each molecule matches the true spectrum (task C). Given the vastness of molecular space, it is highly likely that many structures we are measuring are not yet recorded in databases. We assert that we should decouple the tasks of searching molecule space and ranking candidate molecules by making “molecule candidate generation” into a separate task. While I do not explicitly address molecular generation within this thesis, I see this as an important direction for future work. Some work has arisen already to this end [60], [61], but evaluations are far from clear and also implicitly conflate performance between spectrum-to-molecule prediction (task B) and candidate generation tasks. A path forward will require decoupling the evaluations accordingly and using standardized training and testing splits to evaluate performance, as has been done for molecular generation in other drug discovery settings [62].

Data resources. Training datasets are a key criteria to addressing any one of these tasks with machine learning. Typically, *paired data* containing molecules and their corresponding spectra, is ideal for learning the various tasks above. This thesis has converged on using

two separate datasets: the National Institute of Standards and Technology (NIST) spectral library [63] and the GNPS dataset [28]. Both datasets provide a set of paired molecules and their corresponding spectra. While each dataset contains duplicate entries for the same molecular structures, often measuring spectra with slightly different instrument parameters, I posit that the key axis for dataset quality is not the sheer number of spectra, but rather the number of unique molecules. The NIST dataset has on the order of approximately 30,000 unique molecules, whereas the GNPS dataset has grown progressively to up to approximately 15,000 molecules. We use a subset of the GNPS dataset from 2021 throughout this thesis as prepared by the authors Dührkop, Nothias, Fleischauer, *et al.* [57] that contains on the order of approximately 10,000 unique molecules. The GNPS spectral standard subset that we utilize is referred to as either CANOPUS or NPLIB1 throughout this thesis.

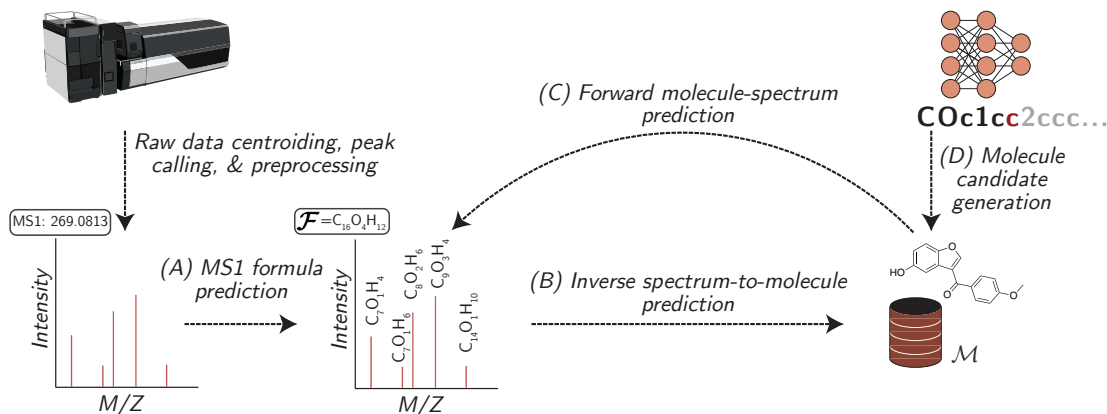


Figure 1.1: **Spectrum prediction workflow tasks.** Raw data is first processed into a mass spectrum with a precursor MS1 mass. Subsequently, a molecular formula can be predicted. Beyond this, three separate classes of models can be used to help identify the molecule structure, and the predictive direction of each model is shown.

1.4 Thesis Organization

This body of this thesis is separated according to the structure of six different project chapters, the first four of which directly address tasks A, B, and C introduced above. The final two content chapters address more general molecular modeling tasks that have implications for how we can consider predicting the functional relevance of newly identified metabolites. Throughout this work, we prioritize open-source code development and providing concrete data and splits for others to compare to and improve upon our work.

- Chapter 2 begins by introducing a new strategy for representing spectra as inputs to neural networks and solving task B for predicting molecule properties from an input spectrum (Figure 1.1). I introduce the Molecular Formula Transformer neural network architecture and our model MIST, and I demonstrate how this model compares favorably to a hand-crafted kernel-based statistical method for the same task. I show how

this model can be applied to clinical healthcare data from patients with inflammatory bowel disease to identify signatures of disease and introduce benchmarks for others attempting to address spectrum representation learning tasks in the future.

- Chapter 3 extends the MIST model to a new model, MIST-CF to address task A of metabolite annotation (Figure 1.1), predicting the molecular formula from an input mass spectrum. This work shows how our same neural network strategy can be adapted into an energy-based modeling framework to cover a wider range of the spectrum annotation and data processing pipeline.
- Chapter 4 switches to addressing the forward task of metabolite annotation, task C, predicting spectra from molecular structures (Figure 1.1). We introduce a model, SCARF, and the idea of decoding mass spectra as prefix trees of molecular formulae.
- Chapter 5 further continues the thread of predicting mass spectra from molecules as introduced by SCARF. We design an improved model, ICEBERG, that instead of outputting spectra as sets of formulae, predicts output spectra as sets of sub-fragments of the input molecule. This more physically-grounded approach proves to be increasingly accurate and beneficial in the task of spectrum annotation, albeit requiring more time at inference than SCARF.
- Chapter 6 is the first chapter to consider problems outside of mass spectrum data processing and instead investigates how to predict the substrate specificity of certain enzyme families using protein and small molecule representation learning. A better understanding of enzyme substrate specificity would be helpful both in industrial biotechnology applications of biocatalysis as well as understanding the pathways by which newly identified metabolites may be synthesized.
- Chapter 7 considers another orthogonal challenge of quantifying uncertainty of neural network predictions in the molecular sciences. This chapter utilizes a recently proposed “evidential uncertainty” method to quantify the uncertainty of neural network model predictions of small molecule properties. This work does not focus on metabolites, but given the generality of these methods, could be extended in future work.
- Chapter 8 concludes this thesis and provides several directions and ideas for future work.

Chapter 2

Predicting Molecular Properties and Substructures from Spectra

This work has previously appeared as S. Goldman, J. Wohlwend, M. Stražar, *et al.*, “Annotating metabolite mass spectra with domain-inspired chemical formula transformers,” *Nature Machine Intelligence*, vol. 5, no. 9, pp. 965–979, 2023. I jointly conceptualized the project and led the execution, implementation, and writing. Co-authors provided support, input, and guidance through all stages.

2.1 Introduction

Untargeted metabolomics is an important tool in advancing our understanding of cellular and environmental biochemistry [8], [14], [18], [65]–[67]. Liquid-chromatography tandem mass spectrometry (LC-MS/MS) is a widely used experimental approach to identify new metabolites, measuring both the mass of an intact molecule (MS1) and its fragments (MS2), measured as a spectrum of peaks [68]. The analysis of these experiments is critically bottlenecked by an inability to accurately annotate observed fragmentation spectra with the chemical structures of the molecules that produced them, with as many as 87% remaining unannotated [29]. In the many cases where the spectrum does not closely resemble a known standard spectrum, imperfect computational tools must be used to infer properties of the unknown molecule. Improving this inference step has the potential to drastically increase the information gleaned across all routine untargeted metabolomics experiments [68].

Despite the recent breakthroughs in deep learning that have “neuralized” the adjacent protein structure prediction field [35], [47], leading computational metabolomics annotation tools still rely upon hand-crafted heuristics and kernel functions (i.e., quantitative functions for comparing properties of mass spectra). Such methods can broadly be grouped into three categories: networking, forward prediction, and inverse prediction. Networking methods such as feature-based molecular networking [69], [70] cluster similar spectra to find neighborhoods of similar compounds (e.g., bile acids [13]). Forward methods (e.g., MetFrag [71], [72] and CFM-ID [51], [73]) augment the library of known metabolite standards by *in silico* fragmentation of molecules to generate artificial data. On the other hand, inverse fragmentation methods (e.g., CSI:FingerID [53] and FingerID [74]) predict molecular properties or

the molecule structure from the spectra directly.

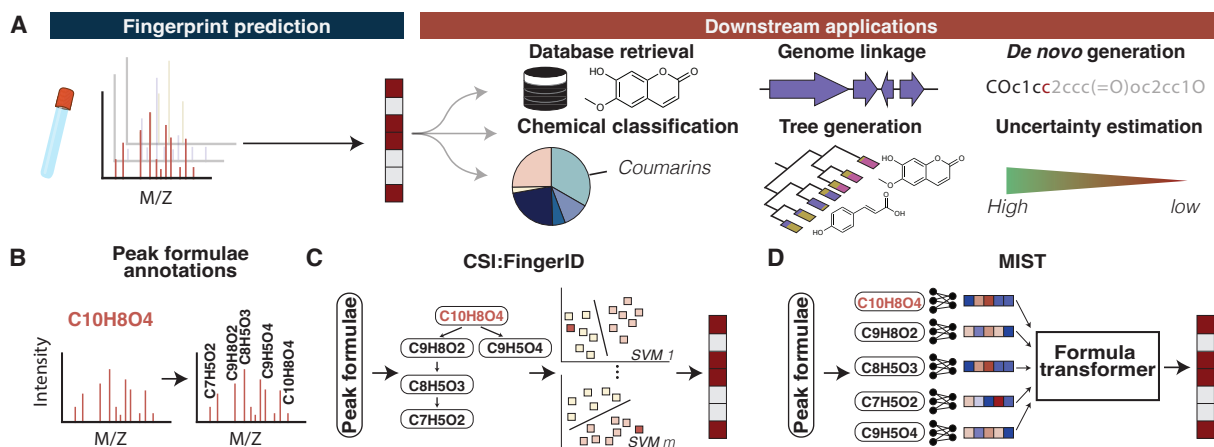


Figure 2.1: **MIST replaces the critical spectrum-to-fingerprint prediction step in the computational workflow.** **A.** A set of fragmentation peaks are extracted from a sample with tandem mass spectrometry. After predicting molecular property fingerprints, several different downstream software tools can be used to annotate molecules from databases, link metabolites to genomes, *de novo* generate plausible candidates, classify the unknown metabolite into compound classifications, generate statistics for all observed spectra, and calibrate certainty of annotations. **B.** The first step in the fingerprint prediction pipeline common to MIST and CSI:FingerID is to annotate the spectrum with a molecular formula for the full metabolite, then label each sub-peak with a molecular formula. **C.** CSI:FingerID arranges annotated sub-peak formulae into a “fragmentation tree”, and subsequently uses many independent support vector machine models to predict fingerprint property bits. **D.** MIST does not require a fragmentation tree, but rather directly utilizes molecular formulae as inputs to a Formula Transformer in order to predict molecular fingerprints.

Inverse prediction of a molecule’s structural fingerprint from its spectrum is an entrenched paradigm in metabolomics, most recently forming the backbone of the winning solution at Critical Assessment of Small Molecule Identification 2022 (CASMI2022) [75], [76] (Figure 2.1A). In particular, CSI:FingerID, the state of the art method, assigns each peak in the fragmentation spectrum a molecular formula and arranges these into a fragmentation tree [77] (Figure 2.1B-C). A set of kernels are then used to train a support vector machine to predict individual molecular fingerprint property bits describing the presence of certain molecular substructures. Predicted fingerprints can then be queried against a database of molecular structures [53] or used for a range of other downstream tasks and applications [57], [61], [78]–[80]. The performance for all such tools is fundamentally limited by the first step of fingerprint prediction, which we aim to improve with this work (Figure 2.1A).

While several inverse approaches leverage deep learning, they have failed to outperform kernel CSI:FingerID when trained with equivalent data. Representation learning of spectra such as Spec2Vec [81], MS2DeepScore [82], and sinusoidal embeddings [83] can be used to learn a more meaningful distance between spectra to facilitate molecular networking. Forward models have incorporated feed forward networks and graph neural networks to predict a

fragmentation spectrum directly from molecular structure [84]–[86]. In the inverse direction, MassGenie [87], Spec2Mol [88], and MetFID [56] directly attempt to generate fingerprints or SMILES strings [89] from mass spectra. Deep kernel learning has also been leveraged to improve CSI:FingerID [90], but this approach still fundamentally relies upon hand-crafted input features and cannot be edited, retrained, or fine-tuned by third parties. We hypothesize that neural network approaches are limited by the lack in-domain knowledge incorporated into their architectures. Neural representations that treat peak masses as discrete binned values are unable to model exact masses of peaks, mass differences between peaks, or the putative atomic compositions of peaks, all of which could help improve model generalization.

Here, we introduce Metabolite Inference with Spectrum Transformers (MIST), a neural network approach capable of outperforming current approaches despite not using hand-crafted kernels. Rather than using binned spectra as inputs, MIST represents a spectrum as the set of molecular formulae of all peaks, borrowing inspiration from CSI:FingerID (Figure 2.1D). We then introduce inductive biases from the mass spectra domain: we implicitly featurize neutral loss relationships between fragments, simultaneously predict the structures of the metabolite and its fragments in each spectrum, use *in silico* forward augmentation to provide more training data for our model, and introduce a novel “unfolding” architecture to progressively increase the resolution of fingerprint predictions.

We further enhance MIST’s spectra annotation capabilities by contrastive representation learning on the penultimate network layer. We show the utility of various model components and the inductive biases described above through thorough ablation studies on a publicly available dataset to provide a foundation for future model development. Finally, to demonstrate MIST on real world data, we analyze clinical samples from an inflammatory bowel disease patient cohort and propose new dipeptide and alkaloid structure annotations for differentially abundant metabolites.

MIST is distributed as an open source tool that can be easily integrated into existing pipelines, with or without retraining, and is freely available under the MIT license <https://github.com/samgoldman97/mist>.

2.2 Results

2.2.1 The MIST method

MIST accurately predicts fingerprints and compound annotations by using a deep neural network to learn a meaningful representation of an input mass spectrum. An input spectrum is composed of a single MS1 precursor mass and a list of mass over charge (m/z) peaks with corresponding intensity values. Rather than discretize the observed peaks into a fixed-length binned representation of intensities, MIST instead initially represents each spectrum peak as a molecular formula, which can be pre-labeled (inferred) based upon its mass-to-charge ratio using the SIRIUS algorithm [50] (Figure 2.2A). MIST then directly transforms the representation at each peak by projecting these molecular formulae and intensities into feature vectors via a shallow feed forward multilayer perceptron (MLP).

MIST learns relationships between multiple peaks and their inferred molecular formulae with successive Set Transformer multi-head attention layers [91], [92]. We explicitly encode

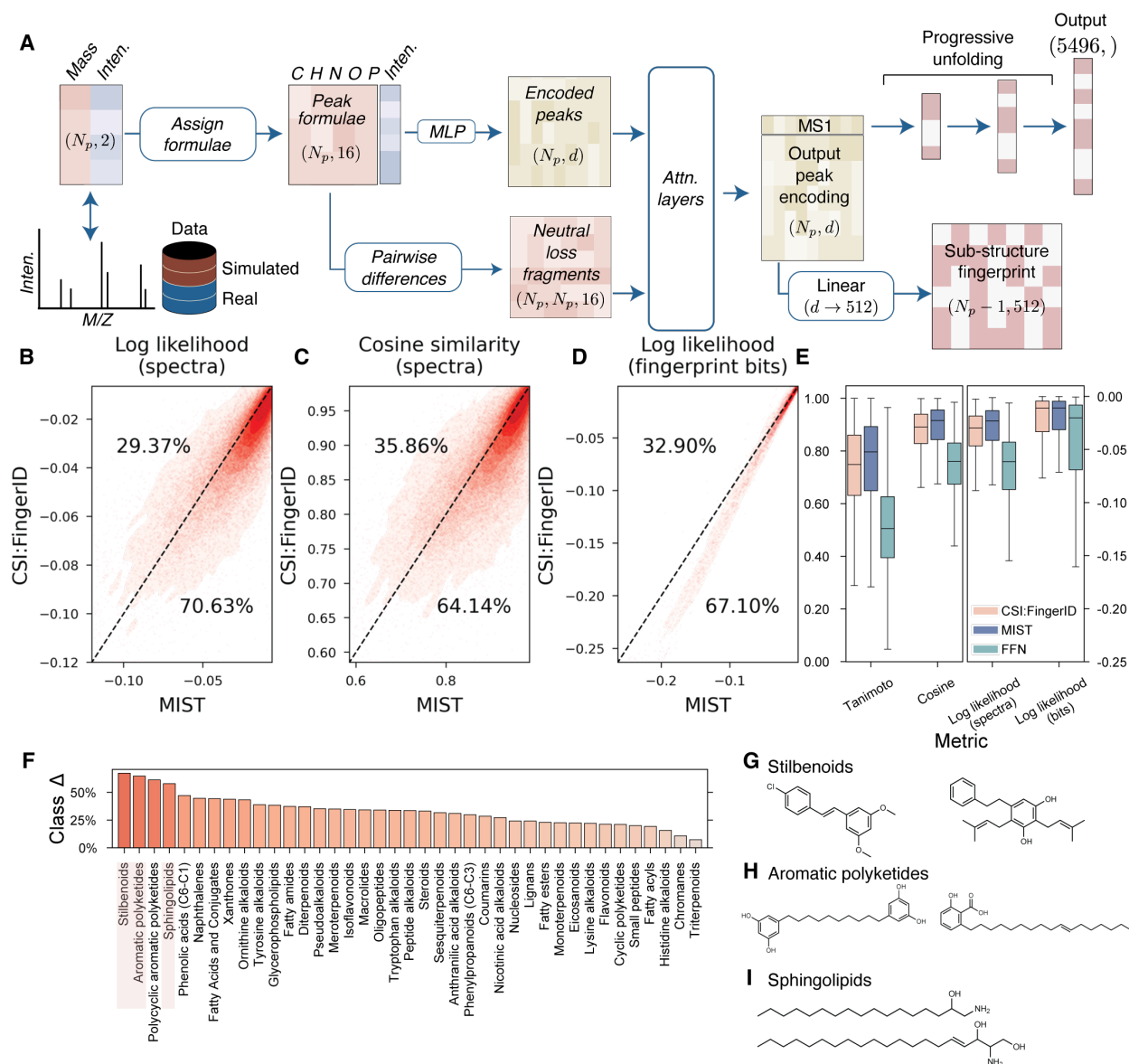


Figure 2.2: MIST accurately predicts compound fingerprints from mass spectra. **A.** Overview of MIST model architecture. An input spectrum of (*mass-to-charge, intensity*) pairs is transformed into molecular formulae vectors, encoded, and fed into a Molecular Formula Transformer along with pairwise neutral losses between formulae. An unfolding module predicts full molecular fingerprints and a separate auxiliary module predicts substructure fingerprints as a secondary training signal. **B,C.** Molecular fingerprints are predicted by MIST and CSI:FingerID for every spectrum in the test set. The performance for each spectrum by log likelihood to the true fingerprint (B) or cosine similarity (C) is evaluated and plotted. Points below the line represent instances where MIST is more performant. **D.** Equivalent evaluation showing the likelihood of predicting each fingerprint bit correctly across all spectra.

Figure 2.2: **E.** The performance of CSI:FingerID [53], MIST, and FFN, a baseline inspired by MetFID [74], are shown. "Tanimoto," "Cosine," and "Log likelihood (spectra)" indicate performance across spectra and "Log likelihood (bits)" indicates performance across bits (higher is better). Median lines are shown; boxes show the interquartile range, with whiskers indicating 1.5x interquartile ranges. Statistics are pooled across $n = 18,700$ test spectra (Tanimoto, Cosine, and Spectra log likelihood) and $n = 5,496$ fingerprint bits (Bit log likelihood). **F.** Performance differences by compound class as assigned by NPCClassifier. All classes with > 40 molecules are shown. **G-I.** Example molecules from stilbenoid, aromatic polyketide, and sphingolipid classes. All results indicate predictions aggregated across 3 separate independent splits and model re-trainings. MIST fingerprints are created by averaging predictions from 5 separately trained models per split.

all pairwise neutral loss relationships [93] as inputs to the attention layer modules. This allows MIST to update the initial representation at each of the peaks based solely on molecular formula with information about all other peaks and the pairwise losses between them. To obtain the final representation of the overall spectrum, we extract the learned representation of the peak corresponding to the precursor formula of the parent compound.

To predict sparse fingerprints from this hidden representation, we introduce a novel "unfolding" layer where MIST progressively predicts the molecular fingerprint with larger and larger resolutions before reaching the final fingerprint prediction. This strategy enables the network to first learn coarse-grained molecular properties before generating the full resolution property vector, similar to how progressive growing has been used to generate high resolution images in computer vision [94].

Simultaneously, we ensure the spectral representation learns to make use of information reflected by each peak in the fragmentation pattern by incorporating a secondary auxiliary loss task. We use the MAGMa algorithm [95] to label chemical substructures for our training data and simultaneously train our model to predict fingerprints for substructures as an additional training signal to improve performance (Figure 2.13).

A final key performance driver to MIST is the use of simulated data, inspired in part by the knowledge distillation strategy employed by AlphaFold [35], [96]. We train a forward neural network to predict spectra from molecules [84] and predict putative new spectra for a database of (unlabeled) biological molecules to randomly add to the training set (Figure 2.9). We refer the reader to Section 2.4.1 for a more complete description of the method.

2.2.2 Formula transformers accurately predict fingerprints

We directly compare MIST to the state of the art fingerprint prediction model, CSI:FingerID, within SIRIUS, to show the benefit of our approach. We prioritize direct comparison with CSI:FingerID as it was the highest performing method at CASMI2022 [75], [76], a prospective competition held in June 2022. MIST is trained on a dataset of positive ion spectra with hydrogen adducts constructed from NIST [63], MONA [97], and the GNPS databases [28], totaling 31,145 unique spectra of 27,797 unique compounds (Section 2.4.8). We train MIST to predict the 5496-dimensional fingerprint used by CSI:FingerID and evaluate performance on 3 separate splits of 20% data holdouts. We ensemble predictions from 5 separately

trained MIST models with different random initializations; we note CSI:FingerID fingerprint prediction is itself an ensemble of 5,496 separately trained models.

Of the 18,700 spectra in the test set holdouts, MIST predictions have higher cosine similarity to the true prediction for 11,994 of the examples (Figure 2.2C), with similarly encouraging results on log likelihood (Figure 2.2B,D). We additionally compare MIST to a more straightforward feed forward neural network (FFN) operating on binned (discretized) spectra, which performs substantially worse than both CSI:FingerID and MIST, indicating the importance of a domain-inspired architecture (Figure 2.2E, 2.8). Even with a single, non-ensembled model, MIST still performs better on over 55% of examples by cosine similarity (Figure 2.14, Table 2.3). This performance advantage from MIST is consistent across chemical classes, with especially strong accuracy on stilbenoids, aromatic polyketides, and sphingolipids (Figure 2.2F-I).

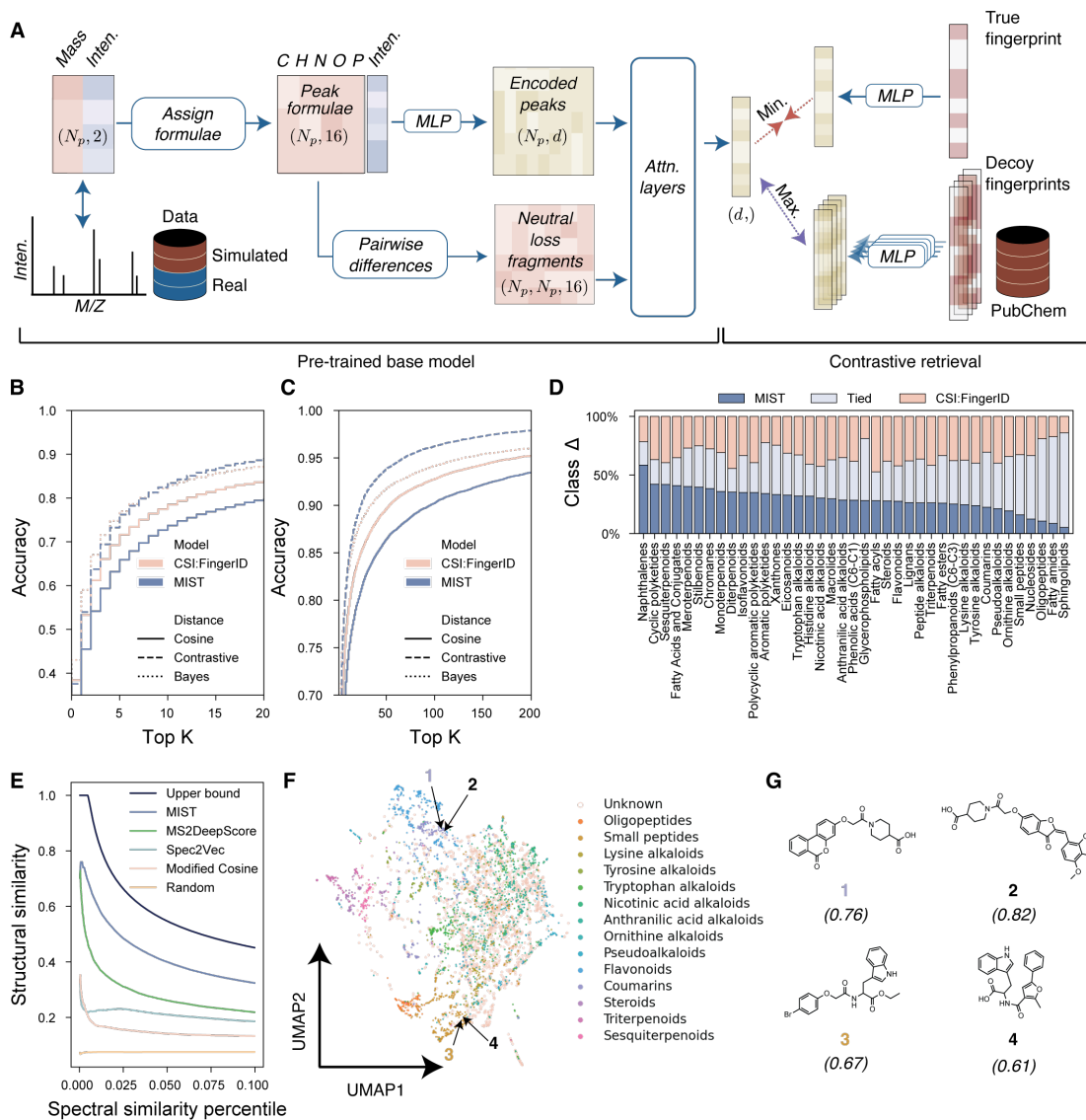


Figure 2.3: **Contrastive fine-tuning improves compound annotation by database retrieval.** **A.** Overview of the contrastive fine-tuning workflow for MIST. Model weights in the formula embedding module and attention layers are initialized from the fingerprint prediction tasks (“Pre-trained base model”) and fine-tuned to project fingerprints into the a latent hidden representation space shared with spectra. The model minimizes distance between the spectra representations and their true structure’s representation and maximizes distance to decoy small molecules. **B,C.** MIST top-k retrieval accuracy assessed against CSI:FingerID by querying the PubChem reference database using predicted fingerprint cosine distances, contrastive latent space distances, or a custom CSI:FingerID fingerprint distance function (“Bayes”). Top-k accuracy is computed by ranking candidate compounds for each spectra and computing the fraction of spectra for which the true compound appears in the top-k rankings, aggregated across 3 test set data folds.

Figure 2.3: **D.** Retrieval accuracy for all folds are segmented into various chemical classes. For each class for which > 40 examples exist, the fraction of examples on which MIST performs better, equivalently, or worse than CSI:FingerID are shown. **E.** MIST latent space spectra distances more effectively cluster high similarity compounds than MS2DeepScore [82], Spec2Vec [81], or cosine distances. Test set compounds pairs are sorted by similarity according to various metrics and the average structural similarity of pairs is computed using Tanimoto distance at various percentiles for a single test fold. **F.** Contrastive spectra latent vectors for a single test fold are projected into 2D space with UMAP and colored by their respective compound class. Compounds 2 and 4 are example spectra for which NPCClassifier fails to return an annotation. **G.** Example compounds from (F) reveal structural similarity of neighboring compounds in the UMAP space derived from MIST’s spectral embeddings. Maximal Tanimoto similarity to any example in the training set is shown in parentheses.

2.2.3 Contrastive fine-tuning improves metabolite retrieval

In principle, one could directly query predicted fingerprints against a database of fingerprints for known biomolecules to retrieve the best matches. CSI:FingerID instead uses custom Bayesian or heuristic distance functions to account for fingerprint bit correlations [98]. We focus on building a learned distance metric by fine-tuning MIST with a contrastive learning objective inspired by noise contrastive estimation [99] (Section 2.4.5) so that molecular fingerprints and spectra can be projected into a shared latent space; similar spectrum-compound pairs are trained to be close in terms of latent space distance (Figure 2.3A).

MIST with contrastive embedding leads to a top-1 annotation accuracy of 37.39% on retrospective datasets, improving upon a 31.70% top-1 annotation accuracy if predicted fingerprints are used to retrieve candidate compounds directly (Figure 2.3B; Table 2.4). Curiously, despite the improved fingerprint accuracy of our model in comparison to CSI:FingerID, we find that on the retrospective dataset, CSI:FingerID has equivalent top-1 accuracy (38.24%) compared to our method and higher top-1 accuracy when retrieving methods using a custom Bayesian fingerprint retrieval method (43.00%) [98]. Nevertheless, MIST achieves higher $k \geq 20$ retrieval accuracy (Figure 2.3C, Table 2.4). This indicates that even when MIST predicts the top-1 molecule incorrectly, it still captures important details of the true compound.

Overall, we find that MIST retrieves the compound on 66% of examples with equivalent or better rank than CSI:FingerID (strictly better on 29% of the data). Analyzing the fraction of examples at which MIST is better by chemical classes, we see that there is not a clear agreement with fingerprint prediction. Whereas MIST seems to have an advantage on fingerprint prediction within sphingolipids (Figure 2.2F), we see that MIST and CSI:FingerID have near equivalent sphingolipid retrieval accuracy (Figure 2.3D). We attribute this to the fact that certain chemical classes have higher top-1% accuracy, possibly due to the bias within the retrieval libraries themselves (e.g., 43/57 sphingolipids are retrieved correctly by both methods). Even so, MIST appears to offer improvements on certain chemical classes such as sesquiterpenoids.

2.2.4 Neural embeddings elucidate compound class clusters

Given the increased retrieval accuracy at higher values of k , we next asked if MIST is capable of meaningfully organizing chemical space in its learned latent space representation of mass spectra. Molecular networking has long been an important tool to cluster spectra that may represent similar molecules, used recently to guide the discovery of new bile-acid conjugate molecules [13], [70]. A challenge of molecular networking is the choice of a meaningful spectral distance and distance cut-off to define network edges.

MIST learns contrastive embeddings with inter-spectra distances that highly correlate with metabolite structure similarity. We compare this effect directly against modified cosine distance as implemented in MatchMS [100], along with MS2DeepScore [82] and Spec2Vec [81], two recent representation learning strategies for embedding spectra (Figure 2.3E). For 38,850,289 pairs of spectra scored by spectral similarity, the average structural similarity of the top 0.1% according to MIST is 0.32 compared to a theoretical maximum of 0.45 using the top 0.1% of most similar pairs, whereas using MS2DeepScore and Spec2Vec yield lower average structural similarities of 0.22 and 0.19 respectively. By this measure, MIST’s spectral similarity metric represents a greater than 45% improvement in structural similarity from the current state of the art.

MIST’s continuous embeddings can also be projected into two-dimensional space for visualization with UMAP, revealing that spectra cluster by chemical class (Figure 2.3F). Even molecules that NPClassifier [101] is unable to annotate are indeed quite structurally similar to their neighbors (Figure 2.3G); while the unlabeled molecule **2** does not share the characteristic coumarin motif with molecule **1**, both molecules share a 4-piperidinecarboxylic acid substructure not detected by chemical classifications.

2.2.5 MIST uncovers putative dipeptides in clinical cohort data

To show the utility of MIST on real world data, we applied MIST to identify unknown metabolites from a recent clinical metabolomics dataset extracted from patients with ulcerative colitis (UC) and Crohn’s disease (CD), two types of inflammatory bowel disease (IBD) as reported by Mills et al. [102]. Mounting evidence has shown that the microbiome mediates IBD treatment response and disease progression, including by the production of novel metabolites [15], [103], [104].

Repeating the authors’ data processing pipeline, we extracted relative abundances for 1,990 tandem mass spectra across the 210-patient clinical cohort (73 UC, 117 CD, 20 healthy) and annotated a total of 1,499 of these spectra using MIST to assign structures from the Human Metabolome Database (HMDB) when possible, falling back on PubChem for molecular formulae without candidate structures in HMDB.

Putative metabolites were grouped into chemical classes using NPClassifier [101] and relative abundances for each chemical class were calculated for each patient. We identified correlation coefficients between the total abundance of each metabolite class and patient disease severity within the UC and CD cohorts. Dipeptides exhibited the highest correlation with disease severity for UC patients of any chemical class ($R = 0.478$, $p_{\text{adjust}} = 0.01$) (Figure 2.4A,B). This is consistent with class-level metabolite analysis from Mills et al. and their finding that increased activity of microbial proteases in dysbiosis led to dipeptide ac-

cumulation [102]. Peptides and amino acids are carbon sources for numerous inflammation-associated microbes, particularly oral cavity residents who have been observed to colonize and thrive in inflamed gut [105]. An understanding of such peptide structures could reveal nutrient niches and footprints of prevalent microbial proteases, possibly suggesting new disease biomarkers or therapeutic targets [106].

A small subset of dipeptide metabolites had higher correlation with disease severity (Figure 2.4C). Curiously, dipeptide metabolites tightly cluster within latent space into different groups (Figure 2.4D); at least four such clusters (compounds highlighted with arrows) contain putative dipeptides that correlate with disease activity in the cohort. Because the chemical class annotation tool is itself a predictive model, certain tripeptides were observed, despite analyzing metabolites labeled as dipeptides.

We further inspected two of the novel annotations to determine whether MIST’s assigned structure can explain the observed spectrum using the CFM-ID [51] *in silico* fragmentation tool at an error threshold of 30 parts-per-million (Figure 2.4F,G). While compound **11** explains 6 high intensity peaks (Figure 2.4F), **13** only explains half of the peaks observed in the spectrum (Figure 2.4G). Given that only four peaks of the six in the spectrum are explained, this example may be a near match, rather than exact annotation, highlighting a potential shortcoming of the model and the importance of manual inspection.

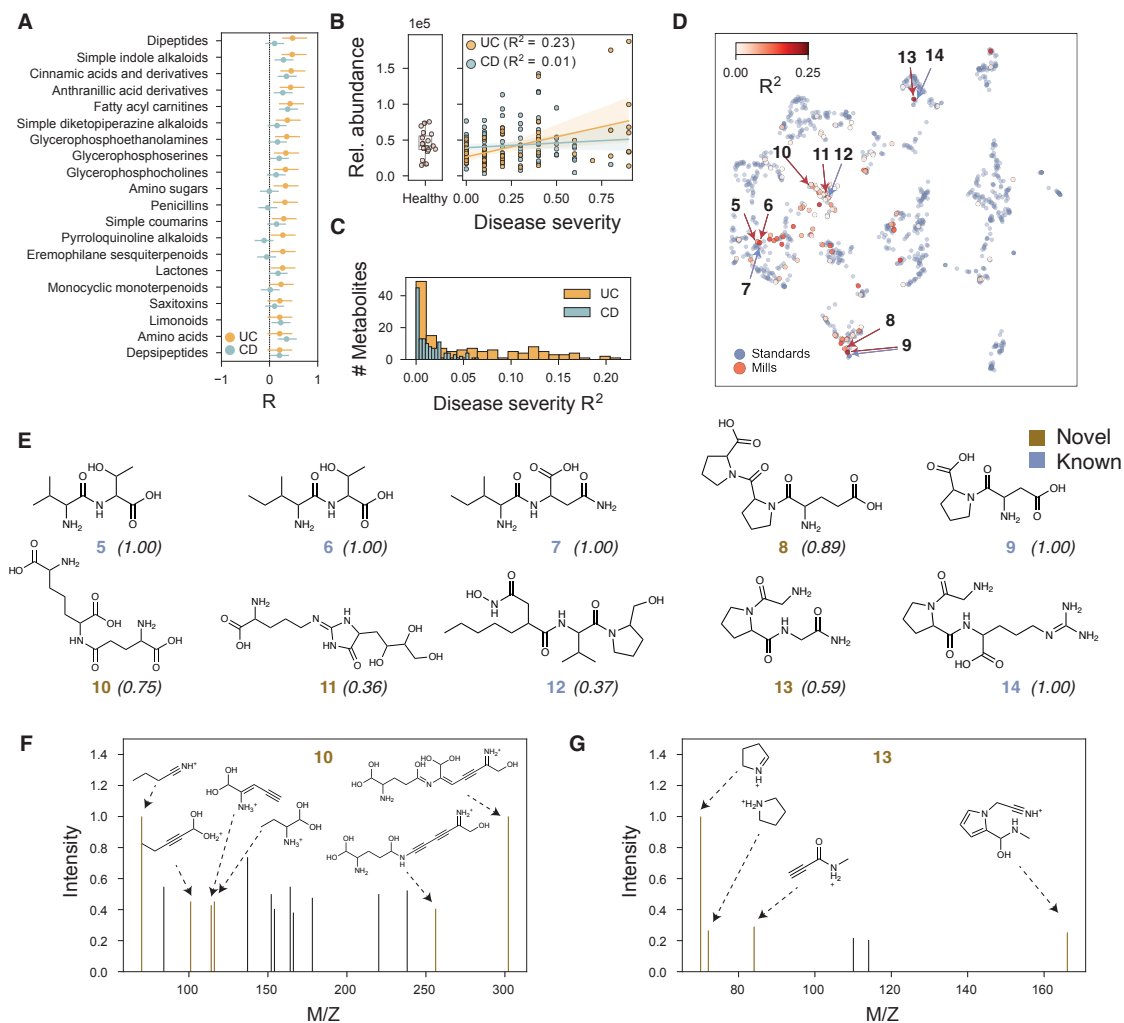


Figure 2.4: MIST annotates putative and clinically relevant dipeptides. **A.** Metabolite classes ranked by their Pearson correlation with IBD disease severity. All putative metabolites are assigned disease classes with NPClassifier [101] and relative abundances are summed across classes. Within the ulcerative colitis (UC) and Crohn’s disease (CD) cohorts, metabolite classes are correlated with disease severity and sorted according to their correlation with UC (top 20 shown; each computed for $n=60$ patients with UC, $n=102$ CD patients). Error bars indicate 95% confidence intervals around the mean. **B.** Dipeptide relative abundances for each patient are plotted against disease severity within the healthy ($n=19$), UC, and CD cohorts. Dipeptides correlate with disease severity for UC patients ($R^2 = 0.23$, $p = 0.01$, two-sided the Wald test, adjusted with Benjamini-Hochberg) as observed by Mills et al. [102]. In the left boxplot, the data median line is shown; boxes show the interquartile range, with whiskers indicating 1.5x interquartile ranges. Regression lines are shown with 95% confidence intervals computed with bootstrapping. **C.** Within the Dipeptide class, the distribution of individual metabolite correlations are shown for both UC and CD cohorts.

Figure 2.4: **D.** All putative dipeptide metabolites (red) are embedded into the latent space of a single, contrastive MIST model alongside embeddings for reference standard spectra from the original model training set and also labeled as dipeptides (blue). All embeddings are projected in 2D space with UMAP and plotted. Metabolites are colored by their coefficient of determination R^2 with UC disease severity as computed in **C.** **E.** Annotations for example metabolites from **D** are shown across different compound clusters. Metabolites with molecular formulae not in standards library are colored brown. **F-G.** Example spectra and their MIST annotated compound are shown. Maximal Tanimoto similarity to any example in the training set is shown in parentheses. Explained subpeaks are annotated using CFM-ID [51] and shown as validation of plausibility. Compound annotations are made using an ensemble of 5 MIST models, contrastive distance retrieval, and the HMDB and PubChem reference databases of molecules.

In addition to finding metabolites that correlate with disease severity, we also searched for metabolite classes that separate the healthy and IBD patients, as immunoregulatory metabolite are continuously being uncovered [107]. We identified two differentially *less* abundant classes that contain N-heterocycles: piperidine (UC $p_{\text{adjust}} = 0.034$, CD $p_{\text{adjust}} = 0.014$) and pyridine (UC $p_{\text{adjust}} = 0.006$, CD $p_{\text{adjust}} = 0.014$) alkaloid classes (Figure 2.5A-C).

Within the piperidine and pyridine alkaloid metabolite classes, we observed a total of 14 and 19 metabolites respectively, the majority of which were annotated as chemically novel structures by MIST (Figure 2.5D). MIST predicted putative compound annotations **15**, **16**, **17**, and **18** as chemically diverse structures with significantly lower abundance than in healthy patients, with compounds **15** and **16** differentially less abundant across both disease cohorts (Figure 2.5E-G). Curiously, not all compounds have the characteristic saturated and unsaturated six-member heterocycles expected for these alkaloid classes; compound **17** instead features a pyrrole ring. We attribute this to the approximate nature of automated chemical class annotation.

This demonstration highlights the position of MIST in common metabolomic workflows. However, we note that our results here are subject to biases both in the retrieval database as well as the molecular formula annotations provided by SIRIUS; no predicted annotations were experimentally verified. We include a more complete discussion of the compound classes for which MIST and SIRIUS share consensus in Table 2.1.

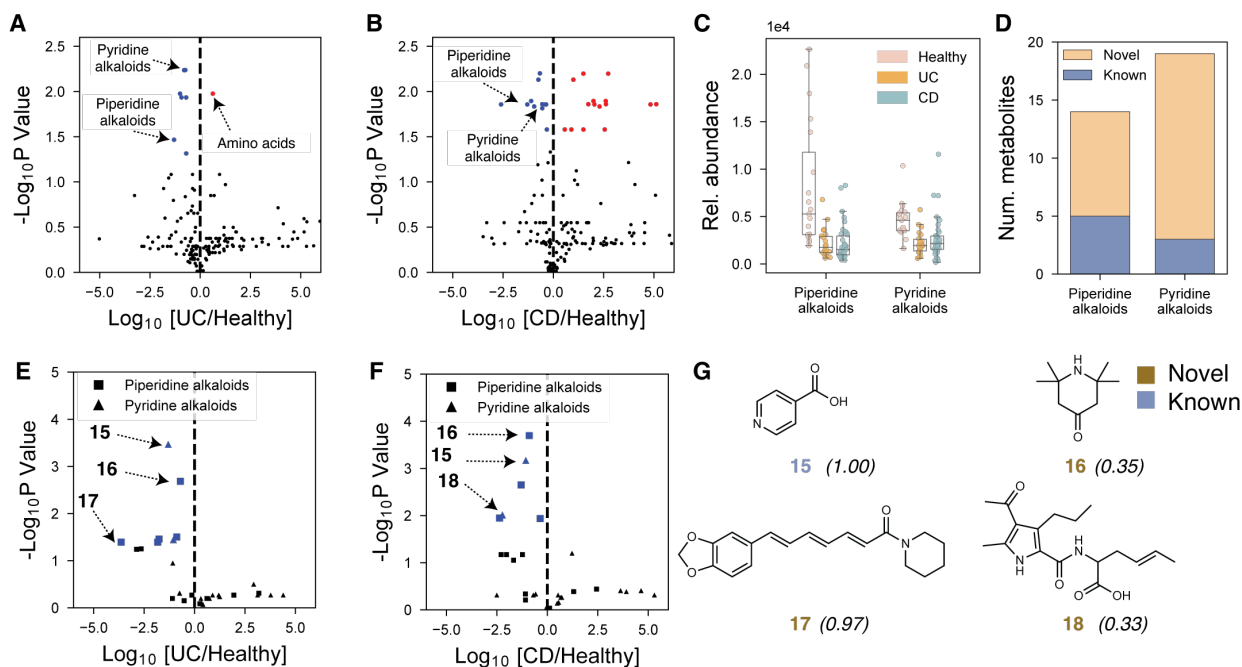


Figure 2.5: Putative alkaloids separate healthy and diseased IBD cohorts. **A,B.** Putative metabolite classes from the Mills et al. cohort [102] are scattered showing their fold change and respective p value comparing and ulcerative colitis (UC) (**A**) and Crohn’s disease (CD) (**B**) cohorts to the healthy cohort. **C.** Distributions of patients’ relative abundances for piperidine and pyridine alkaloids classes are shown for healthy, UC, and CD groups ($n=19$ healthy, $n=22$ UC, $n=38$ CD). Median lines are shown; boxes show the interquartile range, with whiskers indicating 1.5x interquartile ranges. **D.** The total number of putative metabolites in both significant classes are shown, highlighting the number of metabolites with putative novel structures not observed in the standards library. **E,F.** Individual metabolite abundance fold changes and statistical significance are shown comparing the healthy cohort to UC (**E**) and CD cohort (**F**). **G.** Select putative annotated molecule structures are shown for compounds **15** (isonicotinic acid, HMDB:0060665), **16** (triacetonamine, HMDB:0031779), **17** (piperettine, HMDB:0034371), and **18** (unknown). Spectra **15** and **16** are differentially less abundant for both UC and CD. Novel compounds indicate those without spectra standards and are labeled brown; compounds with respective standards are shown in blue. Maximal Tanimoto similarity to any example in the training set is shown in parentheses. Compound annotations are made using an ensemble of 5 MIST models, contrastive distance retrieval, and the HMDB and PubChem reference database of molecules. The top chemical class annotation is found by running the top putative prediction through NPClassifier [101]. P values were computed using independent two-sided t-tests and adjusted for multiple hypothesis testing with the Benjamini Hochberg method. To have equal cohort sizes for healthy and diseased patients, UC and CD cohorts were subsetted to patients with disease severity score > 0.2 . Blue scattered points indicate differentially less abundant compounds or classes; red indicates differentially increased abundance.

2.2.6 Public data benchmarking for future comparison to MIST

A critical impediment to the development of new models for structural elucidation in metabolomics has been the lack of community standards for benchmarking algorithms. While competitions like CASMI [108] provide test set data, they do not specify training datasets. We train a separate set of MIST variants on a publicly accessible and easily downloaded dataset, providing future methods developers with the ability to directly compare to MIST on the prediction of commonplace molecular fingerprints. Using this public dataset, we conduct a more easily reproduced evaluation and ablation study of MIST on a 8030 spectra (7131 unique compounds) subset of the GNPS dataset as prepared by Duhrkop *et al.* [28], [57].

Using this smaller but public dataset, we compare MIST first on the task of fingerprint prediction to variants of MIST that do not include pairwise neutral loss featurization ("MIST - pairwise"), use no auxiliary substructure losses ("MIST - MAGMa"), do not train on simulated fingerprints ("MIST - simulated") or use only a feed forward neural network ("FFN"). We make predictions for 4096-bit circular Morgan fingerprints, which are straightforward to generate using RDKit [109].

As with the larger proprietary dataset, MIST sharply outperforms the FFN model in terms of cosine similarity and log likelihood (Figure 2.6A,B; Table 2.5). Interestingly, we find that outside of using the overall Molecular Formula Transformer, fingerprint unfolding has the largest effect on the performance of cosine similarity (Section 2.5.4). Unlike in the case of CSI:FingerID fingerprint bit predictions (Table 2.6), pairwise embeddings do not appear to improve performance of MIST, which may be a function of the limited public benchmarking dataset size. In further support of dataset limitations, we find that cosine similarity seems to be improving linearly as a function of dataset size in this regime on both the public dataset (Figure 2.6C) and the full, partially proprietary dataset (Figure 2.15A). Additional comparisons to sinusoidal transformer embeddings and different treatments of peak intensities can be found in the (Section 2.5.6; Table 2.7).

We additionally test model performance on retrieval with various scoring functions including the full MIST contrastive function ("MIST contrastive + fingerprint"), only contrastive distance ("MIST contrastive"), only fingerprint distance ("MIST fingerprint"), contrastive distance computed with a model that was not pretrained on fingerprint distance ("MIST contrastive - pretrain"), contrastive distances for a feed forward network ("FFN contrastive") and fingerprints computed with a FFN model ("FFN fingerprint"). The full MIST retrieval using a weighted average of contrastive and fingerprint distances performs better than either distance alone at 30.7% Top 1 accuracy and 75.7% Top-20 accuracy (Figure 2.6E, Table 2.8).

2.3 Discussion

We present a new computational method, MIST, for annotating metabolites in untargeted metabolomics. The MIST architecture borrows insights from years of statistical analysis of mass spectra, such as the use of neutral loss fragments and molecular sub-structures as additional training signals, while remaining sufficiently flexible to adapt to its training data. Because MIST is completely open source, it is straightforward to substitute any proprietary mass spectra training data or target fingerprints to train a new model.

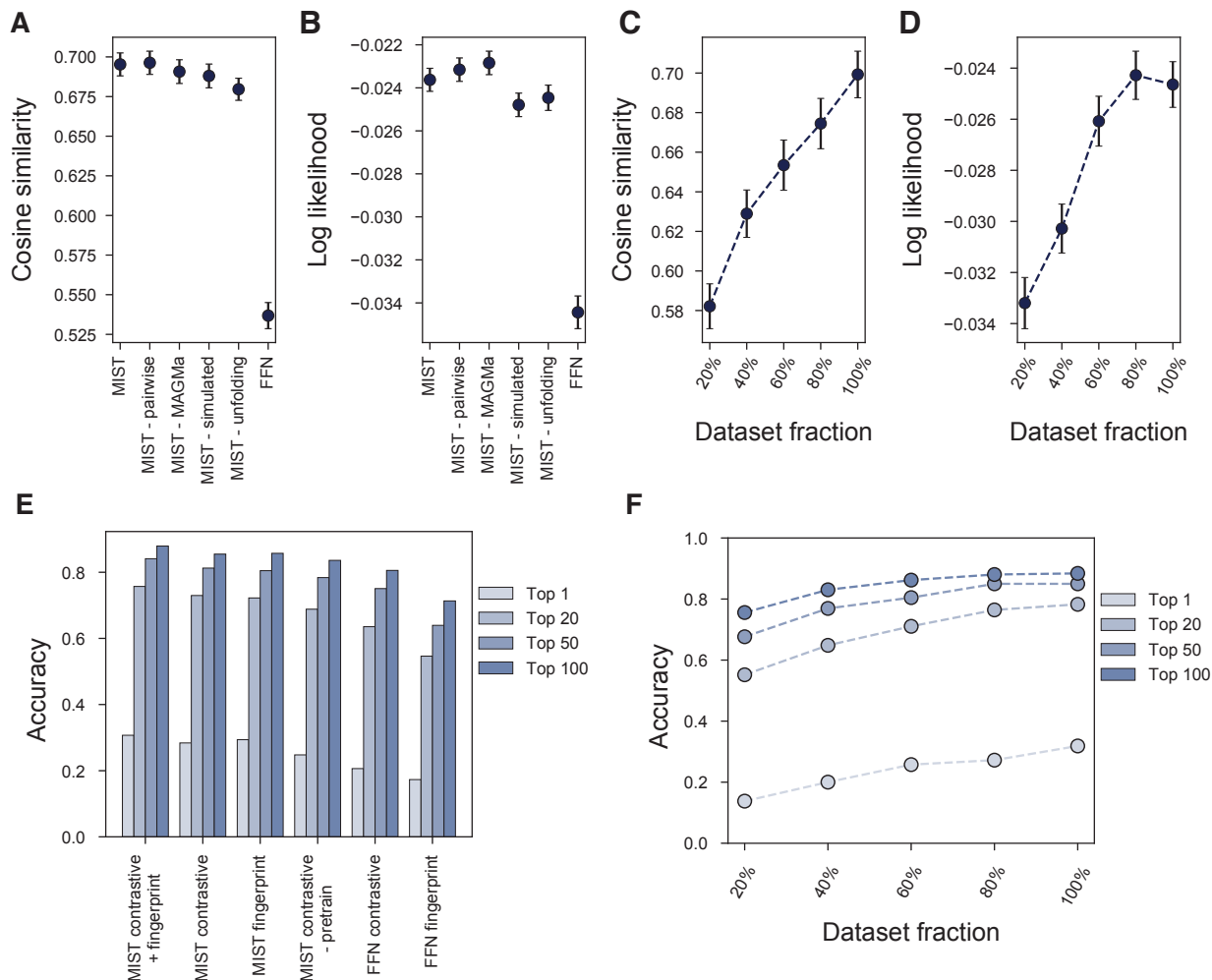


Figure 2.6: **Model ablations affirm the value of domain-inspired model components.** **A,B.** MIST performance on fingerprint prediction compared to ablated variants of MIST without pairwise interactions, without MAGMa substructure supervision, using no simulated spectra, without fingerprint unfolding and using only a feed forward network (FFN) using metrics of **(A)** cosine similarity and **(B)** log likelihood ($n=2,438$ test spectra). **C,D.** MIST fingerprint prediction accuracy improves as a function of dataset size for both **(C)** cosine similarity and **(D)** log likelihood ($n=819$ test spectra). **E.** MIST retrieval performance using a weighted sum of contrastive and fingerprint distance outperforms contrastive distance or fingerprint distance alone, a contrastive version of the model without finetuning, a FFN contrastive model, or FFN fingerprints. **F.** Full MIST contrastive + fingerprint retrieval accuracy improves with training set size. All results are computed on the public GNPS subset data released with CANOPUS. Model ablations are conducted with 3 random splits of the data by molecular structure; dataset ablations are conducted for a single split of the data. Log likelihood values are clamped to a minimum of -5 . Unless otherwise stated, error bars show 95% confidence intervals for the mean.

We demonstrate how MIST can be used both for predicting molecular structure fingerprints and also for learning a meaningful continuous representation, or latent space, with contrastive learning. We showcase these capabilities by applying them in an integrated workflow to annotate metabolites from an untargeted metabolomics study of an IBD patient cohort. Curiously, we document that the improved MIST fingerprint predictions do not lead to strictly better retrieval performance, potentially due to bias introduced by handcrafted CSI:FingerID kernels, the ability of SVM-based kernel methods to predict rare fingerprint bits [110], or other unknown reasons. Increasing our understanding of the retrieval task is an important avenue for future research.

Along these lines, a key difficulty in mass spectrometry model development is the lack of well standardized benchmarks, which have been essential to progress in machine learning tasks. CASMI competitions simultaneously evaluate molecular formula annotation alongside retrieval annotation, making the impact of modeling decisions difficult to fully deconvolute from decisions such as the retrieval database. To facilitate future progress in this field, we provide fully benchmarked model ablations on a small and tractable subset of public GNPS data for both annotation and fingerprint prediction. Such standards will enable better comparisons between models across studies.

A limitation of this work is that MIST is highly dependent on proper molecular formula assignments from MS1 precursor masses. Higher accuracy formula annotations will be synergistic with MIST. In addition, we note that MIST is only trained herein on spectra recorded in positive mode with proton adducts (i.e., H^+), limiting the utility of current trained models. Future work will explore how to model multiple spectra modes, collision energies, and incorporate strategies such as pretraining to further improve prediction quality and expand MIST into a more robust and user-friendly metabolite annotation tool.

MIST provides a competitive neural solution to transform a mass spectrum and predicted molecular formula to a molecular fingerprint or latent space embedding for structural elucidation. Just as protein structure prediction has become powered almost completely by neural network models, we anticipate a similar shift in small molecule structure elucidation pipelines.

2.4 Methods

2.4.1 MIST architecture

MIST model architecture and features are described below. Full hyperparameter searches and exact model parameters are described in the Tables 2.11 and 2.12.

Molecular formula inputs. Each fragmentation spectrum, \mathcal{S} corresponding to a ground truth molecule \mathcal{M} is parsed from a corresponding ‘.ms’ file and contains a precursor mass $M_0(\mathcal{S})$ for the full compound and a list of N_p mass/charge, intensity tuples:

$$\mathcal{S} = [(m_1, I_1), (m_2, I_2), \dots (m_{N_p}, I_{N_p})]$$

MIST assumes that a user has first determined the molecular formula for the precursor compound $M_0(\mathcal{S})$ and removed adduct information (i.e., using SIRUIS, ZODIAC, or other tools [50], [111]). MIST then attempts to assign each fragment peak m_i a molecular formula

subset from the precursor formula, filtering out all masses for which this conversion is not possible. In this work, SIRIUS [50] is used to perform this step but may be replaced with other approaches. Following formula assignment by SIRIUS, peaks are sorted in order of decreasing intensity and only the top 60 most intense peaks are considered for molecular formula labeling.

Formula embedder. Every molecular formula sub-peak is represented as a vector $\mathbf{c} \in \mathbb{R}^{16}$, where each position in the vector is the integer count of atoms for common elements: C, N, P, O, S, Si, I, H, Cl, F, Br, B, Se, Fe, Co, and As, each normalized by the maximum observed atom count of that element in the dataset.

Each formula vector in the spectrum, including the root corresponding to $M_0(\mathcal{S})$, is embedded into a continuous fixed length vector with a shared shallow multi-layer perceptron (MLP) network:

$$\mathbf{h}_i^{(0)} = \text{FormulaEmbedder}(\mathbf{c}_i) \quad (2.1)$$

$$\mathbf{h}_i^{(0)} = \text{Dropout}(\text{ReLU}(\mathbf{W}_2^{\text{emb}} \text{Dropout}(\text{ReLU}(\mathbf{W}_1^{\text{emb}} \mathbf{c}_i)))) \quad (2.2)$$

where $\mathbf{W}_1^{\text{emb}} \in \mathbb{R}^{h \times 16}$ and $\mathbf{W}_2^{\text{emb}} \in \mathbb{R}^{h \times h}$ are learned parameters. Learned bias terms are used but omitted for clarity.

Set Transformer. Each embedded formula is concatenated with its relative intensity (normalized to be maximum 1 for each spectrum) as input to the transformer, $\mathbf{h}_i^{(0)} = [\mathbf{h}_i^{(0)}; I_i]$, where $[a; b]$ denotes a concatenation between vectors/scalars a and b .

A series of multi-headed attention layers [92] as in the transformer neural network architecture are applied to progressively update the hidden representation at each peak.

$$\mathbf{h}_i^{(l+1)} = \text{MultiheadAttention}(\mathbf{h}_i^{(l)}, \{\mathbf{h}_j^{(l)}\}) \quad (2.3)$$

Unlike the original transformer, which uses positional encodings to maintain the ordering of inputs, we remove positional embeddings to make our model equivariant under permutations as in the SetTransformer [91]. Thus, each peak formula attends to other formulae in the spectrum to update its internal representation.

Pairwise attention. The relationship between each peak sub-formula in the spectrum can be explicitly written as the difference in atom counts between the two peaks. We incorporate this pairwise relationship into our model using featurized attention [112], [113].

Before entering the pairwise attention module, an all-by-all formula difference is computed for the input spectrum’s sub-formula list, $\{\mathbf{c}_i\}$, to extract and embed these features using the same FormulaEmbedder module described above. Formula differences are always positive and formula differences are only included when all differences have the same sign (i.e., one sub-formula is a superset of another):

$$\mathbf{f}_{i,j} = \begin{cases} \text{FormulaEmbedder}(\mathbf{c}_i - \mathbf{c}_j) & \mathbf{c}_i > \mathbf{c}_j \\ \text{FormulaEmbedder}(\mathbf{c}_j - \mathbf{c}_i) & \mathbf{c}_i < \mathbf{c}_j \\ \text{FormulaEmbedder}(\mathbf{0}) & \text{otherwise} \end{cases} \quad (2.4)$$

In the original transformer attention layers, attention is computed as

$$\mathbf{K}_j = \mathbf{W}^k \mathbf{h}_j^{(l)}, \quad \mathbf{Q}_i = \mathbf{W}^q \mathbf{h}_i^{(l)}, \quad \mathbf{V}_i = \mathbf{W}^v \mathbf{h}_i^{(l)} \quad (2.5)$$

$$A_{i,j} \propto \mathbf{Q}_i \mathbf{K}_j \quad (2.6)$$

$$\mathbf{h}_i^{(l+1)} = \text{Softmax}(\mathbf{A}_i) \mathbf{V} \quad (2.7)$$

We incorporate attention features by modifying the attention scores to combine formula embeddings with formula difference embeddings as peak and interaction representations, respectively:

$$A_{i,j} \propto (\mathbf{Q}_i + \mathbf{b}^1) \mathbf{K}_j + (\mathbf{Q}_i + \mathbf{b}^2) \mathbf{f}_{i,j} \quad (2.8)$$

Here, \mathbf{b}^1 and \mathbf{b}^2 are trainable biases and $\mathbf{W}^k, \mathbf{W}^q, \mathbf{W}^v$ are learned weight matrices.

Pooling. After iterating with L layers of multiheaded attention, a final hidden representation \mathbf{h}^{out} for the spectrum \mathcal{S} is extracted using the hidden representation $\mathbf{h}_0^{\text{out}} = \mathbf{h}_0^{(L)}$, where the 0th hidden state corresponds to the precursor mass formula. We empirically find that using the hidden representation at the precursor mass formula, rather than mean pooling over all peaks, greatly improves performance.

Loss function. The loss function to train MIST is composed of several different terms related to target fingerprint prediction ($\mathcal{L}_{\text{targ}}$), unfolding of fingerprints ($\mathcal{L}_{\text{unfold}}$), and substructure prediction ($\mathcal{L}_{\text{magma}}$). The relevant secondary terms are weighted by their corresponding hyperparameters λ_{unfold} and λ_{magma} to derive a loss for the full spectra:

$$\mathcal{L}_{\text{MIST}} = \mathcal{L}_{\text{targ}} + \lambda_{\text{unfold}} \mathcal{L}_{\text{unfold}} + \lambda_{\text{magma}} \mathcal{L}_{\text{magma}} \quad (2.9)$$

Each term is defined in depth below.

MAGMa substructure prediction. In addition to training MIST to predict full spectra fingerprints, we train MIST to predict 512-bit Morgan fingerprints for each labeled formula substructure as an additional training signal to regularize the model. For each training example composed of spectrum \mathcal{S} and molecular structure \mathcal{M} , we combinatorially enumerate possible substructures of \mathcal{M} using the MAGMa algorithm to assign putative substructure fragments to peaks in \mathcal{S} (Figure 2.13). We compute a Morgan fingerprint $\mathbf{q}_i \in \{0, 1\}^{512}$ for each of the i^{th} sub-fragments (peaks) [95]. We learn a separate prediction layer $\mathbf{W}^{\text{magma}} \in \mathbb{R}^{512 \times d}$ and apply this to the final representation *at each sub-peak* after the Formula Transformer layers:

$$\hat{\mathbf{q}}_i = \sigma(\mathbf{W}^{\text{magma}} \mathbf{h}_i^{(L)}) \quad (2.10)$$

where $\sigma(\cdot)$ is the sigmoid activation function.

We derive an auxiliary substructure fingerprint loss from these predictions for all labeled sub-structures in each spectrum which empirically enhances performance:

$$\mathcal{L}_{\text{magma}} = \frac{1}{|\{\mathbf{q}_i\}|} \sum_i \mathcal{L}_{\text{fp}}(\hat{\mathbf{q}}_i, \mathbf{q}_i), \quad (2.11)$$

where $\mathcal{L}_{\text{magma}}$ is a clipped cosine loss function, $\mathcal{L}_{\text{fp}}(\hat{\mathbf{y}}, \mathbf{y}) = 1 - \frac{\hat{\mathbf{y}} \cdot \mathbf{y}}{\max(\|\mathbf{y}\| \|\hat{\mathbf{y}}\|, 10^{-8})}$. We note that the number of magma substructure labels is less than the number of peaks, since not all peaks have assigned substructures.

Fingerprint prediction by unfolding. The final representation $\mathbf{h}^{\text{out}} \in \mathbb{R}^d$ is used to predict the target binary fingerprint describing the molecular structure with dimension D :

$\mathbf{y} \in \{0, 1\}^D$. Previous approaches to fingerprint prediction have used simple feedforward networks from some representation of the spectrum, e.g., $\hat{\mathbf{y}} = \sigma(\mathbf{W}\mathbf{h}^{\text{out}})$. However, given the sparsity of our target fingerprint, we instead introduce an unfolding module to progressively grow the fingerprint by predicting ν folded fingerprint intermediates $\hat{\mathbf{y}}^{(\nu)}, \dots, \hat{\mathbf{y}}^{(1)}$, with a fingerprint loss calculation at each step. We deterministically define the folding function to establish targets, $\tilde{\mathbf{y}}^{(1)}, \dots, \tilde{\mathbf{y}}^{(\nu)}$, for each intermediate prediction layer, mirroring how hashed fingerprints like the Morgan fingerprint are folded into fixed-length vectors in the first place:

$$\tilde{\mathbf{y}}^{(0)} = \mathbf{y} \quad (2.12)$$

$$\tilde{\mathbf{y}}^{(i)} \in \{0, 1\}^{\lfloor |\tilde{\mathbf{y}}^{(i-1)}|/2 \rfloor} \quad (2.13)$$

$$\tilde{\mathbf{y}}_k^{(i)} = \min(\tilde{\mathbf{y}}_k^{(i-1)} + \tilde{\mathbf{y}}_{k+\lfloor \tilde{\mathbf{y}}^{(i)} \rfloor}^{(i-1)}, 1) \quad (2.14)$$

In doing so, each progressive target halves in length and bits are merged using the modulo function (i.e., assuming there are K bits in the current fingerprint target, bit k is the 1 if either bit k or bit $k + K$ is 1 in the previous vector). During model training, the network learns to invert this procedure such that the previous intermediate is first expanded then gated using learned weights $\mathbf{W}_i^u \in \mathbb{R}^{|\tilde{\mathbf{y}}^{(i)}| \times |\tilde{\mathbf{y}}^{(i+1)}|}$, $\mathbf{W}_i^g \in \mathbb{R}^{|\tilde{\mathbf{y}}^{(i)}| \times d}$, $\mathbf{W}_\nu^u \in \mathbb{R}^{|\tilde{\mathbf{y}}^{(\nu)}| \times d}$. The lowest resolution target ($\hat{\mathbf{y}}^{(\nu)}$) is predicted first and unfolded progressively:

$$\hat{\mathbf{y}}^{(\nu)} = \sigma[\mathbf{W}_\nu^u \mathbf{h}^{\text{out}}] \quad (2.15)$$

$$\hat{\mathbf{y}}^{(i)} = \sigma[\mathbf{W}_i^u \hat{\mathbf{y}}^{(i+1)}] \cdot \sigma[\mathbf{W}_i^g \mathbf{h}^{\text{out}}] \quad (2.16)$$

We define both a final fingerprint loss and unfolding loss corresponding to each layer:

$$\mathcal{L}_{\text{targ}} = \mathcal{L}_{fp}(\hat{\mathbf{y}}, \mathbf{y}) \quad (2.17)$$

$$\mathcal{L}_{\text{unfold}} = \sum_{i=1}^{\nu} \mathcal{L}_{fp}(\hat{\mathbf{y}}^{(i)}, \tilde{\mathbf{y}}^{(i)}) \quad (2.18)$$

Forward model simulation. A key insight for MIST is to increase the total number of spectra in the dataset and, consequently, the diversity of fingerprints predicted. To do this, we train a "forward" simulator model to predict spectra from candidate molecule fingerprints similar to Wei et al. [84]. Using the same train and test split as MIST, we learn a model mapping from the binary fingerprint \mathbf{y} to a binned representation of the spectrum, $\mathbf{x} \in [0, 1]^{15000}$ that contains 15,000 spectrum intensities from m/z 0 to 1500 with bin spacings of 0.1. The intensity in each bin is normalized to $[0, 1]$ by dividing by the maximum intensity observed in that spectrum (Figure 2.9A). Following Wei et al., we predict both fragment intensities and neutral losses and pool these together (Figure 2.9D):

$$\mathbf{h}^{\text{fwd}} = \text{MLP}(\mathbf{y}) \quad (2.19)$$

$$z_{g,j} = \sigma(\mathbf{W}^o \mathbf{h}^{\text{fwd}})_j, \quad z_{f,j} = \sigma(\mathbf{W}^f \mathbf{h}^{\text{fwd}})_j, \quad z_{r,M_y-j} = \sigma(\mathbf{W}^r \mathbf{h}^{\text{fwd}})_j \quad (2.20)$$

$$\hat{\mathbf{x}} = G(\mathbf{h}^{\text{fwd}}) = \mathbf{z}_g \cdot \mathbf{z}_f + (\mathbf{1} - \mathbf{z}_g) \cdot \mathbf{z}_r \quad (2.21)$$

Here, MLP defines a shallow feed forward network transforming the input fingerprint into a hidden representation. \mathbf{W}^f and \mathbf{W}^r transform the hidden representation into predictions

of peaks and losses. Losses are converted into peaks by subtracting each predicted loss bin from M_y , the bin of the full molecule mass. \mathbf{W}^o transforms the hidden input into a set of gates to balance the forward and reverse prediction terms, \mathbf{z}_f and \mathbf{z}_r .

To further increase the performance of the forward simulator, we use a similar unfolding strategy as utilized in MIST for fingerprint prediction (Figure 2.9C). We apply several repeated unfolding layers to predict γ binned spectra with lower dimensions at each step $\tilde{\mathbf{x}}^{(1)}, \tilde{\mathbf{x}}^{(2)}, \dots, \tilde{\mathbf{x}}^{(\gamma)}$. We define these by binning spectra targets with progressively lower resolution:

$$\tilde{\mathbf{x}}^{(0)} = \mathbf{x}, \quad \tilde{\mathbf{x}}^{(i)} \in [0, \infty]^{\lfloor |\tilde{\mathbf{x}}^{(i-1)}|/2 \rfloor}, \quad \tilde{\mathbf{x}}_j^{(i)} = \tilde{\mathbf{x}}_{2j}^{(i-1)} + \tilde{\mathbf{x}}_{2j+1}^{(i-1)} \quad (2.22)$$

To make predictions, we first predict the lowest resolution $\hat{\mathbf{x}}^{(\gamma)} = G_\gamma(\mathbf{h}^{\text{fwd}})$, where G_γ is a function that combines predictions for both fragment and neutral loss intensity predictions, defined above. Following our strategy from fingerprint prediction, one strategy would be to predict the higher resolution spectrum $\hat{\mathbf{x}}^{(i)}$ using a linear transformation of the coarser prediction $\hat{\mathbf{x}}^{(i+1)}$. However, given the wide dimensionality, such learned transformations would have shape $\mathbf{W} \in \mathbb{R}^{|\tilde{\mathbf{x}}^{(i)}| \times |\tilde{\mathbf{x}}^{(i-1)}|}$. This is prohibitive with large bin sizes. Instead, we alternatively elect to naively expand the representation $\hat{\mathbf{x}}^{(i+1)}$ by setting $\hat{x}_j^{(i)} = \hat{x}_{\lfloor j/2 \rfloor}^{(i+1)}$ without any additional parameterization. We update this prediction as a function of the original hidden representation, once again using the same generalized neural network layers G , defined above, where G_i^p is used to predict $\tilde{\mathbf{x}}^{(i)}$ as a function of \mathbf{h}^{fwd} and G_i^s predicts a multiplicative factor to weight the expanded prediction vs. the new prediction:

$$\hat{x}_j^{(i)} = \sigma(G_i^s(\mathbf{h}^{\text{fwd}}))_j \hat{x}_{\lfloor j/2 \rfloor}^{(i+1)} + (1 - \sigma(G_i^s(\mathbf{h}^{\text{fwd}}))_j) \sigma(G_i^p(\mathbf{h}^{\text{fwd}}))_j \quad (2.23)$$

Spectra are filtered to include peaks at only the positions for which chemical sub-formula are labeled and sufficient intensities are observed. Unlike Wei et al., to train this model we use a binary cross entropy loss for the presence of certain fragments, rather than a cosine similarity loss to reduce the influence of low intensity peaks [84]. We define a binarized spectrum $\mathbf{x}' = \mathcal{K}(\mathbf{x} > 0)$ and also a weighting factor $s(\mathbf{x}) = \text{Softmax}(\mathbf{x})$ to enforce that peaks with higher intensities are weighted more heavily in our loss function:

$$\mathcal{L}_{\text{BCE}}(\hat{\mathbf{x}}, \mathbf{x}) = \frac{1}{D} \sum_{d=1}^D s(\mathbf{x})_d [x'_d \log[(\hat{x}_d)] + (1 - x'_d) \log[1 - \hat{x}_d]] \quad (2.24)$$

We apply this BCE loss at each unfolding layer step to get an unfolding loss and final target loss, weighted by a factor λ_{fwd} to define the full forward loss:

$$\mathcal{L}_{\text{fwd_targ}} = \mathcal{L}_{\text{BCE}}(\hat{\mathbf{x}}, \mathbf{x}) \quad (2.25)$$

$$\mathcal{L}_{\text{fwd_unfold}} = \sum_{i=1}^{\gamma} \mathcal{L}_{\text{BCE}}(\hat{\mathbf{x}}^{(i)}, \tilde{\mathbf{x}}^{(i)}) \quad (2.26)$$

$$\mathcal{L}_{\text{fwd}} = \lambda_{\text{fwd}} \mathcal{L}_{\text{fwd_targ}} + (1 - \lambda_{\text{fwd}}) \mathcal{L}_{\text{fwd_unfold}} \quad (2.27)$$

To verify the utility of the reverse prediction layer and the unfolding layer in the forward model architecture, we conduct a brief ablation study on the public dataset. We evaluate

our model by both the fraction of true spectrum peaks each predicted spectrum covers (“Coverage”) and cosine similarity. We find the full feed forward model outperforms on both metrics, with coverage of 0.58 vs. 0.56 without unfolding and 0.52 without the reverse layer (Table 2.10).

After training this model, we make spectra predictions for a randomly sampled set of biomolecules while excluding the validation and test set molecules to avoid accidentally biasing our model. We note that in additional experiments, including these molecules in the augmentation had a surprisingly strong effect in biasing our model toward the true predictions at test time, highlighting the ease for data leakage to inflate performance in this setting. For each predicted spectrum, we convert the binned spectrum back to sub-formula labels using combinatorial enumeration and using a maximum of 50 labeled sub-formulae with the highest predicted intensity. We use a prediction threshold of 0.2 to call peaks (Figure 2.9B).

During training, in each epoch we include the entire training set of experimental spectra. We augment the training set in that batch to include a fraction of $1 - f_{\text{real}}$ simulated spectra, with f_{real} set during hyperparameter optimization, such that experimental spectra only make up a fraction f_{real} of each epoch.

Spectra noising. To make the model robust to peak shifts, each fragmentation spectrum is noised using a similar protocol to MetFID [56]. During training, with a $p_{\text{noise}} = 0.5$ probability any given spectrum is augmented. Peaks within augmented spectra are randomly selected for either intensity rescaling or removal. Full pseudocode is listed in Section 2.5.7.

2.4.2 Fingerprints

To directly compare against CSI:FingerID, we train models using the custom length 5,496 fingerprints included with SIRIUS as provided by the authors and computed using the Docker image provided for MSNovelist [61]. For ablation datasets on public data, we use RDKit to compute 4096-bit circular Morgan fingerprints with radius 2 [109], [114].

2.4.3 SIRIUS software

Throughout this work, we use the SIRIUS [50] software version 4.9.3 for molecular formula annotation and direct comparison unless explicitly stated otherwise.

2.4.4 Feed forward baseline

We train a feed forward network baseline similar in concept to MetFID that maps binned spectra to fingerprints [56]. We bin mass values for each spectrum into n_{bins} equally spaced bins from 0 to 1,500. n_{bins} is set in hyperparameter optimization. The model is trained using cosine similarity loss between the predicted fingerprint and the true fingerprint as in MIST’s architecture.

2.4.5 Contrastive fine-tuning for database retrieval

After pre-training MIST to predict the true fingerprint \mathbf{y} , we extract the model weights used to produce the hidden representation $\mathbf{h}^{\text{out}} = F(\mathcal{S})$ and treat this internal representation as a contrastive space to conduct retrieval, where F defines the pretrained MIST feature extractor. We learn a separate single layer projection $\mathbf{W}^c \in \mathbb{R}^{d \times |\mathcal{Y}|}$ to map fingerprints into the contrastive space.

We sample a set of n_{decoy} decoy training fingerprints, $\mathbf{y}'_1, \mathbf{y}'_2, \dots, \mathbf{y}'_{n_{\text{decoy}}}$ for each spectrum and learn to project these into the latent space such that the true target fingerprint maps closer to the latent \mathbf{h}^{out} than to the decoys. We utilize a softmax contrastive loss inspired by noise contrastive estimation with temperature parameter τ and cosine similarity function [99], [115]:

$$\text{sim}_\theta(\mathcal{S}, \mathbf{y}) = e^{\tau_s \cos(F(\mathcal{S}), \mathbf{W}^c \mathbf{y})} \quad (2.28)$$

$$\mathcal{L}_{\text{NCE}}(\mathcal{S}, \{\mathbf{y}, \mathbf{y}'_1, \dots, \mathbf{y}'_{n_{\text{decoy}}}\}) = -\log \frac{\text{sim}_\theta(\mathcal{S}, \mathbf{y})}{\text{sim}_\theta(\mathcal{S}, \mathbf{y}) + \sum_{n=1}^{n_{\text{decoy}}} \text{sim}_\theta(\mathcal{S}, \mathbf{y}'_n)} \quad (2.29)$$

To sample decoys, we utilize the PubChem database (April 2022) to identify isomers with high similarity to the ground truth for all spectra in the training set. We calculate the 256 closest examples for each isomer by Tanimoto similarity using Morgan fingerprints with length 4096 and radius 2. In each batch, these are sampled proportional to an exponential of their their Tanimoto similarity to the target:

$$p(\mathbf{y}'_i; \mathbf{y}) \propto e^{\tau_s \text{Tani}(\mathbf{y}'_i, \mathbf{y})}$$

Similarity is weighted by an additional temperature factor $\tau_s = 4$ enforcing that decoys with high Tanimoto similarity are sampled more often.

We include forward simulated spectra in contrastive fine-tuning as well. Rather than sample isomeric decoys from PubChem for these simulated examples, we extract decoy fingerprints from the set of all forward simulated molecules to avoid further additional computational complexity. We again sample decoys using Tanimoto similarity for hard negative mining.

During contrastive learning we train the network using a weighted sum of the contrastive loss and MIST loss with a hyperparameter λ_c :

$$\mathcal{L}_{\text{contrastive}} = \lambda_c \mathcal{L}_{\text{NCE}} + (1 - \lambda_c) \mathcal{L}_{\text{MIST}}$$

2.4.6 Model training

All models are trained with the RAdam [116] optimizer. Learning rate and weight decay hyperparameters are set for each dataset and detailed in Section 2.5.8. A validation set containing 5% of the data is excluded from the training set and utilized for early stopping during training with a patience of 20.

2.4.7 Retrieval

Fingerprint retrieval. Retrieval results are calculated using the PubChem database (April 2022) unless otherwise stated. For each spectrum, a fingerprint is first predicted. All isomers matching the spectrum’s molecular formula are rank ordered by cosine similarity between the isomer fingerprint and the predicted fingerprint. For all ties, the optimistic lower rank of the tied options is chosen. Ties are broken by selecting the minimum rank.

Bayesian retrieval. CSI:FingerID and SIRIUS internally utilize an adjusted distance, rather than cosine similarity, that corrects for correlations between fingerprint bits [98]. We utilize rankings calculated with Bayes as provided for CSI:FingerID predictions by the authors for direct comparison.

Contrastive retrieval. Contrastive retrieval results are computed by projecting fingerprints and spectra into the same continuous latent space using the trained contrastive model. We use a cosine distance metric in the latent space to rank isomers as in fingerprint retrieval. We ensemble contrastive distance and fingerprint distances with hyperparameter $\lambda_r = 0.3$, controlling the relative weights of fingerprint distance $d_{fp}(\mathcal{S}, \mathcal{M})$ and contrastive distance $d_c(\mathcal{S}, \mathcal{M})$: $d(\mathcal{S}, \mathcal{M}) = \lambda_r d_{fp}(\mathcal{S}, \mathcal{M}) + (1 - \lambda_r) d_c(\mathcal{S}, \mathcal{M})$. This parameter was manually tuned on a single data fold.

2.4.8 Training datasets

CSI:FingerID NIST comparison. We directly compare against CSI:FingerID using an $[M+H]^+$ adduct dataset compiled from public data sources including the GNPS, MONA, and NIST20 [28]. This dataset has a total of 31145 spectra and 27797 unique compounds. We use cross validation, training on 3 independent structure-disjoint splits of the data as computed and provided by the CSI:FingerID authors [53]. The structure-disjoint dataset ensures that no two molecules with the same 2D InChiKey will be split across different folds, but does not restrict maximal Tanimoto similarity across splits. An in depth discussion of data splits can be found in Section 2.5.3.

GNPS CANOPUS benchmark. We utilize a public subset of the data including spectra pulled from the GNPS and MONA. This dataset is further filtered to $[M+H]^+$ spectra and contains a total of 8,030 spectra and 7,131 unique compounds. We create structure-disjoint splits as above (i.e., unique at the level of 2D InChiKeys) of this dataset using varying fractions of the dataset in the training set, 0.2, 0.4, 0.6, 0.8. For each of these settings, 5% of the training set is randomly sampled and set aside for validation to perform early model stopping. This dataset is freely available and provided with the CANOPUS software tool [57].

Biomolecules. For forward simulation, we use biomolecules extracted from KEGG, KNAPSACk, HMDB, and others [117]–[119], again prepared by Duhrkop *et al.* [57]. This dataset is only provided with InChiKey values, which cannot be mapped to molecular structures in the absence of a database lookup. We query PubChem to assign structural identities for a total of 997,554 unique structures. We randomly sample this dataset to include 300,000 molecules (10x the full training dataset). Further, in each split, we exclude any molecules also appearing in the test and validation sets.

Compound retrieval libraries. PubChem was downloaded in its entirety [58] in April 2022 in order to build retrieval HDF5 files for chemical isomers. Comparisons to

CSI:FingerID using PubChem as a retrieval library were made using a version of the PubChem database provided by the CSI:FingerID authors at the same time. HMDB 5.0 [59] was similarly downloaded.

2.4.9 Clinical data reanalysis

Stool metabolite samples were previously collected by Mills et al. [102]. Metabolomic data was accessed under MassIVE accession number MSV000084908, including patient metadata at MSV000086509. Mass spectrometry data was extracted and converted into relative abundance tables and MS2 metabolite files using MZMine 3 [120] to be used by SIRIUS (molecular formulae identification) and MIST (structure elucidation). MZMine 3 parameters were set as specified by Mills et al., and an MS2 ppm tolerance of 10 was used for SIRIUS to identify molecular formulae.

MIST annotations were made using ensembles of 5 models. Structural annotations were first retrieved from the HMDB compound library; for all spectra without isomers in HMDB (i.e., no molecules in HMDB share the putative molecular formula), PubChem was utilized as a compound library for annotations.

Correlation coefficients with disease severity were computed within UC and CD cohorts for all patients that were assigned disease activity levels. Healthy-UC and healthy-CD comparisons were computed using independent two sided t-tests and adjusted with the Benjamini-Hochberg method. To utilize more equal sized cohort comparisons, UC and CD were subsetting to include patients with disease severity > 0.2 down to 40 and 22 patients respectively when comparing to healthy control groups.

2.4.10 Chemical classifications

NPClassifier [101] version 1.5 was used to produce chemical classifications of each compound. We query the public and easily accessed GNPS web server endpoint [28] and extract the superclass definition for chemical classification. Compounds with no predicted annotation are labeled as "Unknown."

2.4.11 Latent distance comparisons

Spectra are projected into a shared spectrum-structure latent space using the fine-tuned contrastive MIST model. To assess the extent to which spectral similarity accurately reflects structural (Tanimoto) similarity, we repeat analyses from Spec2Vec [81] within a single fold of test data. We calculate all pairwise spectral similarities using cosine similarity and sort from highest to lowest similarity. At each threshold, we compute the average Tanimoto structural similarity. We compare to an arbitrary theoretical upper bound if paired spectra were sorted by Tanimoto similarity directly.

Spec2Vec. A pretrained Spec2Vec version 0.8.0 [81] model is downloaded and directly used to embed all spectra into latent space. We repeat the same analysis above using Cosine similarity in Spec2Vec latent space.

MS2DeepScore. The pretrained MS2DeepScore version 0.4.0 [82] model is similarly downloaded and directly used to embed all spectra into latent space before computing cosine

distances.

Modified cosine similarity. Modified cosine similarity [70] is computed between spectra themselves using the `MatchMS` Python library [100] version 0.17.0.

UMAP computation UMAPs for projected compounds are computed using the UMAP Python package applied to latent contrastive spectra embeddings in the test set [121].

2.4.12 Code and data availability

Public data used for benchmarking MIST models as processed by Duhrkop et al. [57] can be downloaded alongside our code with full directions included at <https://github.com/samgoldman97/mist>. Data for NIST and head-to-head CSI comparisons is unavailable due to strict licensing rules around the NIST20 [63] dataset. Data to repeat the retrospective study and reanalysis of IBD data can be retrieved from the MASSIVE database at accessions MSV000084908 (raw data) and MSV000086509 (cohort info) and via Zenodo Record 8084088. PubChem (April 2022) and HMDB 5.0 data libraries used for compound retrieval are publicly accessible with exact details for reproduction described alongside released code.

All code to replicate experiments, train new models, and load pre-trained models is available at <https://github.com/samgoldman97/mist>.

2.5 Additional Results

2.5.1 Additional model details

Unfolding layers. Novel unfolding modules are introduced in MIST to progressively grow the fingerprint prediction. Target fingerprints are compressed using modulo operations to derive "lower resolution" targets. These lower dimensional vectors inherently have collisions for properties. That is, rather than each bit indicating the presence of a single property or motif, each bit represents the presence of one of several properties (Figure 2.7).

Binned feed forward baselines. A simple feed forward network inspired by MetFID is utilized as a baseline (Figure 2.8) [56].

2.5.2 CSI:FingerID annotations for Mills et al. cohort

In Section 2.2.5, we conduct a prospective analysis to showcase how MIST can be used with real world data. To evaluate the reasonableness of our annotations against the existing tool, SIRIUS and CSI:FingerID [50], we repeat the annotation process with these tools. Concretely, we use an equivalent MGF spectrum extraction process and pass this into SIRIUS Version 5.6.3 for end-to-end formula, fingerprint, and structure annotation. We use SIRIUS Version 5.6.3 as SIRIUS Version 4.9 was deprecated during the course of this work.

We constrain the compound search to the HMDB database [117] to match our MIST pipeline. SIRIUS 5 stalled out and failed on the full MGF file including larger compounds; correspondence with the authors revealed that larger fragmentation trees can lead to memory errors. As such, we subsetted the dataset to feature only the 1,515 spectra with precursor ions under 500 Da. 545 compounds were annotated with structures, of which 342 had equivalent

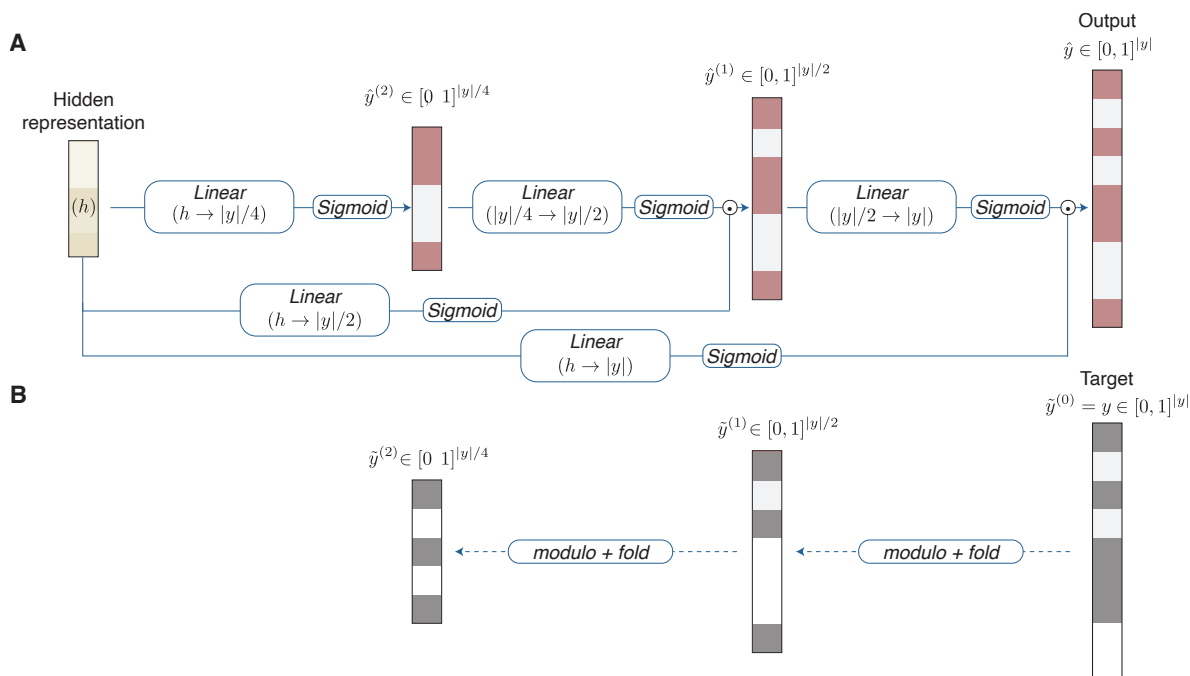


Figure 2.7: **Unfolding fingerprint prediction module.** **A.** The module starts with a hidden representation and subsequently predicts increasing resolutions of the fingerprint. **B.** At each intermediate fingerprint prediction $\hat{y}^{(i)}$, a corresponding target vector is generated by “folding” the previous, higher resolution fingerprint. An intermediate loss is calculated at each of these resolutions.

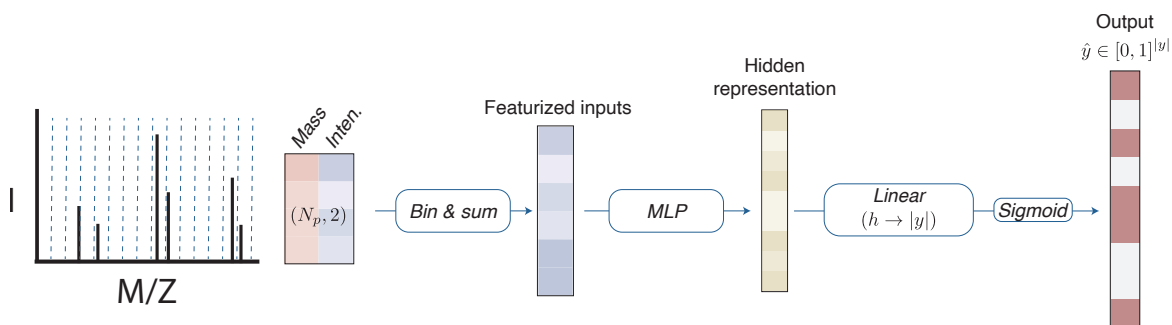


Figure 2.8: **Feed forward network binned spectra baseline.** An input spectrum is binned, a multilayer perceptron (MLP) is applied, and the output is projected to a fingerprint prediction.

molecular formulae to those we extracted with SIRIUS 4.9 earlier. From this set, MIST predicted equivalent InChIKeys on 58% of the spectra; a random selection from HMDB isomers matched SIRIUS on 32%. We include a full table showing drop-off at each step in the processing pipeline (Table 2.1).

Surprisingly, we observe in this workflow that the molecular formulae of compounds predicted by SIRIUS are different from molecular formulae we exported from SIRIUS as inputs

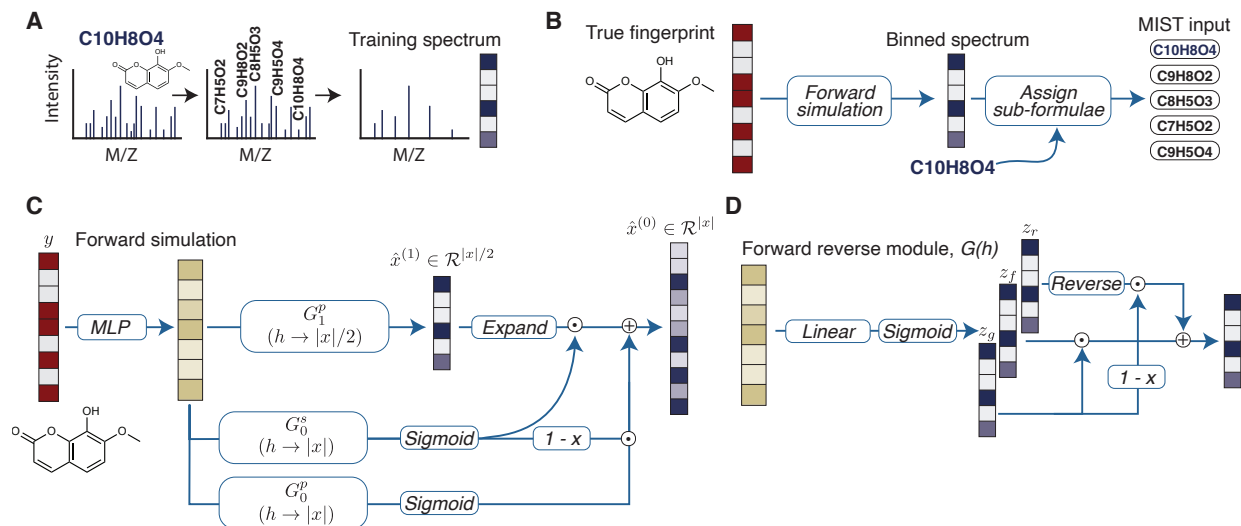


Figure 2.9: **Forward simulation of spectra from molecules.** **A.** An input spectrum, molecule pair is first denoised and cleaned by labeling chemical sub-formulae. This is converted into a binned spectra representation. **B.** To generate new model training examples, molecules are fingerprinted, binned spectra predictions are predicted by a neural network, and sub-formula are assigned to each of the bins such that the spectra can be used as input to MIST. **C.** The forward simulation neural architecture utilizes unfolding layers. A multilayer perceptron (MLP) projects the molecular fingerprint into a hidden representational space. Forward-reverse models, G , are used to progressively grow the binned prediction such that resolution increases at each step, similar to unfolding layers. **D.** Illustration of the forward-reverse module, where intensities (z_f) and neutral losses (z_r) are predicted. Neutral losses are reversed and mapped onto intensity bins using the full mass as input then summed according to a learned gate to get a single prediction of intensities in each bin.

to MIST. This is due to SIRIUS *re-ranking* molecular formulae during the compound retrieval phase of their pipeline. Testing this, we run SIRIUS again with PubChem as a database retrieval and find that of the 1,355 compounds that are annotated with molecular formulae in the initial step, only 625 of these have the same molecular formulae after compound retrieval.

To probe the robustness of our findings at the chemical class level, we repeat the MIST processing pipeline with two changes: (1) we utilize 640 spectra and formula annotations from the final step of SIRIUS’s pipeline, after re-ranking during retrieval from PubChem and (2) utilize PubChem as a retrieval database. Under this analysis, we find 162 of the 640 spectra have equivalent annotations under SIRIUS and MIST, a greater than 25% agreement even when searching the much larger PubChem database.

Using the resulting 640 revised annotations outputted from MIST, we repeat our prospective IBD analysis pipeline and find that dipeptides are still one of the most correlated classes with disease severity (Figure 2.10). Further, when looking for differential abundance, we see that piperidine and pyridine alkaloids, are still differentially abundant in CD patients (Figure 2.11). These class level differences are no longer present for UC patients. This discordance

emphasizes that automated annotation pipelines are not yet robust to changes in processing parameters such as spectrum preprocessing, formula annotations, fingerprint prediction, and retrieval library selection.

In the future, moving beyond fragmentation tree based predictions for formula annotation will enable efficient annotation of larger compounds. Improving automated molecular formula annotation and database retrieval library construction will similarly lead to higher quality and more consistent annotations.

Table 2.1: **Comparison of SIRIUS and MIST for annotating IBD cohort data.**

| | Total spectra |
|--|---------------|
| MIST Pipeline | 1,990 |
| SIRIUS Pipeline (< 500 Da) | 1,515 |
| MIST identifies a structure in HMDB | 705 |
| SIRIUS identifies a structure in HMDB | 724 |
| SIRIUS and MIST both identify a structure in HMDB | 545 |
| SIRIUS and MIST annotate compounds with equiv. formula | 342 |
| SIRIUS and MIST find same InChIKey | 200 |
| SIRIUS and Random-HMDB find same InChIKey | 110 |

2.5.3 Dataset splits

For both the proprietary NIST data and public benchmarking data, care is taken to ensure that no two molecules with the same InChIKey appear across separate dataset split folds. We can quantify and demonstrate this effect by computing the nearest neighbor compound (according to Tanimoto similarity) in the training set for all molecules in the dataset. By plotting the distributions of the “maximum similarity neighbor” for both datasets, we can see that the training set features equivalent molecules (i.e., a neighbor with similarity 1.0) with a distribution mean of 0.69 and 0.71 Tanimoto similarity across the private and public datasets. In comparison, the test set distributions are more normal with means of 0.63 and 0.66 across the two datasets respectively (Figure 2.12).

2.5.4 Unfolding layer ablation

MIST contains unfolding layers for fingerprint prediction. We find that this leads to improved performance via an ablation study. In our ablation study, we replace the unfolding module with a single linear layer output to predict fingerprints. To verify that the unfolding mechanism (rather than just more layers) is driving performance, we add a shallow 3 layer FFN to the top of the *MIST - unfolding* model to match the 4 unfolding layers in the Unfolding module, yielding a total of 15.1M parameters vs. 3.1M parameters for *MIST* and 2.8M parameters for *MIST - unfolding*. We differentiate the two MIST ablation models as “MIST - unfolding (linear)” and “MIST - unfolding (FFN)” respectively. To avoid high computational cost of additional model features, we train these models without the auxiliary MAGMa loss and without simulated spectra. As can be seen in Table 2.2, we find that without unfolding,

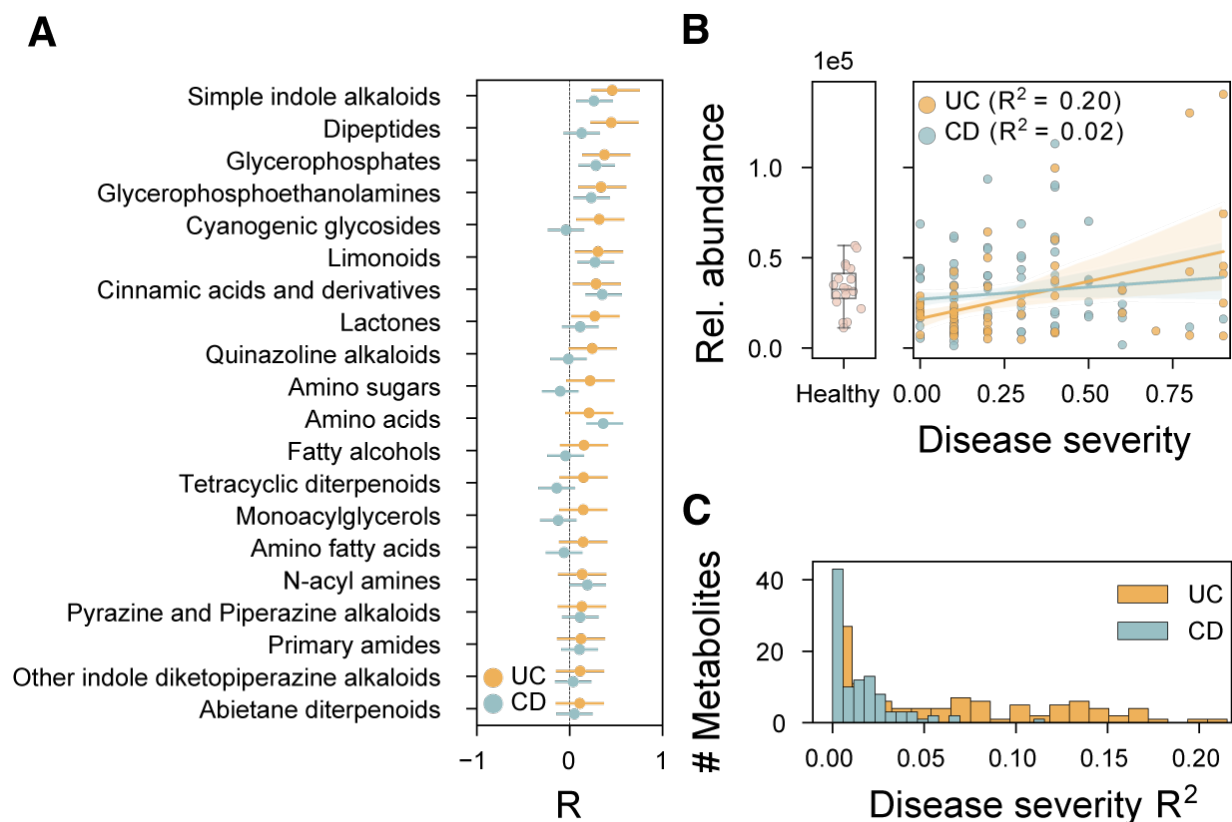


Figure 2.10: **Repeated analysis of clinically relevant dipeptides with the PubChem retrieval library.** **A.** Metabolite classes ranked by their Pearson correlation with IBD disease severity. All putative metabolites are assigned disease classes with NPClassifier [101] and relative abundances are summed across classes. Within the ulcerative colitis (UC) and Crohn’s disease (CD) cohorts, metabolite classes are correlated with disease severity and sorted according to their correlation with UC (top 20 shown; each computed for $n=60$ patients with UC, $n=102$ CD patients). Error bars indicate 95% confidence intervals around the mean. **B.** Dipeptide relative abundances for each patient are plotted against disease severity within the healthy ($n=19$), UC, and CD cohorts. Dipeptides correlate with disease severity for UC patients ($R^2 = 0.20$, $p = 0.02$, two-sided the Wald test, adjusted with Benjamini-Hochberg) as observed by Mills et al. [102]. In the left boxplot, the data median line is shown; boxes show the interquartile range, with whiskers indicating 1.5x interquartile ranges. Regression lines are shown with 95% confidence intervals computed with bootstrapping. **C.** Within the dipeptide class, the distribution of individual metabolite correlations are shown for both UC and CD cohorts. Compound annotations are made using an ensemble of 5 MIST models, contrastive distance retrieval, and PubChem reference databases of molecules. Input spectra are subsetted to be under 500 Da, have $[M+H]^+$ adducts, and utilize molecular formula as output by CSI:FingerID at the end of the compound identification step.

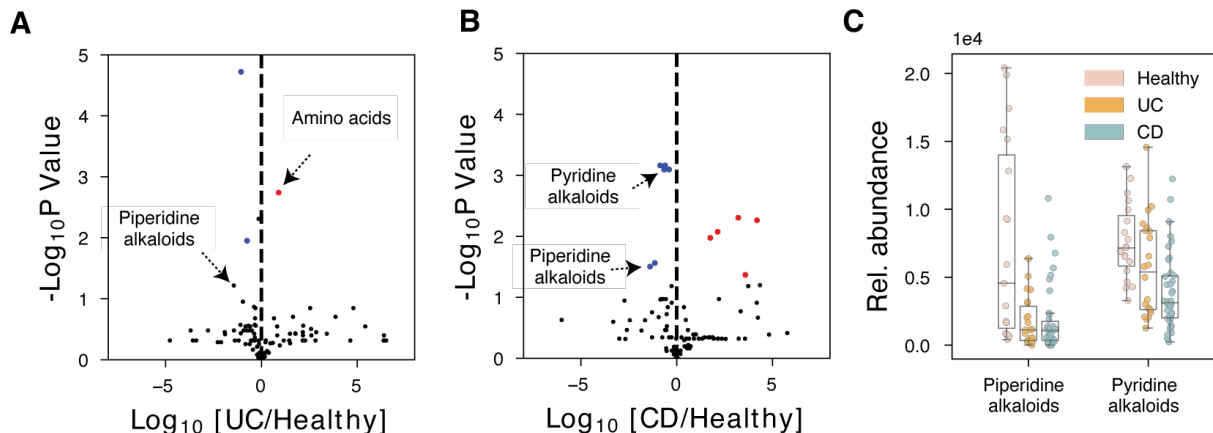


Figure 2.11: **Reanalyzing putative metabolite class differences between healthy and diseased IBD cohorts using the PubChem database.** **A,B.** Putative metabolite classes from the Mills et al. cohort [102] are scattered showing their fold change and respective p value comparing and ulcerative colitis (UC) (**A**) and Crohn’s disease (CD) (**B**) cohorts to the healthy cohort. P values were computed using independent two-sided t-tests and adjusted for multiple hypothesis testing with the Benjamini Hochberg method. To have equal cohort sizes for healthy and diseased patients, UC and CD cohorts were subsetted to patients with disease severity score > 0.2. Blue scattered points indicate differentially less abundant compounds or classes; red indicates differentially increased abundance. **C.** Distributions of patients’ relative abundances for piperidine and pyridine alkaloid classes are shown for healthy (n=19), UC (n=22), and CD (n=38) groups. Median lines are shown; boxes show the interquartile range, with whiskers indicating 1.5x interquartile ranges. All compound annotations are made using an ensemble of 5 MIST models, contrastive distance retrieval, and the PubChem reference database of molecules. The top chemical class annotation is found by running the top putative prediction through NPClassifier [101]. Input spectra are subsetted to be under 500 Da, have $[M+H]^+$ adducts, and utilize molecular formula as output by CSI:FingerID at the end of the compound identification step.

the linear and FFN output models are capable of achieving 0.6701 and 0.6617 cosine similarity vs. 0.6818 from the unfolding model. This validates the improvements offered by the unfolding layer.

2.5.5 MAGMa substructure labels

The MAGMa algorithm is used to take pairs of (molecule, spectrum) from the trained dataset and attribute molecule substructures to MS2 peaks [95]. In brief, the algorithm works by progressively removing atoms from the original structure to create substructures. Each time an atom is removed, the resulting subfragments are given a tolerance of $\pm H$ to account for hydrogen rearrangements. We highlight several spectra from the public GNPS dataset to illustrate this substructure labeling (Figure 2.13B-D).

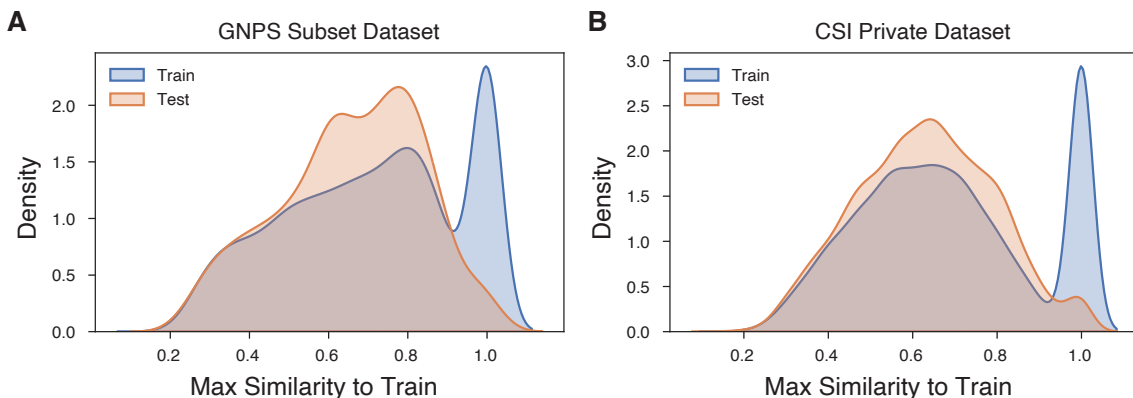


Figure 2.12: **Dataset split similarity.** Distribution plot of the most similar molecule in the training set according to Tanimoto similarity for molecules in the train and molecules in the test sets. Distributions are computed on the public GNPS CANOPUS dataset (A.) and private proprietary CSI dataset (B.).

Table 2.2: **Fingerprint prediction accuracy on CANOPUS GNPS subset dataset for MIST top model variations trained.** All models are trained without any auxiliary MAGMa loss and without any simulated spectra added to the training set. “MIST - unfolding (linear)” indicates a model without unfolding layers featuring only a single linear layer to predict fingerprint outputs. “MIST - unfolding (FFN)” indicates a MIST model trained without unfolding layers, but instead using a dense feedforward network of 3 layers. Each experiment was conducted on a single split of the data. Results are shown \pm the standard error of the mean.

| Model | Cosine similarity | Log likelihood |
|---------------------------|---------------------------------------|--|
| MIST | 0.6818 ± 0.0066 | -0.0248 ± 0.0005 |
| MIST - unfolding (linear) | 0.6701 ± 0.0061 | -0.0250 ± 0.0005 |
| MIST - unfolding (FFN) | 0.6617 ± 0.0063 | -0.0255 ± 0.0005 |

2.5.6 Alternative transformer baselines

Contemporaneously with this work, ref. [83] have proposed to utilize a Transformer architecture applied to m/z values. In their architecture, m/z values are first transformed into frequency encodings using a sinusoidal embedding, inspired by the original positional encodings of the Transformer architecture [92]. Because the code is not available at the time of this work, we develop our own in-house implementation of this approach and hyperparameter optimize model parameters on our public-domain dataset, settling upon a 256 hidden size, 4 transformer layers, 0.3 dropout, 0 weight decay, and 0.0002 learning rate.

We find that in our data regime, the multi-scale Transformer performs equivalently to the FFN baseline, as shown in Table 2.5. We hypothesize that this is due to the relatively small amount of data in our training set compared to the original authors and also the lack of explicit information about neutral loss relationships between peaks.

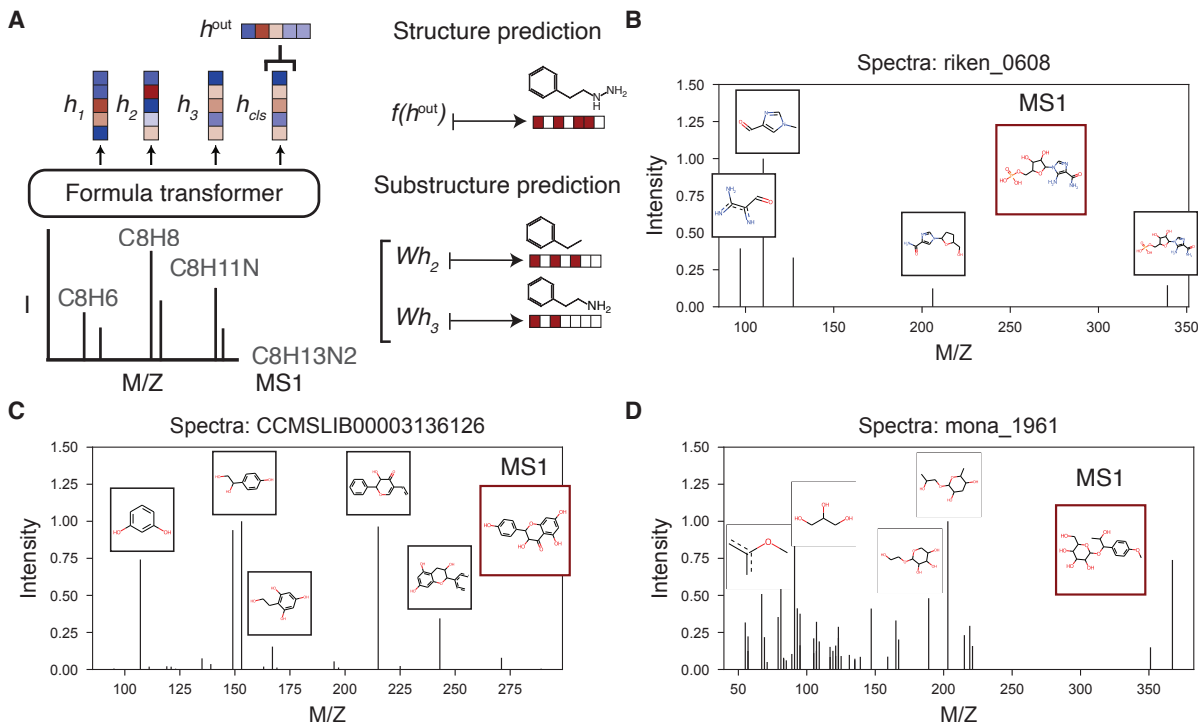


Figure 2.13: **Labeling data substructures.** **A.** MIST is trained to build representations of each peak in the spectrum. The model uses the precursor peak representation to predict a fingerprint of the full molecule, and when substructures can be labeled for other peaks, the model learns to predict these as well. **B-D.** Example substructure peak labels in the training set computed with MAGMa. The full compound is outlined in red. Other substructures for the top 6 highest intensity peaks are outlined in black.

2.5.7 Spectra noising

To increase generality as in MetFID [56], we introduce an algorithm to stochastically add noise to input training spectra (Algorithm 2.5.1). Spectra are selected for augmentation with a 50% probability. For each spectrum selected for augmentation, we define two augmentation operations: **remove**, in which a peak is fully removed, and **intensity**, in which a peak’s intensity is multiplied by a randomly sampled scalar value $s_i \sim \mathcal{N}(1, 1)$, truncated to a minimum of 0. The probability of removing each peak, p_j is proportional to its intensity value to ensure high intensity peaks are less likely to be removed:

$$p_j(\text{remove}) = 1 - p_j(\text{keep}), \quad p_j(\text{keep}) \propto \exp[I_j] \quad (2.30)$$

For each noised spectrum, we sample the number of peaks to remove from a binomial distribution with $p_{\text{remove}} = 0.5$ and we sample a second binomial distribution for the number of peaks to rescale with probability $p_{\text{intensity}} = 0.1$.

Algorithm 2.5.1: Noising a training spectrum

Data: \mathcal{S} ▷ Example spectrum
Result: Noised version of \mathcal{S}
 $peaks \leftarrow |\mathcal{S}|$ ▷ Get the number of peaks
 $I \leftarrow \text{get_intens}(\mathcal{S})$ ▷ Get all intensities
if $\text{rnd}() \geq 0.5$ **then**
 $\text{num_remove} \leftarrow \text{Binomial}(peaks, 0.5)$
 $\text{num_rescale} \leftarrow \text{Binomial}(peaks, 0.1)$
 for $j \in \mathcal{S}$ **do**
 $\text{modify_probs}_j \leftarrow \exp(I_j)$
 $\text{remove_inds} \leftarrow \text{Sample}(peaks, \text{num_remove}, \text{modify_probs})$
 $\text{rescale_inds} \leftarrow \text{Sample}(peaks, \text{num_rescale}, \text{modify_probs})$
 $\mathcal{S} \leftarrow \text{remove_peaks}(\mathcal{S}, \text{remove_inds})$
 $\mathcal{S} \leftarrow \text{rescale_peaks}(\mathcal{S}, \text{rescale_inds})$

2.5.8 Hyperparameter optimization

The tree-structured Parzen estimator hyperparameter search algorithm [122] as implemented by Optuna [123] was used to identify the best set of hyperparameters for each model. Hyperparameters were optimized once to maximize performance on a single validation split of the data, excluding auxiliary loss terms. RayTune [124] was used to enable efficient parallelization across 3 RTX 3090 GPUs and a HyperBand scheduler was used to prune unpromising hyperparameters for efficiency. We ensured the same compute resources were available for tuning our models and each of the baselines and therefore capped resources at 100 trials on 3 GPUs. We conducted two separate hyperparameter sweeps to identify model parameters best for predicting fingerprints from CSI:FingerID (NIST dataset) and also the Morgan 4096 bit fingerprints (CANOPUS benchmark dataset). Parameters are described in Table 2.11 and Table 2.12.

Table 2.3: Full performance metrics for CSI:FingerID, FFN, MIST, and a MIST ensemble of 5 models are shown averaged across 3 structural splits of the data. Log likelihoods are clamped such that they have a minimum of -5 . Results are shown \pm standard error with an arbitrary number of significant figures.

| | Tanimoto similarity (spectra) | Cosine similarity (spectra) | Log likelihood (spectra) | Log likelihood (bits) |
|--------------|---------------------------------------|---------------------------------------|--|--|
| CSI:FingerID | 0.7382 ± 0.0011 | 0.8759 ± 0.0006 | -0.0352 ± 0.0002 | -0.0352 ± 0.0008 |
| FFN | 0.5108 ± 0.0012 | 0.7418 ± 0.0009 | -0.0710 ± 0.0003 | -0.0710 ± 0.0016 |
| MIST | 0.7475 ± 0.0012 | 0.8774 ± 0.0007 | -0.0345 ± 0.0002 | -0.0345 ± 0.0007 |
| MIST (5x) | 0.7629 ± 0.0012 | 0.8892 ± 0.0006 | -0.0309 ± 0.0002 | -0.0309 ± 0.0007 |

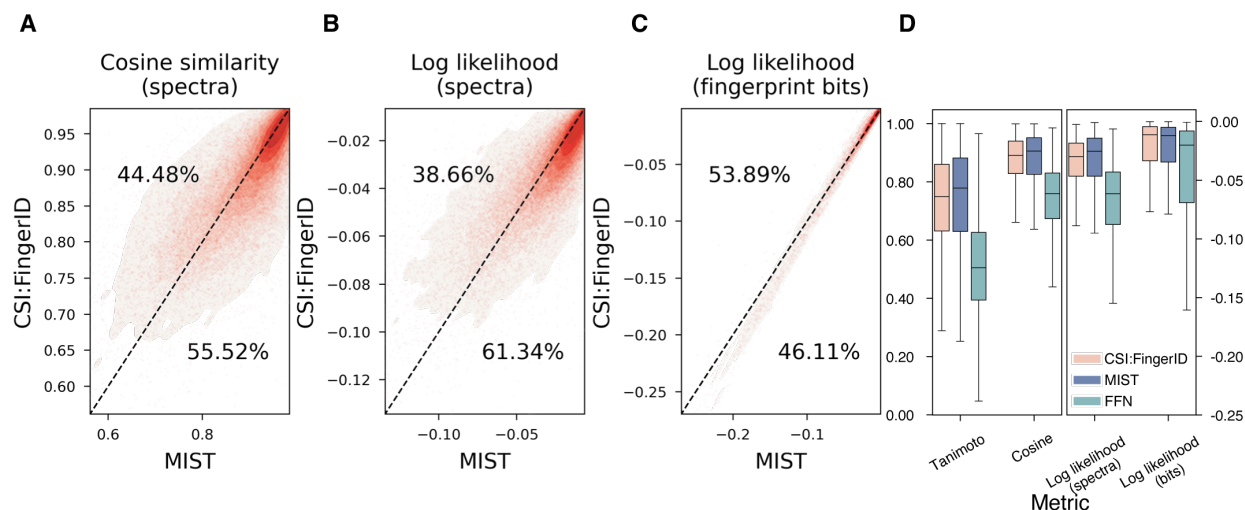


Figure 2.14: **Single MIST model comparison to CSI:FingerID.** Molecular fingerprints are predicted by MIST and CSI:FingerID for every spectrum in the test set using a single MIST model as in Figure 2. The performance for each spectrum by cosine similarity to the true fingerprint (A) or log likelihood (B) is evaluated and plotted. Points below the line represent instances where MIST is more performant. C. Equivalent evaluation showing the likelihood of predicting each fingerprint bit correctly across all spectra. D. The performance of CSI:FingerID, MIST, and FFN, a baseline inspired by MetFID, are shown. “Tanimoto,” “Cosine,” and “Log likelihood (spectra)” indicate performance across spectra and “Log likelihood (bits)” indicates performance across bits (higher is better). Median lines are shown; boxes show the interquartile range, with whiskers indicating 1.5x interquartile ranges. Statistics are pooled across $n = 18,700$ test spectra (tanimoto, cosine, and spectra log likelihood) and $n = 5,496$ fingerprint bits (bit log likelihood). All results are aggregated across 3 splits of the data.

Table 2.4: **Full retrieval accuracy performance metrics for CSI:FingerID, FFN, MIST, and a MIST ensemble of 5 models are shown averaged across 3 structural splits of the data.**

| | Top 1(%) | Top 5(%) | Top 10(%) | Top 20(%) | Top 50(%) | Top 100(%) | Top 200(%) |
|-------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| CSI:FingerID + Cosine | 38.24 | 68.88 | 77.07 | 83.20 | 89.12 | 92.23 | 94.76 |
| CSI:FingerID + Bayes | 43.00 | 74.77 | 82.15 | 87.23 | 91.70 | 94.13 | 96.00 |
| MIST + Cosine | 29.72 | 59.56 | 69.26 | 73.60 | 84.10 | 88.13 | 91.84 |
| MIST + Contrastive | 34.45 | 69.28 | 79.04 | 85.75 | 91.63 | 94.56 | 96.49 |
| MIST (5x) + Cosine | 31.70 | 62.85 | 72.33 | 79.05 | 85.92 | 89.76 | 93.00 |
| MIST (5x) + Contrastive | 37.39 | 72.90 | 82.03 | 88.20 | 93.32 | 95.75 | 97.40 |

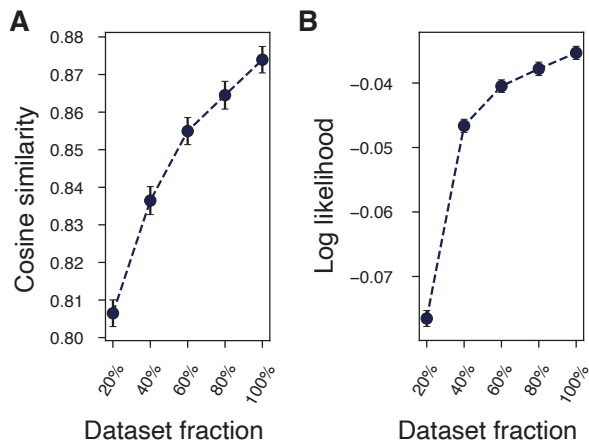


Figure 2.15: **Dataset ablations on the CSI2022 dataset validate that model performance is limited by data set size.** MIST fingerprint prediction accuracy improves as a function of dataset size in terms of both (A) cosine similarity and (B) log likelihood. Model performance is computed on the partially proprietary CSI2022 dataset. Dataset ablations are conducted for a single split of the data. Log likelihood values are clamped to a minimum of -5 . Error bars show 95% confidence intervals for the mean ($n=3,116$ test spectra). For this ablation, MIST models are trained without simulated data and without MAGMa auxiliary loss to reduce model training time.

Table 2.5: **Full fingerprint prediction accuracy on CANOPUS GNPS subset dataset for various model ablations.** All model performances are computed on a merged dataset of fingerprint predictions computed for 3 independent random structural split hold outs of the CANOPUS dataset. FFN: Feed forward neural network. Transformer: A partial reimplement of [83]. Results are shown plus or minus the standard error of the mean.

| Model | Cosine similarity | Log likelihood |
|------------------|---------------------------------------|--|
| FFN | 0.5368 ± 0.0042 | -0.0344 ± 0.0004 |
| Transformer [83] | 0.5267 ± 0.0039 | -0.0333 ± 0.0003 |
| MIST - unfolding | 0.6796 ± 0.0035 | -0.0244 ± 0.0003 |
| MIST - simulated | 0.6880 ± 0.0038 | -0.0248 ± 0.0003 |
| MIST - MAGMa | 0.6907 ± 0.0038 | -0.0228 ± 0.0003 |
| MIST - pairwise | 0.6963 ± 0.0037 | -0.0232 ± 0.0003 |
| MIST | 0.6953 ± 0.0037 | -0.0236 ± 0.0003 |

Table 2.6: **Fingerprint prediction accuracy on primary CSI2022 dataset to compare models with and without pairwise featurizations.** Model performance values are computed on a merged dataset of fingerprint predictions computed for 3 independent random structural split test sets. All MIST models described were trained without simulated data or MAGMa auxiliary loss features for quicker model runs. Log likelihoods are clamped such that they have a minimum of -5. Results are shown \pm standard error with an arbitrary number of significant figures.

| Model | Cosine similarity | Log likelihood |
|-----------------|---------------------|----------------------|
| MIST + pairwise | 0.8693 \pm 0.0007 | -0.0368 \pm 0.0002 |
| MIST - pairwise | 0.8688 \pm 0.0008 | -0.0369 \pm 0.0002 |

Table 2.7: **Fingerprint prediction accuracy on CANOPUS GNPS subset dataset for MIST models with various featurizations of peak intensities.** All models are trained without any auxiliary MAGMa loss and without any simulated spectra added to the training set. “MIST + float intensity” indicates a concatenation of a single floating point value of intensity between 0 and 1; “MIST + float intensity & pooling” indicates a concatenation of a floating point intensity and selective intensity-weighted pooling of hidden states in the transformer; “MIST + log intensity” indicates a concatenation of a single log-transformed intensity value; “MIST + cat. intensity” indicates a concatenation of a one-hot vector of intensity transformed into one of 10 evenly spaced categorical bins; “MIST + zero intensity” indicates a concatenation of a constant float of zero rather than the true intensity. Each experiment was conducted on a single split of the data and average results are shown \pm the standard error of the mean.

| Model | Cosine similarity | Log likelihood |
|----------------------------------|---------------------------------------|--|
| MIST + float intensity | 0.6818 \pm 0.0066 | -0.0248 \pm 0.0005 |
| MIST + float intensity & pooling | 0.6775 \pm 0.0063 | -0.0262 \pm 0.0005 |
| MIST + log intensity | 0.6743 \pm 0.0065 | -0.0264 \pm 0.0005 |
| MIST + cat. intensity | 0.6743 \pm 0.0066 | -0.0252 \pm 0.0005 |
| MIST + zero intensity | 0.6726 \pm 0.0063 | -0.02707 \pm 0.0005 |

Table 2.8: **Full retrieval accuracy on CANOPUS GNPS subset dataset for various model ablations.** All model performance values are computed on a merged dataset of retrieval predictions computed for 3 independent random structural splits of 3 random structural splits of the CANOPUS dataset. FFN: Feed forward neural network.

| | Top 1(%) | Top 5(%) | Top 10(%) | Top 20(%) | Top 50(%) | Top 100(%) | Top 200(%) |
|--------------------------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
| FFN fingerprint | 17.309 | 37.121 | 45.611 | 54.634 | 63.946 | 71.329 | 77.359 |
| FFN contrastive | 20.632 | 44.996 | 54.389 | 63.536 | 75.062 | 80.558 | 86.095 |
| MIST contrastive - pretrain | 24.774 | 50.656 | 60.624 | 68.827 | 78.384 | 83.593 | 88.146 |
| MIST contrastive | 28.384 | 55.373 | 65.217 | 72.970 | 81.255 | 85.480 | 89.377 |
| MIST fingerprint | 29.368 | 55.332 | 63.536 | 72.231 | 80.476 | 85.726 | 89.418 |
| MIST contrastive + fingerprint | 30.703 | 58.120 | 68.927 | 75.709 | 84.094 | 87.916 | 92.355 |

Table 2.9: **Average molecular similarity between spectra pairs with the 0.1% highest similarity according to the specified distance function.**

| Method | Average Tanimoto Similarity |
|-----------------------|-----------------------------|
| Upper bound | 0.451 |
| MIST | 0.324 |
| MS2DeepScore [82] | 0.219 |
| Spec2Vec [81] | 0.186 |
| Modified Cosine [100] | 0.133 |
| Random | 0.076 |

Table 2.10: **Forward model ablations demonstrate the utility of unfolding and reverse mass predictions.** Three variants of the forward simulation model are trained on a single split of the public GNPS dataset: “FFN” (full model), “FFN - unfolding” (full model a single linear layer to map to the 15,000 dimension output), and “FFN - reverse” (full model without predicting mass differences as in [84]). The average cosine similarity and coverage (i.e., fraction of overlapping peaks between the top 100 predicted peaks and ground truth spectra) between the binned spectra and true spectra are shown. Values are plotted \pm the standard error of the mean.

| Model | Cosine similarity | Coverage |
|-----------------|---------------------|---------------------|
| FFN | 0.3839 ± 0.0056 | 0.5815 ± 0.0086 |
| FFN - unfolding | 0.3502 ± 0.0058 | 0.5555 ± 0.0088 |
| FFN - reverse | 0.2583 ± 0.0058 | 0.5160 ± 0.0088 |

Table 2.11: Model hyperparameters searched and set for CSI fingerprint prediction.

| Model | Parameter | Grid | Value |
|--|--|-------------------------|---------|
| Forward | learning rate | $[1e-4, 1e-3]$ | 0.00086 |
| | scheduler | [True, False] | False |
| | learning rate decay | [0.7, 0.999] | - |
| | dropout | {0.1, 0.2, 0.3} | 0.2 |
| | hidden size, d | {64, 128, 256, 512} | 512 |
| | layers, l | {1, 2, 3} | 2 |
| | unfolding | [True, False] | True |
| | unfolding layers, γ | [1, 2, 3, 4] | 3 |
| | unfolding weight, λ_{fwd} | $[1e-4, 1]$ | 0.003 |
| | batch size | - | 64 |
| FFN | learning rate | $[1e-5, 1e-3]$ | 0.00087 |
| | scheduler | [True, False] | False |
| | learning rate decay | [0.7, 0.999] | - |
| | dropout | {0.1, 0.2, 0.3} | 0.0 |
| | hidden size, d | {64, 128, 256, 512} | 512 |
| | layers, l | {1, 2, 3} | 2 |
| | spectra bins | [1000, 2000, ... 15000] | 11000 |
| | batch size | - | 64 |
| MIST | learning rate | $[1e-5, 1e-3]$ | 0.00026 |
| | scheduler | [True, False] | True |
| | weight decay | $\{1e-6, 1e-7, 0\}$ | $1e-7$ |
| | learning rate decay | [0.7, 0.999] | 0.97 |
| | dropout | {0.1, 0.2, 0.3} | 0.1 |
| | hidden size, d | {64, 128, 256, 512} | 256 |
| | layers, l | {1, 2, 3, 4, 5} | 3 |
| | fraction real data f_{real} | {0.1, 0.2, ..., 1.0} | 0.5 |
| | unfolding layers, ν | {1, 2, 3, 4, 5} | 5 |
| | unfolding loss weight, λ_{unfold} | {0.1, 0.2, ..., 1.0} | 0.6 |
| | magma loss weight, λ_{magma} | {1, 2, ..., 15} | 2 |
| | probability noised spectrum, p_{noise} | - | 0.5 |
| | probability remove peak, p_{remove} | - | 0.5 |
| probability rescale peak, $p_{\text{intensity}}$ | - | 0.1 | |
| batch size | - | 128 | |
| Contrastive | learning rate | $[1e-5, 1e-3]$ | $6e-5$ |
| | scheduler | [True, False] | True |
| | weight decay | $\{1e-6, 1e-7, 0\}$ | 0 |
| | learning rate decay | [0.7, 0.999] | 0.801 |
| | contrastive weight, λ_{c} | {0.1, 0.2, ..., 1.0} | 0.6 |
| | fraction real data f_{real} | {0.1, 0.2, ..., 1.0} | 0.5 |
| | probability noised spectrum, p_{noise} | - | 0.5 |
| | probability remove peak, p_{remove} | - | 0.5 |
| | probability rescale peak, $p_{\text{intensity}}$ | - | 0.1 |
| | batch size | - | 32 |

Table 2.12: Model hyperparameters searched and set for Morgan fingerprint prediction.

| Model | Parameter | Grid | Value |
|--|--|-------------------------|---------|
| Forward | learning rate | $[1e-4, 1e-3]$ | 0.00054 |
| | scheduler | [True, False] | False |
| | learning rate decay | [0.7, 0.999] | - |
| | dropout | {0.1, 0.2, 0.3} | 0.3 |
| | hidden size, d | {64, 128, 256, 512} | 512 |
| | layers, l | {1, 2, 3} | 2 |
| | unfolding | [True, False] | True |
| | unfolding layers, γ | [1, 2, 3, 4] | 3 |
| | unfolding weight, λ_{fwd} | $[1e-4, 1]$ | 0.0016 |
| | batch size | - | 64 |
| FFN | learning rate | $[1e-5, 1e-3]$ | 0.00087 |
| | scheduler | [True, False] | False |
| | learning rate decay | [0.7, 0.999] | - |
| | dropout | {0.1, 0.2, 0.3} | 0.3 |
| | hidden size, d | {64, 128, 256, 512} | 512 |
| | layers, l | {1, 2, 3} | 2 |
| | spectra bins | [1000, 2000, ... 15000] | 11000 |
| | batch size | - | 64 |
| MIST | learning rate | $[1e-5, 1e-3]$ | 0.00077 |
| | scheduler | [True, False] | False |
| | weight decay | $\{1e-6, 1e-7, 0\}$ | $1e-7$ |
| | learning rate decay | [0.7, 0.999] | - |
| | dropout | {0.1, 0.2, 0.3} | 0.1 |
| | hidden size, d | {64, 128, 256, 512} | 256 |
| | layers, l | {1, 2, 3, 4, 5} | 2 |
| | fraction real data f_{real} | {0.1, 0.2, ..., 1.0} | 0.6 |
| | unfolding layers, ν | {1, 2, 3, 4, 5} | 4 |
| | unfolding loss weight, λ_{unfold} | {0.1, 0.2, ..., 1.0} | 0.4 |
| | magma loss weight, λ_{magma} | {1, 2, ..., 15} | 8 |
| | probability noised spectrum, p_{noise} | - | 0.5 |
| | probability remove peak, p_{remove} | - | 0.5 |
| probability rescale peak, $p_{\text{intensity}}$ | - | 0.1 | |
| batch size | - | 128 | |
| Contrastive | learning rate | $[1e-5, 1e-3]$ | 0.00057 |
| | scheduler | [True, False] | True |
| | weight decay | $\{1e-6, 1e-7, 0\}$ | $1e-7$ |
| | learning rate decay (10k steps) | [0.7, 0.999] | 0.74 |
| | contrastive weight, λ_{c} | {0.1, 0.2, ..., 1.0} | 0.6 |
| | fraction real data f_{real} | {0.1, 0.2, ..., 1.0} | 0.2 |
| | probability noised spectrum, p_{noise} | - | 0.5 |
| | probability remove peak, p_{remove} | - | 0.5 |
| | probability rescale peak, $p_{\text{intensity}}$ | - | 0.1 |
| | batch size | - | 32 |

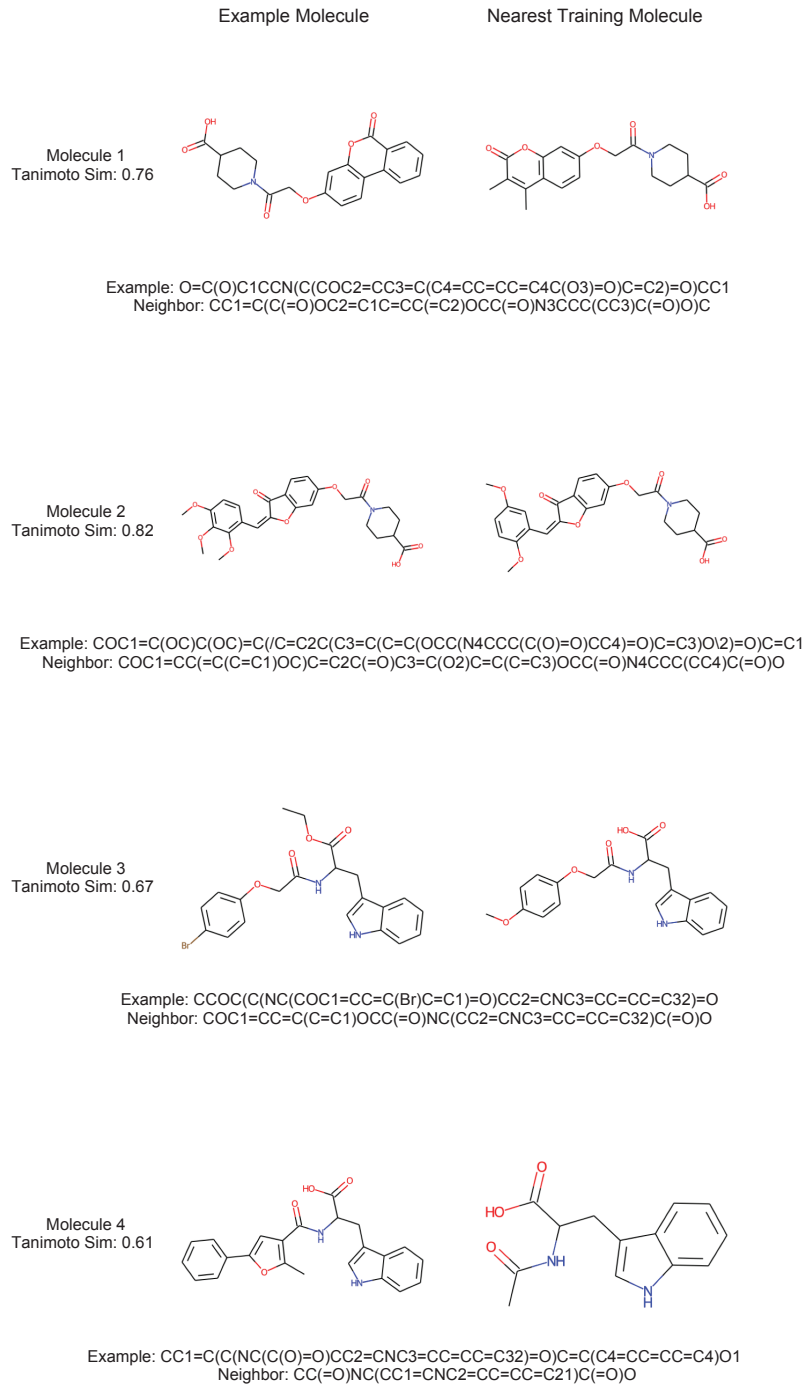


Figure 2.16: **Example molecules and their training precedents for Figure 2.3.** Each molecule is shown along side its nearest neighbor in the training set, with its Tanimoto similarity listed.

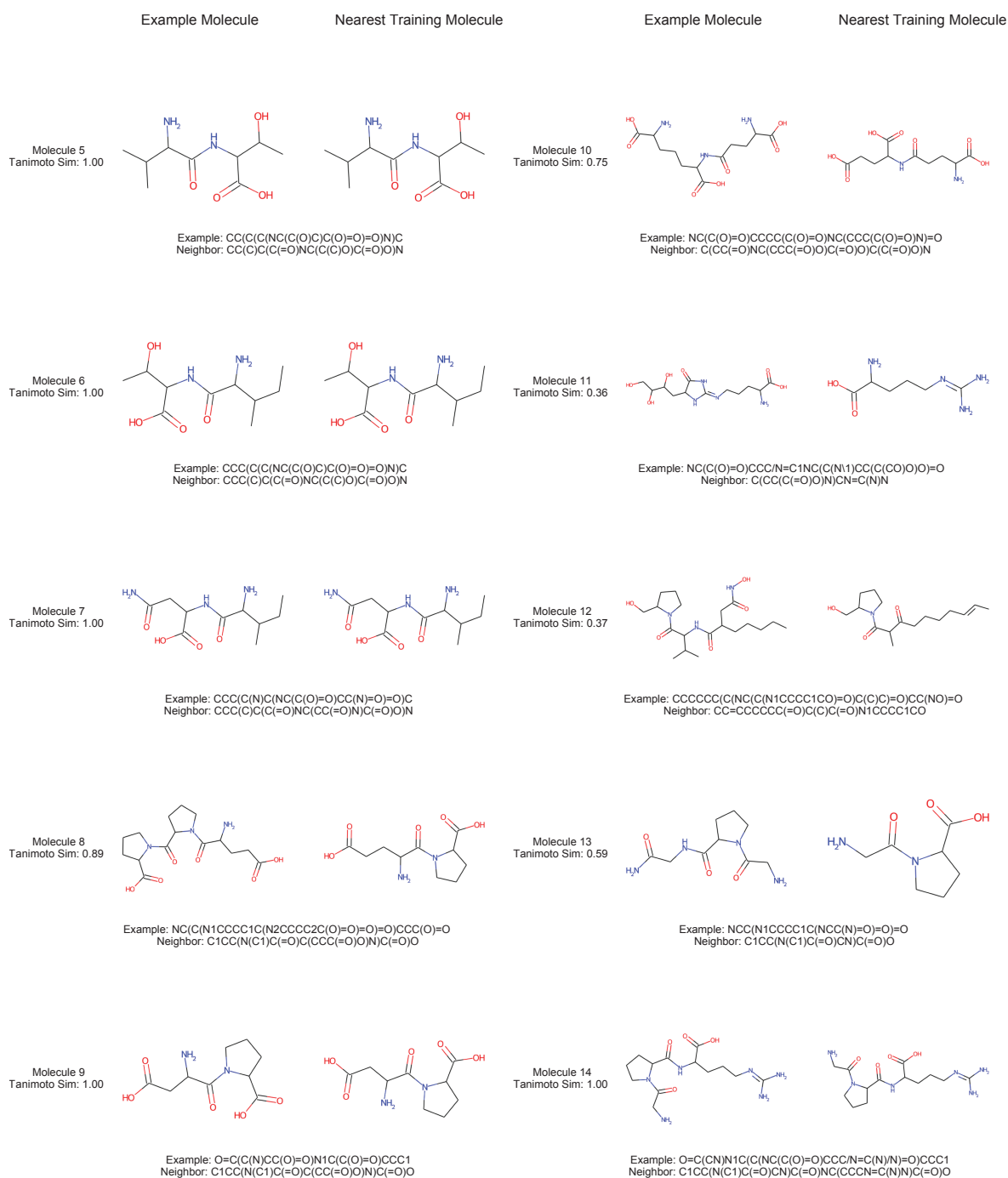


Figure 2.17: **Example molecules and their training precedents for Figure 2.4.** Each molecule is shown along side its nearest neighbor in the training set, with its Tanimoto similarity listed. Even in cases for which the proposed molecular structure appears in the training set, it is predicted on the basis of a new experimental spectrum that does not exactly match the reference spectrum.

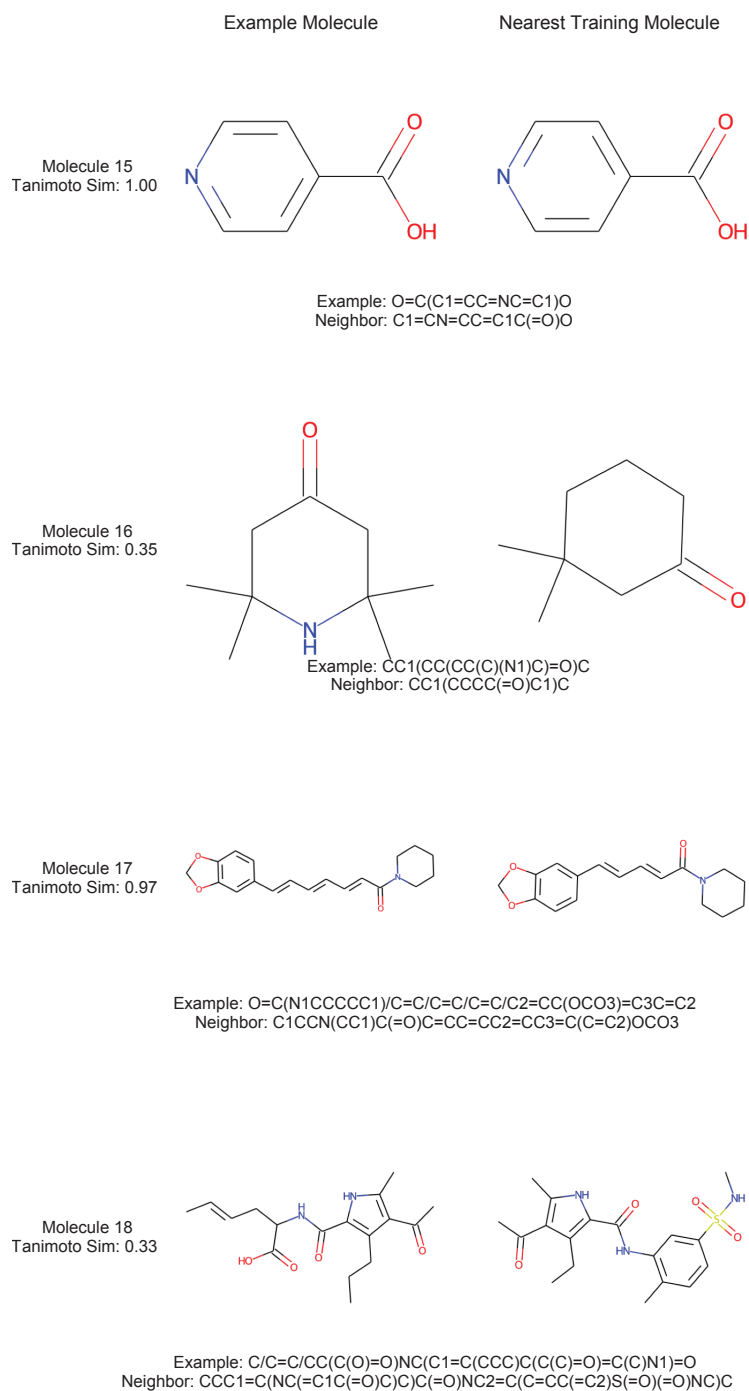


Figure 2.18: **Example molecules and their training precedents for Figure 2.5.** Each molecule is shown along side its nearest neighbor in the training set, with its Tanimoto similarity listed. Even in cases for which the proposed molecular structure appears in the training set, it is predicted on the basis of a new experimental spectrum that does not exactly match the reference spectrum.

Chapter 3

Inferring the Precursor Mass Molecular Formula from Spectra

This work has previously appeared as S. Goldman, J. Xin, J. Provenzano, *et al.*, “MIST-CF: Chemical Formula Inference from Tandem Mass Spectra,” *Journal of Chemical Information and Modeling*, 2023. I jointly led the conceptualization, execution, and writing for this work. I contributed to all of these tasks equally alongside the second author, Jiayi Xin, an undergraduate student whom I supervised and mentored through the duration of this project. Other authors contributed to this work and all processes in various ways.

Through peer review and initial publication of this work, we referred to the precursor formula defining the elemental composition of a molecule as a “chemical formula” rather than a “molecular formula.” It has since come to our attention that “chemical formula” is ambiguous and can also apply to skeletal formula representations of molecules, rather than the elemental composition as we intended. In this thesis, I have adjusted “chemical formula” to “molecular formula” in all such instances for disambiguation, with the one exception being the title of the method MIST-CF.

3.1 Introduction

The discovery of previously unknown small molecules in biological samples is rapidly expanding our knowledge of plant chemistry [66], [126], cancer biology [8], [15], host-microbiome interactions [13], [14], [127], and other metabolite-mediated human biology [67]. Similar small molecule discoveries in the environmental sciences have led to new insights regarding the exposome and pollutant effects, resolving mysteries such as high salmon mortality rates. [18], [19] Increasing our ability to detect and identify the so-called “dark metabolome” with analytical chemistry techniques represents an exciting opportunity in experimental and computational chemistry.

Tandem mass spectrometry (MS/MS) is a particularly well-suited analytical technique for this, as it allows for the high throughput characterization of small molecules from complex mixtures. [68] In an MS/MS experiment, both an intact ionized mass (MS1) and a set of fragment peak masses (MS2) can be measured for each unknown molecule. This fragmentation spectrum serves as a structural representation of the molecule, ideally allowing

the practitioner to match a small molecule structure to each resulting spectrum based upon database spectra matches. Due to the vastness of chemical space, however, many observed spectra have no precedent; in a large public MS/MS database, 87% of observed spectra remain unannotated [29].

In such instances without spectral matches, we must rely on informatics and predictive modeling to identify the likely molecular structure. This pipeline almost always begins by inferring a molecular formula from the observed spectrum (Figure 3.1A). There are many formula options for each observed MS1 value, especially for higher mass compounds. Specifying the molecular formula (e.g., $C_6H_{12}O_6$, $C_9H_{11}NO_3$, etc.) constrains the space of potential compound candidates to a set of isomers, whereas the MS1 alone only constrains the space of candidates to those with similar masses. Automated assignment is far from trivial; while the the maximum formula annotation accuracy was 94% in in the recent Critical Assessment of Metabolite Identification in 2022 (CASMI20220), [75] the median score was 71%, *and* these percentages were calculated only for the submitted predictions, rather than the total number of spectra tested. Improving the automated prediction accuracy of this step promises to improve, simplify, and speed up downstream analyses.

Molecular formula annotation tools can be grouped into two categories: database dependent and database independent. Database dependent searches place restrictions on potential formulae, querying the candidate mass and spectrum against databases including NIST [63], GNPS [28], HMDB [26], or large compound libraries such as PubChem [128]. Relying on databases inherently limits annotations to formulae that have already been observed. On the other hand, database independent (*de novo*) molecular formula annotation considers all possible molecular formulae, though the task becomes more challenging due to the larger number of candidates. A recent bottom-up computational approach, BUDDY [52], is a hybrid of the two approaches and assigns database formulae to MS2 peaks and neutral losses. By combining peak and neutral loss annotations, BUDDY generates potential candidates not present in databases. However, overall annotation performance is a function both of how the candidate space is defined and an algorithm’s ability to rank them; this limits our ability to perform a direct comparison to or analysis of BUDDY. Herein we focus specifically on *de novo* formula annotation for maximal flexibility using only MS/MS information.

Notably, both MZmine [120], [129] and SIRIUS [50], [54], [77] have developed widely used methods for scoring molecular formula candidates using MS/MS information for *de novo* annotation. While MZmine evaluates each candidate formula based on the number of MS/MS peaks it can explain, [129] SIRIUS scores them through a more expressive fragmentation tree strategy [50], [54], [77]. SIRIUS first proposes candidate formulae through an exhaustive enumeration step up to a certain mass error from the observed MS1, labels the MS2 peaks with potential “subformulae” of the candidate MS1 annotation, performs an optimization to arrange these subformula annotations into a fragmentation tree, and finally computes a likelihood of the molecular formula based upon the tree and isotopic patterns. Despite their success and widespread use, these methods offer room for improvement in terms of both accuracy and speed. We recently observed program timeouts using SIRIUS for larger molecules (i.e., over 800 Da) during fragmentation tree calculations [64]. An additional, lesser appreciated component of SIRIUS is that the method re-ranks candidate molecular formula based upon compound scores in the structure annotation step; this phase *implicitly* reposes SIRIUS as a database dependent search and led to changes in formula annotations

in our previous study, particularly with respect to adduct assignments [64]. It is unclear which pipeline steps lead to the high empirical formula annotation rates and the extent to which the tree score can be improved.

In this work, we present MIST-CF: Metabolite Inference with Spectrum Transformers for Chemical Formula prediction, an energy-based modeling approach to improving the database-independent, *de novo* molecular formula assignment step conditioned on both the MS1 mass and MS/MS spectrum. We previously demonstrated that Formula Transformers can be used to replace both the SVM module and fragmentation tree kernels used by CS:FingerID during the annotation step. [64] This study now demonstrates that fragmentation trees can be replaced throughout the MS/MS processing pipeline for equally accurate, fast, and robust predictions. As part of this, we develop a simple peak subformula assignment routine, thus circumventing the need for fragmentation tree construction. We rigorously evaluate the performance of MIST-CF on two datasets: a public dataset subset from the GNPS [28] we term NPLIB1 [57], [130] and a variant of NPLIB1 including spectra from the commercial NIST20 dataset [63]. By training and evaluating our model both with and without data from NIST20, we enable reproducible evaluation of select trained models even in the absence of a NIST20 license. MIST-CF achieves equivalent formula annotation accuracy on the positive mode CASMI2022 challenge spectra to the winning SIRIUS solution and outperforms the out-of-the-box SIRIUS assignments by a margin of 19% (i.e., 86.8% vs. 67.8% accuracy). Altogether, this demonstrates a path forward for replacing fragmentation trees in MS/MS data processing with a fully integrated deep learning processing structure annotation pipeline.

We release MIST-CF as an open source tool that can be easily integrated into existing pipelines, with or without retraining, and is freely available under the MIT license at ref. Zenodo 8151513 and <https://github.com/samgoldman97/mist-cf>.

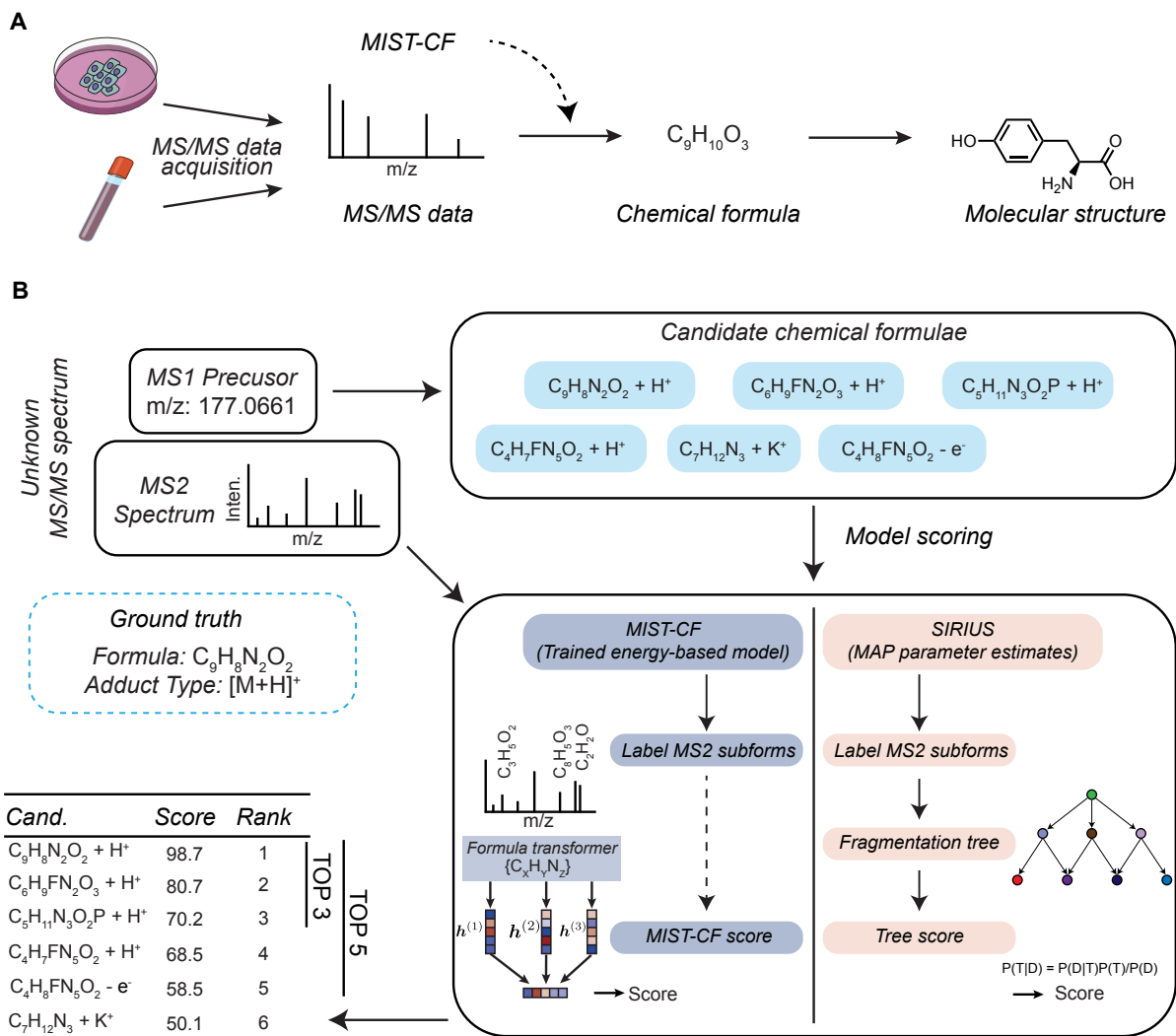


Figure 3.1: **MIST-CF** and **SIRIUS** both address the MS/MS molecular formula annotation problem. **A**. Input samples are first processed, recording a tandem mass spectrum. Before assigning full molecular structures, the molecular formula can first be inferred to constrain structure assignment. **B**. Methodological similarities and differences between MIST-CF and SIRIUS. A candidate precursor mass is first decomposed into plausible molecular formulae and adduct pairs. MIST-CF (left) learns in a data-driven fashion to assign scores and circumvents the need of fragmentation tree construction compared to SIRIUS (right). Both methods rely on assigning subformulae (“subforms”), which are molecular formula subsets of the candidate precursor formula that match individual MS/MS peak masses.

3.2 Methods

3.2.1 Preliminaries

In an MS/MS experiment, an input small molecule \mathcal{M} is ionized, often by the addition of an adduct (denoted by one-hot vector, \mathcal{A}), measured, and fragmented in order to produce an MS/MS spectrum \mathcal{S} . The spectrum is composed of peaks that can each be represented as mass, intensity pairs:

$$\mathcal{S} := \{(M^{(1)}, I^{(1)}), (M^{(2)}, I^{(2)}), \dots, (M^{(|\mathcal{S}|)}, I^{(|\mathcal{S}|)})\} \quad (3.1)$$

Herein we consider only positively charged ions and peaks and assume that each fragment carries only a single positive charge; m/z and mass are used interchangeably when describing observed peaks.

In addition to the observed MS/MS spectrum \mathcal{S} , we also have access to the MS1 measurements denoted as the precursor mass M . The core challenge of formula annotation is to determine the molecular formula of \mathcal{M} given M and \mathcal{S} . We denote the target molecular formula as a vector of integers, \mathcal{F} .

We represent the target molecular formula as a vector \mathcal{F} , in which each index $(\mathcal{F})_j \in \mathbb{Z}^+$ corresponds to the integer count for the observed j^{th} element, where the parenthetical denotes an index into the vector. In total, we consider a set of common elements, "C", "N", "P", "O", "S", "Si", "I", "H", "Cl", "F", "Br", "B", "Se", "Fe", "Co", "As", "Na", and "K" for formula vectors of size 18. We consider common positive mode adducts when generating molecular formula candidates: $[M+H]^+$, $[M+Na]^+$, $[M-H_2O+H]^+$, $[M+NH_4]^+$, $[M]^+$, and $[M-2H_2O+H]^+$. The adduct candidate is specified by a one-hot vector \mathcal{A} . For *de novo* molecular formula annotation, M can be decomposed into potential formula and adduct candidate options exhaustively using a highly efficient dynamic programming algorithm [54] within small 1 to 5 parts-per-million (ppm) mass tolerances. The generated (formula, adduct) candidates $\{(\mathcal{F}_{(1)}), \mathcal{A}_{(1)}), (\mathcal{F}_{(2)}), \mathcal{A}_{(2)}), \dots, (\mathcal{F}_{(k)}), \mathcal{A}_{(k)})\}$ can be further filtered based on presence in a database or using various heuristics such as number of ring double bond equivalents if desired [131].

In addition to inferring the molecular formula for the full molecule, a step common to SIRIUS fragmentation tree generation [77] and CSI:FingerID [53] is to derive corresponding subformula annotations for the MS/MS peaks. In this way, a spectrum can be considered as a set of (subformulae, intensity) pairings, in which the set of peak masses $\{M^{(1)}, M^{(2)}, \dots, M^{(|\mathcal{S}|)}\}$ is replaced with a set of molecular formula vectors $\{\mathbf{f}^{(1)}, \mathbf{f}^{(2)}, \dots, \mathbf{f}^{(N_p)}\}$, where each molecular formula is a subset of the precursor formula, i.e., $\mathbf{f}^{(j)} \subseteq \mathcal{F}$; peaks for which no formula can be assigned are excluded from the list. Because the fragment peaks are also charged, the observed masses are a summation of the mass of the subformula *and* its adduct. We define the MS2 peak adduct for the i^{th} peak $\mathbf{a}^{(i)}$, the mass of which must be subtracted when considering molecular formula assignment. The parts-per-million difference at each peak between the adduct-adjusted peak mass and assigned subformula is referred to as $\epsilon^{(i)}$ to indicate measurement error.

3.2.2 An energy-based model for formula annotation

As exhaustive molecular formula candidate generation can be solved via dynamic programming, [54] the key challenge in *de novo* formula annotation is to *score* how well each candidate formula matches the observed MS1 and MS2 spectra. Taking a probabilistic lens and following previous work in metabolomics and proteomics, [54], [132] our goal is to learn the probability of a candidate molecular formula and adduct, $p_{\theta}(\mathcal{F}_{(i)}, \mathcal{A}_{(i)} | \mathcal{S}, M)$. We assume the MS1 is useful only for candidate generation, and so the problem is simplified to approximating $p_{\theta}(\mathcal{F}_{(i)}, \mathcal{A}_{(i)} | \mathcal{S})$.

Energy based models (EBM) are a probabilistic modeling framework drawing inspiration from physics in which a probability distribution is defined by an energy function, E_{θ} . [133] Mathematically, EBMs take the form:

$$p_{\theta}(x) = \frac{e^{-E_{\theta}(x)}}{Z(\theta)}, \quad (3.2)$$

where the denominator $Z(\theta) = \int_x e^{-E_{\theta}(x)} dx$ is referred to as a partition function and serves to normalize the energy to a valid probability, but is typically intractable to evaluate exactly. EBMs have reemerged in recent years, with work across the chemical sciences for reaction prediction, [134] retrosynthesis, [135] and scoring protein side-chain positions. [136] These models are naturally suited for ranking applications. In our case, we factorize the probability of the candidate formula conditioned on the spectrum as an EBM of the form:

$$p_{\theta}(\mathcal{F}_{(i)}, \mathcal{A}_{(i)} | \mathcal{S}) = \frac{e^{G_{\theta}(\mathcal{F}_{(i)}, \mathcal{A}_{(i)}, \mathcal{S})}}{Z(\theta)}, \quad (3.3)$$

where $G_{\theta}(\mathcal{F}, \mathcal{A}, \mathcal{S})$ defines an arbitrary neural network energy function that takes as input the candidate formula, adduct, and full fragmentation spectrum. For any differentiable G_{θ} , the energy function can be learned via a softmax loss function, aggregated over minibatches, and minimized via stochastic mini batch gradient descent:

$$\mathcal{L}(\mathcal{F}_{(\text{True})}, \mathcal{A}_{(\text{True})}, \mathcal{S}) = -\log \frac{e^{G_{\theta}(\mathcal{F}_{(\text{True})}, \mathcal{A}_{(\text{True})}, \mathcal{S})}}{e^{G_{\theta}(\mathcal{F}_{(\text{True})}, \mathcal{A}_{(\text{True})}, \mathcal{S})} + \sum_{k=1}^K e^{G_{\theta}(\mathcal{F}_{(k)}, \mathcal{A}_{(k)}, \mathcal{S})}}, \quad (3.4)$$

where $\{(\mathcal{F}_{(k)}, \mathcal{A}_{(k)})\}_{k=1}^K$ defines a set of ‘‘decoy’’ formulae for which G_{θ} will learn to assign a low score. In practice, we sample these decoy formulae during data preprocessing from the space of formulae with equivalent masses within a 10 ppm tolerance. We further filter candidate decoys to a maximum of 256 decoys per spectrum using the FastFilter model described below in order to sample ‘‘harder’’ decoys. The trained model can be applied independently to each candidate formula at inference time to yield a ranked list of assignments.

Our approach is conceptually similar to SIRIUS [50], [54], [77], but differs in that SIRIUS uses a heuristic maximum a posteriori (MAP) estimator to score fragmentation trees. This requires manually setting parameters such as the frequency of observing various fragments in the data. By using a flexible energy function trained via supervised learning, MIST-CF does not require manual parametrizations nor necessitate fragmentation tree generation. SIRIUS also allows for the incorporation of isotopic information from the MS1 to help identify molecular formula candidates. We focus solely on the MS2-related score rather than isotopic

MS1 information, as this can be subsequently added and is often excluded from entries in spectral database such as the GNPS [28].

3.2.3 The MIST-CF architecture

MIST-CF parametrizes the energy function using a Molecular Formula Transformer [64] we denote as $G_{\theta}^{(\text{MIST})}(\mathcal{F}, \mathcal{A}, \mathcal{S})$. For the input $(\mathcal{F}, \mathcal{A}, \mathcal{S})$, we first attempt to label each spectrum peak $(M^{(i)}, I^{(i)})$ with a plausible chemical subformula $\mathbf{f}^{(i)}$ such that the spectrum can be represented by a set of subformulae (i.e., $\mathbf{f}^{(1)}, \mathbf{f}^{(2)}, \dots, \mathbf{f}^{(3)}$). Peaks are sorted by intensity and a maximum number, N_p are retained, set to 20 by default. MIST-CF then subsequently encodes these subformula peaks with a Molecular Formula Transformer, a Set Transformer variant [92], [137] we previously defined as part of the MIST architecture. [64] The full MIST-CF pipeline is illustrated in Figure 3.1B, compared side-by-side to the SIRIUS model that requires a fragmentation tree calculation. We review the exact structure of this model below by describing each modeling step.

Subformula annotation The first step in scoring a candidate formula for a spectrum is to assign a subformula to as many MS2 peaks as possible. We begin by subtracting the precursor adduct mass from all mass values in the MS2. This approach assumes that the subpeaks have the same adduct ionization as the precursor ion, but could be modified in future iterations of MIST-CF. All possible subformulae are enumerated and filtered with a ring double bond equivalent heuristic [131] to remove implausible formulae and generate a candidate set $\{\mathbf{f} | \mathbf{f} \subseteq \mathcal{F}, \text{RDBE}(\mathbf{f}) \geq 0\}$. The mass for each subformula is compared to all adduct-adjusted spectrum masses and the peak masses are assigned their nearest subformula match within a variable ppm tolerance for each instrument type (Unknown/Ion Trap: 15 ppm; Q-ToF: 10 ppm; Orbitrap/FTICR: 5 ppm). An important contribution of this work with respect to the original MIST model [64] is that this subformula assignment is now achieved with a compact and open source NumPy [138] module, reducing the reliance on SIRIUS for high quality MS2 subformula annotations.

Formula embeddings The subformula-annotated peaks are treated as integer vector inputs. Rather than pass these into a neural network directly, we first featurize each integer vector count into a sinusoidal embedding vector [130], [139] and concatenate the resulting output. Briefly, each integer element count, v , in the input formula is encoded by our counts-based encoder into the vector:

$$\text{Abs} \left(\left[\sin \left(\frac{2\pi v}{T_1} \right), \sin \left(\frac{2\pi v}{T_2} \right), \sin \left(\frac{2\pi v}{T_3} \right), \dots \right] \right), \quad (3.5)$$

where the periods (T_1, T_2 , etc.) are set at increasing powers of two up to 256 to discriminate all possible element counts given in the input, and $\text{Abs}(\cdot)$ is the absolute value function that results in only non-negative embeddings. While some information may be lost with this activation, empirically, using only positive embeddings was found to be helpful in prior work and has been maintained for consistency [130]. Integer encodings of dimension 8 ($\log_2(256)$) for all 18 considered elements within a formula are flattened and concatenated, leading to a full formula encoding $\text{Enc}(\mathcal{F}) \in \mathbb{R}^{144+}$.

Spectra context In addition to the encodings, the instrument type is considered as a one-hot vector covariate distinguishing between Ion Trap, Q-ToF, Orbitrap, and FTICR instruments. We denote this instrument type one-hot as \mathbf{c} .

Transformer architecture Each peak subformula is first encoded into a hidden state vector that is then passed into the Transformer, with the precursor formula included as a special formula annotation $\mathbf{f}^{(0)} := \mathcal{F}$. In addition to encoding the formula, we add additional features for the formula *difference* between the MS1 formula and each MS2 formula candidate vector (also converted into a sinusoidal embedding), a floating point scalar value for the MS2 peak intensity (set to 1 for the MS1 precursor formula candidate), the observed mass error between the observed adduct-adjusted mass and the monoisotopic mass of the MS2 subformula candidate (set to 0 for the MS1 precursor formula candidate), a one-hot encoding for the adduct type (assumed to be the same within a single spectrum), the instrument type one-hot vector, a Boolean flag set to 1 only at the precursor MS1 formula, and the number of total annotated peaks with subformula. These representations are embedded into hidden dimension d with a shallow single hidden layer MLP:

$$\tilde{\mathbf{h}}^{(i)} = [\text{Enc}(\mathbf{f}^{(i)}), \text{Enc}(\mathcal{F} - \mathbf{f}^{(i)}), I^{(i)}, \epsilon^{(i)} \mathbb{I}(\mathbf{f}^{(i)} \neq \mathcal{F}), \mathbf{a}^{(i)}, \mathbf{c}, \mathbb{I}(\mathbf{f}^{(i)} = \mathcal{F}), N_p] \quad (3.6)$$

$$\mathbf{h}^{(i)} = \text{MLP}(\tilde{\mathbf{h}}^{(i)}) \quad (3.7)$$

A standard multi-head Transformer, with a slight modification to include featurized attention between peaks as described previously in the MIST architecture [64], is then used to transform these peak representations into a score:

$$G_{\theta}^{(\text{MIST})}(\mathcal{F}, \mathcal{A}, \mathcal{S}) = \text{MLP}(\text{Pool}(\text{Transformer}([\mathbf{h}^{(0)}, \mathbf{h}^{(1)}, \dots, \mathbf{h}^{(N_p)}])), \quad (3.8)$$

where $\text{Pool}(\cdot)$ denotes the conversion of the variable length output of the transformer into a fixed length vector by selecting the output representation at only the special $\mathbf{h}^{(0)}$ position. Due to certain training set examples not including MS1 masses, we take care to avoid inputting the relative mass difference between the assigned precursor formula and MS1 mass which is artificially set to 0. This helps avoid additional machine bias, as many of the training spectra have theoretical MS1 mass, rather than a measured MS1.

Model training All models are trained using the aforementioned loss calculated using decoy formula. We sample minibatches of b spectra, where there are up to $k = 32$ decoys sampled for each spectrum in the minibatch, resulting in a total of $O(bk)$ candidate formula encoded per batch. To sample the most plausible ‘‘hard’’ negative decoys in each batch, we utilize a FastFilter module described in detail below. All models are implemented in Python version 3.8 using PyTorch version 1.9 and PyTorch Lightning [140] version 1.6 and trained with the Adam optimizer. [141]. Hyperparameters are optimized with Ray Tune [124] version 2.0. Each model was trained on a single NVIDIA RTX A5000 GPU with training times taking under 3 hours of wall time.

3.2.4 Generating formula candidates

By default, MIST-CF utilizes the mass decomposition algorithm embedded within the SIRIUS software [50], [142], [143] (independent module from the rest of their pipeline) to decompose MS1 precursor masses into candidate formula options. We do not restrict the number of common elements C, N, O, or H. We set the maximum number of S and P atoms to 5 and 3, respectively, and limit each halogen (i.e., F, Cl, Br, I) to a maximum of one per formula, allowing for the recovery of 96% of formulae in the public NPLIB1 dataset. This constraint can be changed at inference time as desired. The chemical filter option “COMMON” is used to generate formulae for energy-based model training, and the “RDBE” filter is applied during inference. The default mass deviation is set to 10 ppm during training for all spectra by default, unless otherwise specified. For retrospective analysis on the NPLIB1 and NIST20 datasets, we resample individual parent masses using a truncated Gaussian centered around the exact mass with a standard deviation of 1/5 the instrument-specific ppm tolerance (Unknown/Ion Trap: 15 ppm; Q-ToF: 10 ppm; Orbitrap/FTICR: 5 ppm) as defined in BUDDY [52]. We utilize SIRIUS version 5.6.3.

During inference, the user can choose to avoid this step in if they have their own, more narrow list of potential formulae candidates (e.g., generated by database search, knowledge about the chemical space being measured, or external tools such as BUDDY). In such cases, the exhaustive formula candidate generation step can be skipped and MIST-CF’s energy function can be used to directly evaluate the input candidate list.

3.2.5 Downselecting candidate formulae with a learned filter

Due to its energy-based formulation, the computational cost of MIST-CF scales linearly with the number of candidate precursor formulae. The evaluation of each formula requires assigning subformulae to peaks within a fragmentation spectrum, which is far faster than inducing a tree structure over subformulae but not negligible. To limit the space of candidate formulae, in addition to utilizing heuristics such as COMMON or RDBE as mentioned earlier, we train a data-driven filter we term FastFilter (Figure 3.2A) to further narrow the option set.

$G_{\theta}^{(\text{FF})}$ is a feedforward neural network that takes as input an encoded precursor formula candidate, $\text{Enc}(\mathcal{F})$, and learns, in the same fashion as MIST-CF, to predict an energy value based solely on the formula—not information about the spectrum—to approximate a non-normalized likelihood $p_{\theta}(\mathcal{F}_{(i)}|M)$. Because no spectrum information is needed, we train a single FastFilter model using a large database of biologically-relevant molecules, containing 200,388 molecular formulae extracted from various sources such as KNApSAcK, [118] HMDB, [26] KEGG, [119], and PubChem [128], as prepared by Duhrkop et al. [57] To avoid dataset leakage, all molecular formulae masses that appear in spectral libraries used for model training/testing are excluded; the remaining masses are split in an 80%/10%/10% ratio for training, validation, and testing of FastFilter. We use $G_{\theta}^{(\text{FF})}$ to select the top k candidates during inference (256 by default), or decoys during training, to further score with $G_{\theta}^{(\text{MIST})}(\mathcal{F}, \mathcal{A}, \mathcal{S})$. Formulae are represented as input to the positive sinusoidal embedding vectors as in MIST-CF, and exact training parameters including learning rate, hidden size, and layers are described in Section 3.6.

3.2.6 Datasets

We evaluate MIST-CF in terms of its ability to predict precursor molecular formulae for MS/MS spectra from NPLIB1, a public natural products dataset extracted from the GNPS database [28]. NPLIB1 is prepared as in Goldman, Bradshaw, Xin, *et al.* [130] and extracted from Duhrkop *et al.* [57] The dataset contains positive mode MS/MS spectra for compounds under 1,500 Da containing a predefined set of elements and adducts. In total NPLIB1 contains 10,709 spectra, 8,553 unique structures and 5,433 unique molecular formulae. We employ molecular formula splits in which 20% of the remaining 5,433 molecular formulae are selected randomly and added—with their corresponding spectra—to the test set. 10% of the remaining data is used for validation and early stopping.

In addition to the public data, we use the commercial NIST20 library [63] to supplement the training dataset. We extract all Orbitrap high resolution positive mode mass spectra containing common elements and adducts. Examples with molecular formulae found in the test set are excluded to avoid biasing the model. In total, the combined dataset has 45,838 unique spectra, 30,950 unique 2D molecular structures, and 15,315 unique molecular formula. By using identical public test sets, we are able to report performance metrics with and without commercial library inclusion to enable replication studies and future methodological improvements.

3.2.7 Baseline models

In addition to learning G_θ using the MIST architecture, we select three separate baseline neural architectures: a feed forward network (FFN) inspired by MetFID [56] that acts on a binned representation of the spectrum, $G_\theta^{(\text{FFN})}$; ‘MS1 Only’, a variant of $G_\theta^{(\text{FFN})}$ in which the binned spectrum is set to $\mathbf{0}$ for all spectra; and a Transformer model that utilizes multiscale sinusoidal embeddings (Transformer) [83], $G_\theta^{(\text{TF})}$. These models are trained and hyperparameter optimized equivalently to MIST-CF. In the FFN baseline, the binned spectrum representation is concatenated to the encoded MS1 candidate formula, the one-hot adduct representation, and the context vector, then passed into a multilayer perceptron module. The Transformer baseline concatenates the MS1 candidate formula encoding and context vector to the sinusoidal embedding at each of the top 100 most intense peaks, along with the intensity. An additional ‘cls’ token is added to the peaks containing the mass of the MS1 candidate and intensity of 2. These are subsequently passed into a set of multi-head attention Transformer layers. The output is pooled at a special ‘cls’ token. A single linear layer then predicts a scalar energy value.

3.3 Results

3.3.1 FastFilter constrains candidate formulae with high precision

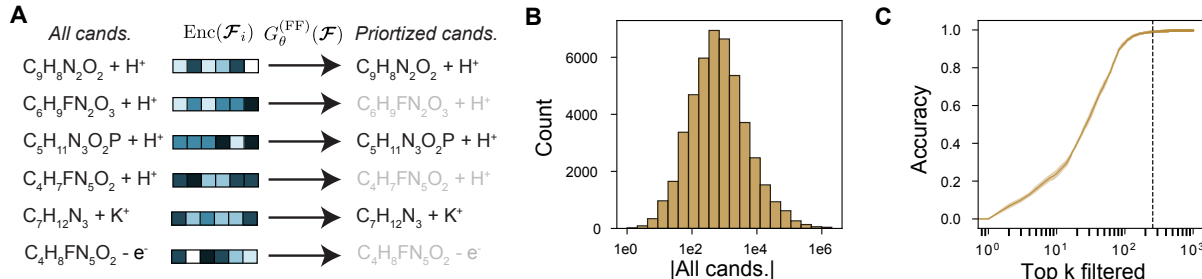


Figure 3.2: **Large numbers of formula candidates can be quickly filtered with a simple feed forward neural network, FastFilter.** **A.** Generated candidate formulae for a spectrum can be prioritized with a learned model. A scalar score is generated for each candidate formula given an MS1 value and mass tolerance via a feed forward neural network. Only the top ranked formulae are selected for consideration with MIST-CF. **B.** The distribution of the number of candidates for all recorded MS1 spectra in our spectra libraries NPLIB1 and NIST20. All candidate molecular formulae are generated for all 6 adduct types considered with an allowed mass deviation of 10 ppm. **C.** FastFilter model accuracy at various top k cutoff values. The top formula is nearly always recovered within the top 256 candidates. The 95% confidence interval of the mean for recovery is shown across the 3 different formula splits considered in this work. All results are computed for a single trained FastFilter model on a large database of biologically-relevant molecules.

A key limiting factor in *de novo* formula identification is the large size of the candidate space, particularly for molecules with higher masses. We first quantify the size of the candidate space. We process the full set of training spectra, including both NPLIB1 and NIST20, into molecular formula candidates with the permissive RDBE filter. Over 15% of spectra have greater than 5,000 formula candidates (Figure 3.2B). By training a light-weight model, FastFilter, to predict how likely each formula is to appear in a biologically-relevant database of small molecules, we are able to filter the candidates down to a smaller subset. We are able to recover the true formula 99% of the time within the top 256 candidates (Figure 3.2C). This guarantees the computational tractability of the MIST-CF pipeline, as subformula labeling would be prohibitively expensive for spectra with hundreds of thousands or millions of candidate formulae.

3.3.2 Molecular Formula Transformers provide meaningful representations

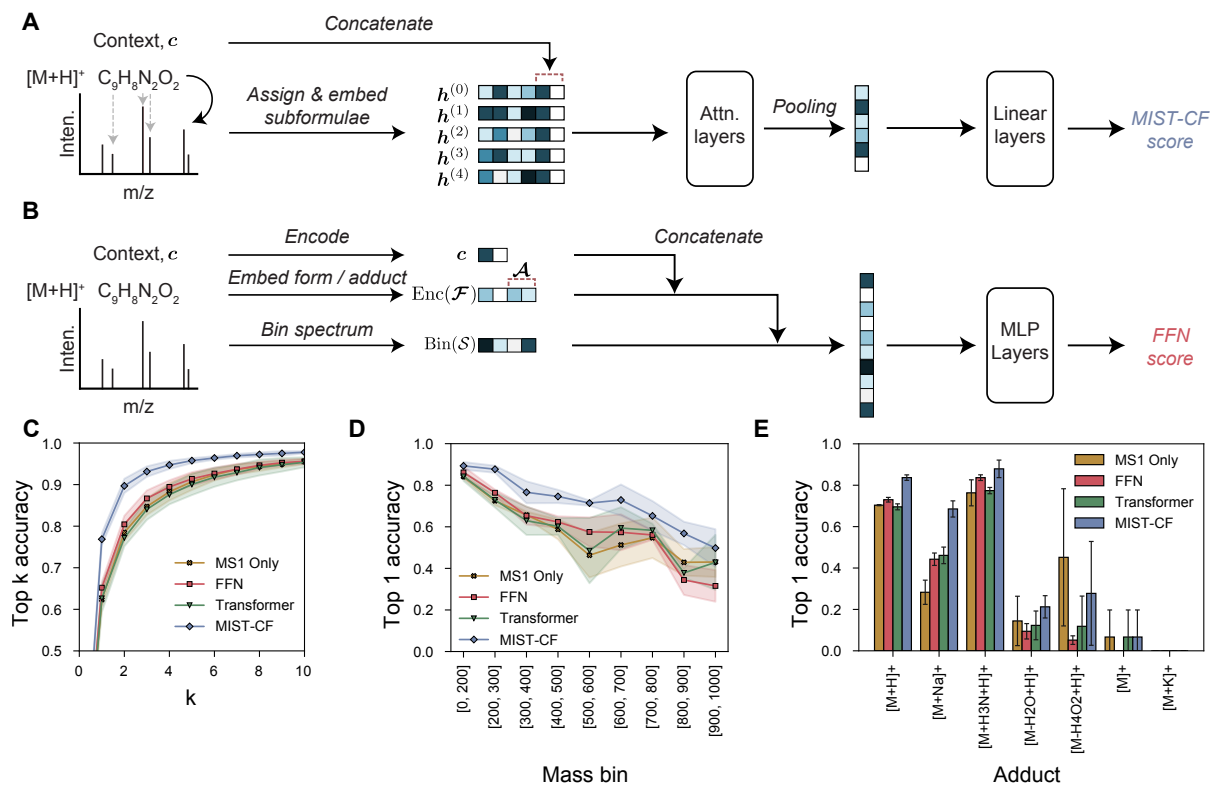


Figure 3.3: MIST-CF is a highly-effective architecture for learning to rank plausible MS1 formula annotations **A**. The MIST-CF architecture uses the candidate molecular formulae to generate subformulae and encode these with a Formula Transformer. **B**. The baseline feed forward network (FFN) separately encodes the spectra and formula before feeding their concatenation into a multilayer perceptron (MLP). **C**. For all spectra in the test set, the fraction recovered at various top k values for all methods is computed and shown. **D-E**. The top 1 accuracy for the methods is grouped by the mass of the MS1 precursor or adduct type. All results are computed for MIST-CF and the following baselines: MS1 only (a feed forward network utilizing only the molecular formula and context vector), a feed forward network, and a Transformer. All results are computed over 3 random formula splits and respective training runs as described in Section 3.2, where the NIST20 is included in the training sets. All spectra include up to 256 candidates to select from as selected with the “COMMON” filter [50] and FastFilter. Error bars and shaded regions show 95% confidence intervals of the mean.

We evaluated four different architectures for encoding spectra and formulae to determine the best neural network architecture for the energy-based modeling framework described in Section 3.2. MIST-CF uses a Transformer model to convolve upon MS2 subformula with context information concatenated to each formula (e.g., adduct, instrument type) (Figure

3.3A). We compare this model to a feed forward neural network (Figure 3.3B), a standard Transformer applied to sinusoidal embeddings of each m/z value [83], and a variant of the feed forward network that does not embed the spectrum (“MS1 only”). Importantly, all baseline neural network models are provided with the same context vectors, training set, and hyperparameter optimization schemes to enable fair comparison.

We find that MIST-CF outperforms all other architectures tested, likely due to its ability to explicitly combine the representation of the spectrum and formula candidate, rather than via intra-architecture concatenations (Figure 3.3B). We test the model by holding out a number of spectra and molecular formulae from model training. For each test spectrum, up to 256 formula and adduct candidates are re-ranked by each model. MIST-CF outperforms the next best model at top 1 accuracy by a $> 10\%$ margin (Figure 3.3C; Table 3.1). Curiously, the feed forward network and Transformer network show little difference from each other, indicating that neural network architecture (i.e., FFN vs. Transformer) has less impact than the input information provided to the model (i.e., subformula labels)—the main strength of MIST-CF. All models are able to perform at least as well as the MS1 only model, consistent with our intuition that the models are learning more than just database bias.

All models decrease in accuracy as the masses of compounds increase (Figure 3.3D) due to the growing number of plausible candidates. Similarly, models struggle on adducts that appear less often in the training dataset, such as potassium adducts (Figure 3.3E). This is one such area where manually parametrized models such as SIRIUS [50] may be better suited to generalization. Data-driven methods such as MIST-CF are empirically better at correctly retrieving formulae with common proton or sodium adducts.

Integrating the higher quality Orbitrap training spectra from NIST20 leads to improved performance on the same test set. MIST-CF models trained on NPLIB1 alone achieve a top 1 accuracy of 0.741 compared to a top 1 accuracy of 0.769 for models that train on the NIST20 in addition (Table 3.1). This $> 2\%$ absolute improvement underscores that formula prediction accuracy may be enhanced with the upcoming release of new spectral libraries and higher quality data. [29]

Table 3.1: **MIST-CF outperforms comparable neural network baselines at molecular formula annotation from MS/MS.** Models were trained to predict held out NPLIB1 test examples using training sets consisting of “NPLIB1” or “NPLIB1 + NIST20.” The best value in each column is typeset in bold. Models were evaluated using three independent formula splits. Values are shown \pm standard errors of the mean.

| Training dataset | NPLIB1 | | | NPLIB1 + NIST20 | | |
|------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|
| Top k | 1 | 2 | 3 | 1 | 2 | 3 |
| MS1 Only | 0.609 \pm 0.002 | 0.773 \pm 0.003 | 0.833 \pm 0.003 | 0.623 \pm 0.007 | 0.785 \pm 0.008 | 0.847 \pm 0.008 |
| FFN | 0.635 \pm 0.009 | 0.786 \pm 0.009 | 0.847 \pm 0.007 | 0.652 \pm 0.006 | 0.804 \pm 0.009 | 0.867 \pm 0.008 |
| Transformer | 0.639 \pm 0.006 | 0.791 \pm 0.009 | 0.851 \pm 0.006 | 0.626 \pm 0.014 | 0.772 \pm 0.008 | 0.840 \pm 0.011 |
| MIST-CF | 0.741 \pm 0.010 | 0.878 \pm 0.009 | 0.919 \pm 0.007 | 0.769 \pm 0.006 | 0.897 \pm 0.007 | 0.931 \pm 0.005 |

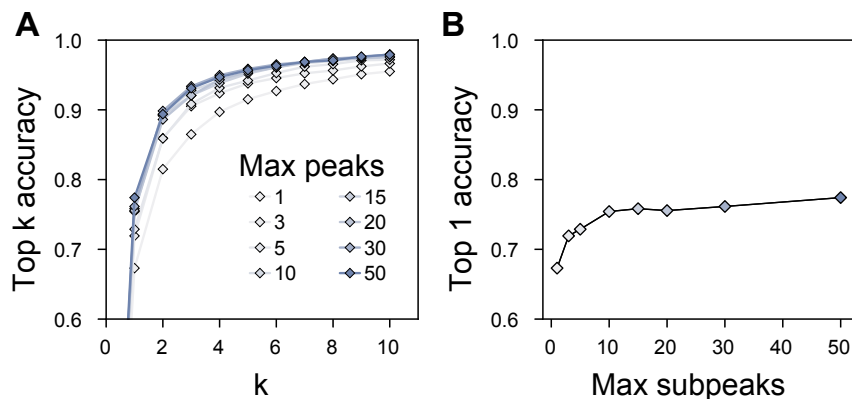


Figure 3.4: **MS/MS subpeaks drive MIST-CF performance.** **A.** The (sub)Formula Transformer in MIST-CF operates on only the N_p highest intensity MS2 peaks. Here, from that set of labeled MS2 peaks, the number of peaks used as input to the spectrum is limited. As the number of included peaks increases, the formula retrieval accuracy does too. **B.** The top 1 retrieval accuracy is shown for all maximum subpeak numbers. All results are shown for MIST-CF trained on the joint NPLIB1 and NIST20 dataset and tested on a single test split of the data.

One of the core methodological decisions in MIST-CF is the application of a Formula Transformer to the set of labeled subformulae in the MS2 spectrum. To verify that these peaks are informative for the model, we repeated benchmarking experiments for a single split of the data, this time modulating the maximum number of peaks (rank-ordered by intensity) viewable by MIST-CF.

Model performance rapidly increases as the maximum number of included peaks increases, sharply rising from 0.673 to 0.719 for 1 and 3 spectrum peaks respectively, where $N_p = 0$ is functionally equivalent to the MS1 Only baseline with an accuracy of 0.623 (Table 3.1). There are diminishing returns of including lower intensity peaks, with 1, 5, 10, 20, 50 maximum formula peaks achieving top 1 accuracies of 0.673, 0.729, 0.754, 0.756, and 0.774 respectively (Figure 3.4). Making the model aware of a greater number of fragments enables more generalizable and accurate predictions. This builds confidence that the model is learning more than database biases, drawing information from even minor peaks to inform predictions. As can be seen, the absolute difference in performance appears to level off beyond 20 with only more marginal increases toward $N_p = 50$. Given that the performance benefit is marginal for many more peak annotations, we maintain our peak default $N_p = 20$ unless otherwise stated to reduce the runtime of the method.

3.3.3 MIST-CF compares favorably to existing formula annotation tools

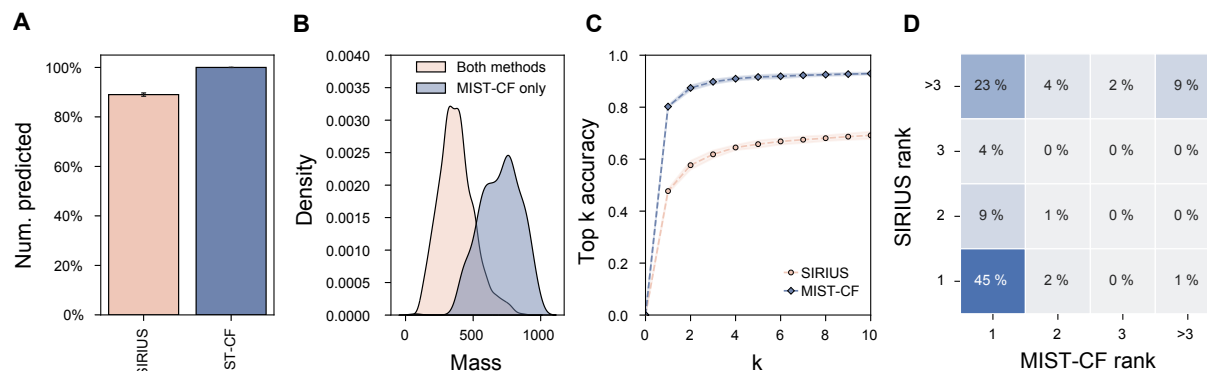


Figure 3.5: **MIST-CF more frequently assigns correct molecular formulae than the SIRIUS formula assignment module on the NPLIB1 test set.** **A.** The number of spectra for which each method is able to predict a formula—regardless of accuracy. **B.** The distribution of molecular masses for the spectra both methods are able to annotate—regardless of accuracy— compared to the distribution of molecular masses for spectra on which only MIST-CF succeeds. **C.** Top k accuracy for both methods is shown. **D.** The rank at which each method is able to recover the true molecular formula. This visualization highlights that there are many spectra for which MIST-CF achieves rank 1 that SIRIUS does not predict in its top 3. MIST-CF models are trained on the joint NIST20 and NPLIB1 dataset. SIRIUS is executed with a compound and tree timeout of 300 seconds. All values are shown for 3 random formula splits of the NPLIB1 data. Error bars and confidence intervals show 95% confidence intervals for the standard error of the mean. Accuracy is computed with respect to the *total* MS1 composition (i.e., summed adduct and formula) to avoid biasing results against SIRIUS.

The widely used SIRIUS tool [50] is the *de facto* state of the art for the task of molecular formula annotation. We therefore perform a head-to-head comparison of MIST-CF and SIRIUS on the same NPLIB1 test set. Herein, MIST-CF considers all candidate MS1 formulae within 10 ppm of the recorded MS1 mass and utilizes FastFilter to down select to 256 candidates for full evaluation. We run the SIRIUS formula module from the command line with tree and compound timeouts of 300 seconds to avoid excessive execution times. With these constraints, MIST-CF is able to predict a formula for every spectrum, whereas SIRIUS fails for 10.77% of the spectra (Figure 3.5A). Failed spectra are associated with larger masses on average (700.80 Da) than successful ones (378.06 Da).

Model accuracy is evaluated in terms of the summed elemental composition of the formula and adduct. SIRIUS does not distinguish tree scores with different adduct assignments as we do in MIST-CF by design (e.g., $[\text{C}_6\text{H}_{12}\text{O}_6+\text{H}]^+$ and $[\text{C}_6\text{H}_{14}\text{O}_7+\text{H}-\text{H}_2\text{O}]^+$ would appear to have equivalent tree scores). On the NPLIB1 test set when using full RDBE decoys (i.e., candidate formula not constrained with the “COMMON” filter as in Table 3.1), MIST-CF successfully predicts molecular formulae with a 71% top 1 accuracy (80% top 1 when considering joint

formula and adduct accuracy), compared to SIRIUS’s 48% accuracy (evaluated on joint formula and adduct predictions). This represents a >20% improvement in absolute top 1 prediction accuracy (Figure 3.5C). On a per-spectrum basis, SIRIUS rarely makes correct predictions MIST-CF does not, with a total of approximately 3% of spectra falling into this category, whereas MIST-CF appears to predict 36% of spectra correctly when SIRIUS cannot (Figure 3.5D).

Beyond the improvement in accuracy, MIST-CF required approximately one-third the wall time of SIRIUS (evaluated on compounds under 700 Da for which SIRIUS does not time out) when both methods are run on a single CPU core (Section 3.6).

3.3.4 MIST-CF is competitive on the CASMI2022 challenge

We selected the dataset released as part of the Comparative Assessment of Small Molecule Identifications 2022 (CASMI2022) [75] for additional validation. Our training datasets were generated prior to the competition announcement, minimizing the risk of training dataset bias and simulating a prospective use case. Herein, we focus only on the formula identification challenge rather than the full task of structural elucidation.

We extract all positive mode MS/MS files from the provided mzML files using MZmine 3 [120] and apply both MIST-CF and SIRIUS to predict molecular formulae for each of the extracted 304 spectra. Using our constraints on element types, 296 of the 304 formulae are recoverable (97%), excluding species such as iodixanol ($C_{35}H_{44}I_6N_6O_{15}$). We use an MS1 tolerance of 5 ppm.

Model performance was evaluated on three key metrics: the accuracy of predicting formula correctly, the accuracy of predicting the adduct correctly, and the accuracy of predicting the total elemental composition of the formula and adduct combined. As discussed above, we use this last metric because evaluating the accuracy of the molecular formula alone would bias the comparison for MIST-CF.

We consider three variants of SIRIUS: SIRIUS, SIRIUS (CSI:FingerID) and SIRIUS (Submission). SIRIUS (CSI:FingerID) provides formulae as re-ranked by their CSI:FingerID score when searched against PubChem. SIRIUS (Submission) uses predictions submitted by the SIRIUS authors at the latest competition under file name “duehrkop_CASMI2022.csv.” This submission reportedly used a mix of ion identity molecular networking [144], which can be used to more accurately resolve adduct types, along with manual curation to improve performance.

Encouragingly, we find that MIST-CF is competitive with SIRIUS (Submission), achieving a nearly equivalent collective joint formula and adduct accuracy of 0.862, despite our automated classification and no additional manual curation (Table 3.2). If evaluating SIRIUS (Submission) on the subset of 272 spectra for which a prediction was submitted, a higher accuracy is achieved compared to MIST-CF, as MIST-CF predicts *all* test set spectra; of the 32 spectra on which MIST-CF outperforms SIRIUS (Submission), SIRIUS (Submission) only submitted predictions for 7. This is in contrast to the 34 spectra on which SIRIUS outperforms MIST-CF. Nevertheless, when using default parameters, MIST-CF reaches a top 1 formula accuracy of 0.842 (including adduct prediction) compared to an accuracy of 0.516 for SIRIUS (CSI:FingerID). These results illustrate the competitiveness of MIST-CF and demonstrate that accurate, prospective formula annotation does not require computing

full fragmentation trees. In addition to testing our default MIST-CF model, we also test a variant with $N_p = 50$ as a comparison. Performance is equivalent for the joint accuracy of formula and adduct pairs. However, accuracy in predicting the formula alone increases, as the model appears to be marginally better at adduct assignment.

Table 3.2: **Model accuracy of MIST-CF and SIRIUS evaluated on CASMI2022.** “Accuracy (\mathcal{F})” indicates accuracy of predicting the exact formula correctly. “Accuracy (\mathcal{A})” indicates the fraction of spectra for which the first predicted (formula, adduct) pair includes the correct adduct. “Accuracy ($\mathcal{F} + \mathcal{A}$)” indicates the fraction of spectra for which the total elemental composition of the top predicted formula and adduct matches the elemental composition of the summed true formula and adduct. “Predicted” indicates the number of spectra for which a formula and adduct prediction was made. “SIRIUS” makes predictions with the default command line tool described above; “SIRIUS (CSI:FingerID)” predicts formula, adduct pairs using CSI:FingerID rankings against PubChem; SIRIUS (Submission) is taken directly from the previous CASMI2022 results. MIST-CF (20 peaks) uses default 20 subformula peaks, whereas MIST-CF (50 peaks) utilizes 50 subformula peaks. Accuracy (\mathcal{F}) and Accuracy (\mathcal{A}) are not reported for SIRIUS because the method does not distinguish between formula, adduct pairs that sum to the same molecular formula, and our metrics default to the worst-case ranking for ties (i.e., somewhat unfairly penalizing methods such as SIRIUS that have more ties). All accuracies are reported for the top 1 prediction and divided by the total number of spectra (304), not the total number of predicted spectra. Numbers typeset in bold indicate the best result in the column.

| Method | Accuracy (\mathcal{F}) | Accuracy (\mathcal{A}) | Accuracy ($\mathcal{F} + \mathcal{A}$) | Predicted |
|-------------------------------------|-------------------------------|-------------------------------|---|------------|
| SIRIUS version 5.6.3 | - | - | 0.641 | 274 |
| SIRIUS version 5.6.3 (CSI:FingerID) | 0.516 | 0.543 | 0.678 | 254 |
| SIRIUS (Submission) | 0.865 | 0.855 | 0.868 | 272 |
| MIST-CF (20 peaks) | 0.842 | 0.901 | 0.862 | 304 |
| MIST-CF (50 peaks) | 0.822 | 0.885 | 0.868 | 304 |

3.4 Conclusion

We have introduced a data-driven neural network model, MIST-CF, for inferring molecular formulae from MS/MS spectra, trained using an energy-based modeling framework. We benchmark this model extensively to show how our recent Molecular Formula Transformer architecture is uniquely suited to this task of integrating formula and spectra information. MIST-CF outperforms other learning-to-rank neural network architectures and ties the winning solution at the CASMI 2022 competition within the positive mode category, despite using zero MS1 isotopic information or manual prediction refinement.

This work defines a clear problem formulation for learning to rank MS1 candidate formula from MS/MS data, including open source code and data splits. To address this task, we release open source and trained models with low memory and run-time concerns, even for large molecules. In addition to these models, we develop a smaller and light-weight neural network

formula filtering model to help prioritize biologically-relevant molecular formula candidates. This work continues to demonstrate the efficacy of our recent Molecular Formula Transformer network architecture in thorough benchmarking comparisons [64]. As part of this, we have greatly simplified and improved the the Molecular Formula Transformer implementation to now utilize custom formula embeddings, a simple and open source subformula assignment routine (i.e., no longer fragmentation trees), additional model inputs such as instrument type, and multiple adduct types. Applying these same changes to the Molecular Formula Transformer architecture for fingerprint prediction is likely to yield similar improvements.

There are many avenues for improving upon MIST-CF. We have trained models only for positive-mode data; we do not directly address the integration of MIST-CF scores with MS1 isotopic scores; we do not consider adduct switching in MS2 subformula assignment; we still rely upon SIRIUS’s algorithmic decomposition of exact masses into formula candidates due to their fast implementation; and we have not yet explored the use of forward structure-to-spectrum models [130], [145] for data augmentation. There is also an opportunity to combine MIST-CF with the recently reported BUDDY [52] by re-ranking formulae generated by BUDDY rather than relying on the FastFilter.

Altogether, we are optimistic about the potential to integrate this model into existing pipelines for small molecule metabolite identification. By addressing the task of molecular formula annotation, this work moves us one step closer to our vision of an integrated neural network driven metabolite annotation pipeline.

3.5 Data and Software Availability

All code to replicate experiments, train new models, and load pre-trained models is available at <https://github.com/samgoldman97/mist-cf>. The code to parse NIST can be found in <https://github.com/samgoldman97/nist-parser> but this optional subset of training data requires a license and cannot be shared publicly. The exact repository version used in this work has been archived at Zenodo record 8151490 (data) and Zenodo record 8151513 (code).

3.6 Additional Results

3.6.1 SIRIUS configuration

SIRIUS version 5.6.3 [50] is used as a baseline throughout the work. For consistency, we utilize a 300 second (5 minute) timeout for each compound and set equivalent adducts and formula as used in MIST-CF. An example call signature using our parameters, as we use in CASMI2022 is shown in Listing 3.1.

Listing 3.1: An example call to SIRIUS for annotating the formula for an input MGF file, \$INPUT_MGF

```
$SIRIUS_PATH \
  --cores $CORES \
  --output $OUTFOLDER1 \
  --input $INPUT_MGF \
  formula \
  -i "[M+H]+,[M+K]+,[M+Na]+,[M+H-H2O]+,[M+H-H4O2]+,[M+NH4]+,[M]+" \
  -e "C[0-]N[0-]O[0-]H[0-]S[0-5]P[0-3]I[0-1]Cl[0-1]F[0-1]Br[0-1]" \
  --ppm-max 5.0 \
  --tree-timeout 300 \
  --compound-timeout 300 \
  fingerprint \
  structure \
  write-summaries \
  --output $OUTFOLDER2
```

3.6.2 Timing experiments

To demonstrate that MIST-CF has reasonable runtime properties, we conduct a simple timed experiment. 10 random sample test set spectra under 700 Da from the NPLIB1 were selected and set aside in an MGF file. Using the command line call signature for formula prediction on a Linux workstation, we test the wall time to predict formula for these spectra with all pre-specified elements, a maximum ppm of 10 from the MS1, and all adduct possibilities considered within MIST-CF. We repeat the same procedure for SIRIUS, using a compound and tree timeout of 1 minute to avoid excessively large times. We conduct these experiments using only 1 CPU core. We find that in three separate trials, MIST-CF completed in 39.31 ± 0.49 seconds (mean \pm standard error of the mean), whereas SIRIUS completed in 112.61 ± 0.21 seconds (mean \pm standard error of the mean).

3.6.3 Hyperparameter setting

To enable fair comparison across models, hyperparameters were tuned for MIST-CF, the FFN binned prediction baseline, the transformer binned prediction baseline, and the FastFilter model. Parameters were tuned using RayTune [124] with Optuna [123] and an ASHAScheduler. Each model was allotted 50 different hyperoptimization trials for fitting. MIST-CF, the FFN baseline, and the Transformer baseline were hyperparameter optimized on a smaller 10,000 spectrum training subset of the combined NPLIB1 and NIST20 datasets. MIST-CF was hyperparameter optimized using a maximum of 10 peaks per spectrum. The FastFilter model was hyperparameter optimized directly on the biological molecules dataset. Parameters are listed in Table 3.3.

Table 3.3: **Hyperparameter settings for MIST-CF, FFN, Transformer, and the FastFilter model to select batches of molecular formula.**

| Model | Parameter | Grid | Value |
|-------------|------------------------------|--------------------------------|---------|
| MIST-CF | Learning rate | [1e-4, 1e-3] | 0.00045 |
| | Learning rate decay fraction | [0.7, 1.0] | 0.8830 |
| | Weight decay | {1e-6, 1e-7, 0.0} | 0.0 |
| | Hidden size | {32, 64, 128, 256} | 128 |
| | Dropout | [0, 0.5] | 0.1 |
| | Number of layers | [1, 2, 3, 4] | 2 |
| | Batch size | {4, 8, 16} | 4 |
| FFN | Number of bins | {1000, 2000, 3000, ..., 10000} | 5000 |
| | Learning rate | [1e-4, 1e-3] | 0.0005 |
| | Learning rate decay fraction | [0.7, 1.0] | 0.8477 |
| | Weight decay | {1e-6, 1e-7, 0.0} | 0 |
| | Hidden size | {32, 64, 128, 256, 512} | 512 |
| | Dropout | [0, 0.5] | 0.4 |
| | Batch size | {16, 32, 64} | 64 |
| Transformer | Learning rate | [1e-4, 1e-3] | 0.00025 |
| | Learning rate decay fraction | [0.7, 1.0] | 0.8601 |
| | Weight decay | {1e-6, 1e-7, 0.0} | 1e-7 |
| | Hidden size | {32, 64, 128, 256, 512} | 32 |
| | Dropout | [0, 0.5] | 0.2 |
| | Number of layers | {1, 2, 3, 4} | 2 |
| | Batch size | {8, 16} | 8 |
| FastFilter | Learning rate | [1e-4, 1e-3] | 0.0003 |
| | Learning rate decay fraction | [0.7, 1.0] | 0.8642 |
| | Weight decay | {1e-6, 1e-7, 0.0} | 0 |
| | Hidden size | {32, 64, 128, 256, 512} | 256 |
| | Dropout | [0, 0.5] | 0.1 |
| | Number of layers | {1, 2, 3, 4} | 3 |
| | Batch size | {16, 32, 64} | 64 |

Chapter 4

Generating Mass Spectra as Subformula Sets with Prefix-Trees

This work has previously appeared as a conference publication, S. Goldman, J. Bradshaw, J. Xin, *et al.*, “Prefix-Tree Decoding for Predicting Mass Spectra from Molecules,” in *Advances in Neural Information Processing Systems 37*, 2023. I jointly conceptualized and subsequently led this project. John Bradshaw contributed to the experiments, implementation of various features of the model, and was especially helpful in shaping the writing and presentation of this work.

4.1 Introduction

As the primary tool to discover unknown small molecule structures from biological samples, tandem mass spectrometry (MS/MS) experiments have enabled the identification of numerous important molecules implicated in health and disease [13], [67], [146]. Tandem mass spectrometers are capable of isolating, fragmenting, and measuring the resulting fragment masses of small molecules from a sample, producing a signature (a mass spectrum) for each detected molecule (Figure 4.1, top).

Computationally predicting mass spectra from molecules *in silico* (Figure 4.1, bottom) is thus a longstanding and important challenge. Not only does this assist practitioners in better understanding the fragmentation process, but it also enables the identification of molecules from newly observed spectra by comparing an observed spectrum to virtual spectra generated from a database of candidate molecules. While a large library of empirical mass spectra could theoretically serve the same purpose, the size of such libraries is limited by the slow and expensive process of acquiring pure chemical standards and measuring their spectra, motivating computational prediction.

We argue that there are three core, interrelated desiderata for a forward molecule-to-spectrum simulation model, or “spectrum predictor”. An ideal spectrum predictor should be (i) *accurate*, being able to predict the exact set of fragment masses and intensities with a precision comparable to experimental measurements; (ii) *physically inspired*, to avoid making physically nonsensical (“invalid”) suggestions and to provide interpretations of the chemical species responsible for each peak for the benefit of human expert chemists; and (iii) *fast*, such

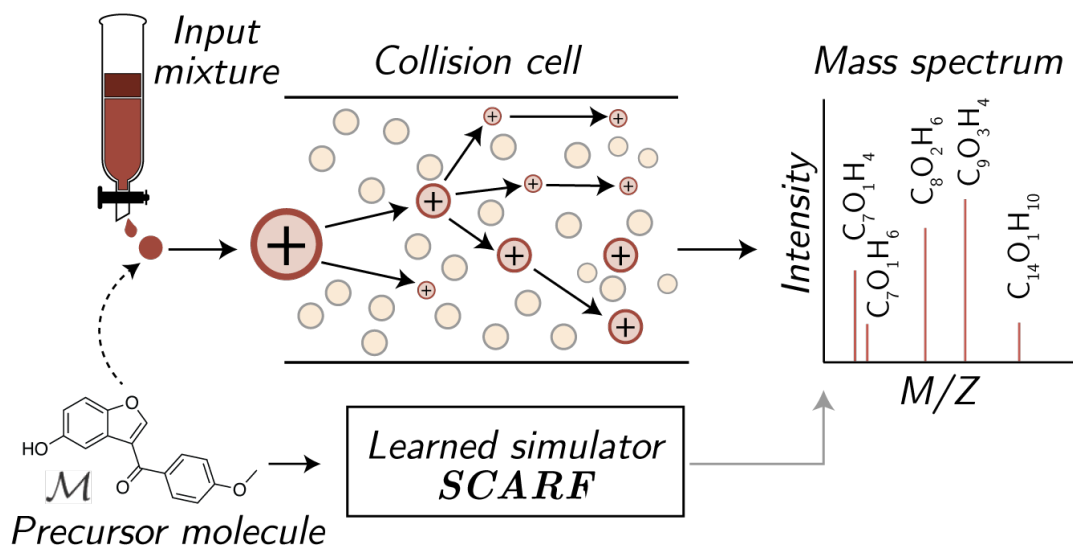


Figure 4.1: **Tandem mass spectrometers measure fragmentation patterns of molecules, resulting in characteristic peaks that are indicative of their structure.** SCARF simulates these fragmentation patterns *in silico*.

that it is computationally inexpensive to predict spectra for many (e.g., millions) hypothetical molecules.

Unfortunately, many existing spectrum predictors do not meet these criteria. Methods to date have tended to follow one of two approaches: (a) physically motivated fragmentation approaches or (b) molecule-to-vector (or “binned”) approaches (Figure 4.2A-B). Fragmentation approaches (e.g., [71], [73], [95], [147]; Figure 4.2A) take an input molecule and suggest bonds that may break, creating fragments that are scored by ML algorithms or curated rulesets. While interpretable, these methods are often slow and restrictive; certain mass spectrum peaks are generated by complex chemical rearrangements within the collision cell that cannot be approximated by bond breaking alone. That is, the bonds in observed fragments are not a subset of those in the original molecule [148], [149]. On the other hand, binned prediction approaches (e.g., [84], [86], [150]; Figure 4.2B) are less physically grounded, using neural networks to directly learn a mapping from molecules to vectors representing discretized versions of the spectra. These methods, while fast, lack interpretability and due to discretization have a mass precision lower than that of most modern spectrometers, limiting their accuracy.

We propose to address the shortcomings of previous work by predicting mass spectra from molecules at the level of molecular formulae (e.g., $C_xN_yO_zH_w\dots$) and introduce a new method, Subformulae Classification for Autoregressively Reconstructing Fragmentations (SCARF) to do so. Because the molecular formula for each input molecule is known, each subformula in the predicted set of peaks is constrained to contain a subset of the atoms in the original formula. Our primary contributions are:

- posing mass spectrum prediction as a two step process: first generating the set of molecular formulae for the fragments, then associating these formulae with intensities;

- overcoming the combinatorial subformula option space by learning to generate formula prefix trees;
- demonstrating the empirical benefit of SCARF in predicting experimental mass spectra quickly and accurately using two separate datasets, providing a benchmark for future work.

4.2 Background

We provide a short introduction of tandem mass spectrometry suitable for a general machine learning audience, detail previous approaches to modeling this process as they relate to our proposed approach SCARF, and explain how such tools can be utilized to discover molecules from new spectra. We refer interested readers to [151] for further details on the physical process of mass spectrometry.

4.2.1 Tandem mass spectrometry

Tandem mass spectrometers (MS/MS) measure fragmentation patterns of molecules in a multi-stage process. The input to the process is a solution containing a *precursor* molecule, $\mathcal{M} \in \mathcal{X}$, associated with a molecular formula, \mathcal{F} , defining the counts of each element present; for instance $\mathcal{F} = \text{C}_{16}\text{O}_4\text{H}_{12}$ for the precursor molecule shown in Figure 4.1. The precursor molecule is first ionized (i.e., made charged), often by bonding or associating with an *adduct* (e.g., a proton, H^+) present in the solution. The charged product is then measured by a mass analyzer (MS1), where its mass-to-charge ratio (m/z) is measured.

This precursor ion is then filtered into a *collision cell*. Here, through interactions with an inert gas, the precursor ion is broken down into a set of one or more *product ions*, each of which is associated with a new chemical formula; for example, one might be $\mathbf{f}^1 = \text{C}_7\text{OH}_4$ for the process shown in Figure 4.1. Finally, this set of product ions is measured by a second mass analyzer (MS2), along with the set of their intensities, $y^i \in \mathbb{R}^+$ (i.e., their relative frequencies over several repetitions of this process), creating for each ion what is referred to as a *peak*. The collection of all peaks makes up a molecule’s *mass spectrum*, and is commonly represented as a plot of intensities versus m/z (Figure 4.1, right).

4.2.2 Predicting mass spectra from molecules (spectrum predictors)

Fragmentation prediction. A complex but physically grounded strategy is to model the bond breakage processes occurring in the collision cell (Figure 4.2A). Examples include MetFrag [71], MAGMa [95], and CFM-ID [73], which recursively fragment molecules (either bond or atom removals) to generate fragment predictions. These methods combine expert rules and local scoring methods to enumerate molecular fragmentation trees to predict spectra. CFM-ID [73] learns subsequent fragmentation transition probabilities between fragments with an expectation maximization algorithm to determine intensities at each fragment. Rule-based methods and full tree enumeration reduce the flexibility of these approaches, and

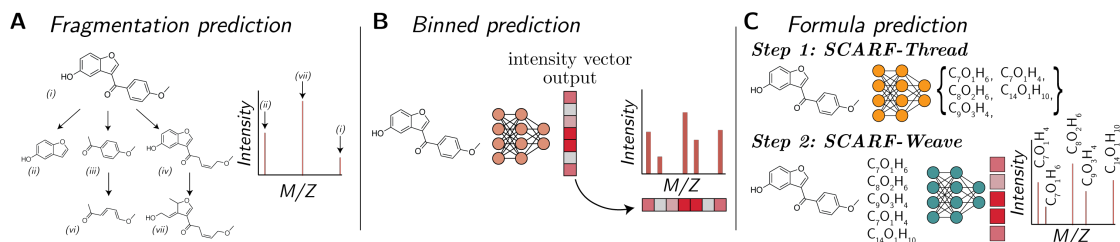


Figure 4.2: **Overview of various approaches to spectrum prediction.** **A.** Fragmentation prediction approaches use heuristics and scoring rules to break down the molecule into fragments and their associated intensities. **B.** Binned prediction approaches discretize the possible mass-to-charge values and predict intensities for each possible bin. **C.** Formula prediction approaches predict spectra as sets of molecular formulae and intensities. Our model SCARF utilizes a two stage approach, first by predicting the product formulae present (constrained by the precursor formula), which defines the x-axis locations of the peaks, before secondly assigning intensities to these formulae (defining the peaks’ y-axis values).

along with the inherent ambiguity in the fragmentation process, limit this strategy’s overall accuracy and speed.

Binned prediction. An increasingly popular and straightforward approach to spectra prediction is to map molecules to discretized 1D mass spectra from either molecular fingerprint [84] or graph inputs [86], [150] (Figure 4.2B). Specifically, these methods divide the m/z axis into fixed-width “bins” and predict an aggregate statistic of the peaks found in each bin (such as their maximum or summed intensity). While more flexible and end-to-end than fragmentation-based approaches, these methods do not impose the same physical constraints or shared information across fragments, making them less interpretable and susceptible to making invalid predictions. Further, discretizing the input spectrum inherently restricts the precision of such models compared to exact-mass predictions.

Formula prediction. We introduce the strategy of predicting spectra at the level of molecular formulae, an intermediate between binned and fragmentation prediction (Figure 4.2C). Simultaneous to our work, two groups have separately explored formula prediction strategies [152], [153]. However, to generate plausible subformulae candidates, they either generate a fixed vocabulary of formulae [152] or restrict their model to molecules under 48 atoms for exhaustive enumeration [153], which is smaller than many compounds of interest. We overcome the combinatorial problem of formula generation using prefix trees, allowing our method to scale and eliminating the need for large, fixed vocabularies.

4.2.3 Mass spectrum libraries

One important use of spectrum predictors is in building large *in silico* libraries of molecule spectra to augment the small size of existing, experimentally derived databases (on the order of 10^4) which are expensive to curate. These spectra libraries are then leveraged downstream in different ways, for example for training molecular property predictors directly

from mass spectra [55]. Another common application of spectra libraries is to infer an unknown molecule’s structure from a newly observed spectra – a particularly hard problem, with only 13% of spectra measured from clinical samples identifiable using current elucidation tools [29]. In this problem, spectra libraries are used as part of a process called *retrieval*: The newly observed spectra is compared with the existing spectra in the library using a fixed or learned spectral distance function, such as cosine distance [81], [82], [100], [154], and the molecules associated with the closest spectra are returned as possible matches. In practice, the retrieval process is constrained to choosing among *isomers* (i.e., molecules with the same molecular formula, and therefore molecular weight, but with different bond configurations) due to the high resolution of modern mass spectrometers (i.e., absolute errors on the order of 10^{-4} to 10^{-3} m/z for MS1 measurements) [50], [52], [111].

Given the varied use cases of spectra libraries, we focus on evaluating spectrum predictors in terms of both (a) their prediction accuracies (Section 4.4.2), using metrics such as “cosine similarity”, and (b) their use in generating virtual spectral libraries to assist with retrieval (Section 4.4.3).

4.3 Model

Here, we describe our model, SCARF, for predicting mass spectra from precursor molecules via first predicting subformulae of the precursor molecule, referred to as *product formulae*. Building upon the notation introduced in the previous section, we continue to denote precursor molecules¹ as $\mathcal{M} \in \mathcal{X}$, and their associated formula vector as $\mathcal{F} \in \mathbb{N}_0^e$, defining at each position, $j \in \{1, \dots, e\}$, the count of each possible chemical element present, \mathcal{F}_j (with zero indicating none of that chemical element is present). Likewise, we define the set of n product formulae as $\{\mathbf{f}^i\}_{i=1}^n$, and associate with each an intensity, y^i . Note that the mass² corresponding to a given formula (and, as such, the x-axis location of the peak on a mass spectrum) is determined deterministically from the counts of each elements present.

At a high level, SCARF generates mass spectra through the composition of two learned functions:

$$\{(\mathbf{f}^i, y^i)\}_{i=1}^n = g_{\theta}^{Weave} \left(g_{\theta}^{Thread}(\mathcal{M}), \mathcal{M} \right), \quad (4.1)$$

first mapping from the original molecule to a set of product formulae, $g_{\theta}^{Thread} : \mathcal{M} \mapsto \{\mathbf{f}^i\}_{i=1}^n$, and then mapping from this set of formulae (and the original molecule) to the respective intensities, $g_{\theta}^{Weave} : (\{\mathbf{f}^i\}_{i=1}^n, \mathcal{M}) \mapsto \{(\mathbf{f}^i, y^i)\}_{i=1}^n$. The particularities of both functions are described in detail below. The specific architectures and hyperparameters used are deferred to the appendix; model code can be found at <https://github.com/samgoldman97/ms-pred>.

¹We model and discuss uncharged molecules and formulae, despite mass spectrometry measuring the masses of adduct *ions*. In practice, we reduce all molecules to uncharged candidates by simply shifting all the spectra weights by the m/z of their respective adducts, which we assume to be equal to the (known) adduct of the parent molecule.

²We assume singly charged adducts (as is common practice, [73]), such that masses and mass-to-charge ratios are interchangeable.

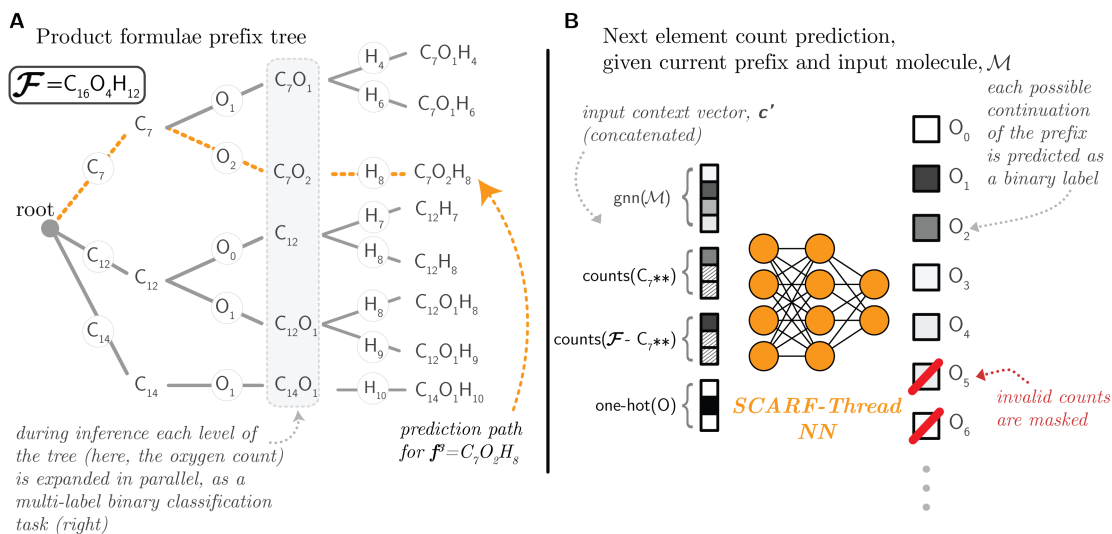


Figure 4.3: **Illustration of the SCARF-Thread architecture.** **A.** The formulae of the product fragments can be represented using a prefix tree. SCARF-Thread predicts this tree for new molecules at test time. It does so by expanding each node at a given depth in parallel, treating the counts of subsequent elements as dependent only on the counts of elements predicted so far (i.e., the prefix) and the original molecular structure. **B.** The SCARF-Thread predictive task at the C_7 node from the prefix tree diagram shown in A. Here the network takes as input (i) an embedding of the overall molecule; (ii) a vector representing the counts of each element in the prefix so far (counts yet to be predicted are represented using a special token), (iii) the difference of the counts predicted so far from the precursor molecule, and (iv) a one-hot representation of the element for which the counts are currently being predicted. The network predicts which counts are valid next nodes in the prefix tree (where counts that are greater than those in the original precursor molecular formula are automatically masked out as invalid). See also Algorithm 9.

4.3.1 SCARF-Thread: Generating product formulae using prefix trees

SCARF-Thread is tasked to learn a mapping to the set of product formulae, $\{\mathbf{f}^i\}_{i=1}^n$, given the original molecule. Naively, one might try to define this model autoregressively, predicting the set formula by formula, chemical element by chemical element. However, such an approach soon runs into a number of problems as (i) the predictions are not invariant to set and ordering permutations; (ii) the time complexity of prediction would scale poorly, being proportional to both the number of elements and number of product formulae (i.e., $\mathcal{O}(e \times n)$); and (iii) the predictions would likely contain duplicates.

We therefore take a different approach using the insight that the set of all product formulae can be compactly represented as a prefix tree (Figure 4.3A). In this tree, edges at a given depth represent valid counts of a particular chemical element, which are often identical across multiple product formulae (shown in the circles). By following each path from the root node to the different leaf nodes, we can reconstruct each product formula (as the orange dashed path does for a single product formula).

We thus propose SCARF-Thread as an autoregressive generator to define a probability distribution over such a prefix tree (Algorithm 9). We assume that each product formula is a subset of the precursor formula, meaning that the precursor formula sets an upper bound on the maximum number of each element³. At each node in the tree (corresponding to a prefix $\mathbf{f}'_{<j}$), we pose the prediction of the set of child nodes (corresponding to the set of valid counts of the subsequent element) as a multi-label binary classification problem (Figure 4.3B). Concretely, we use a neural network module for this task, giving it as input a context vector representing the node being expanded:

$$\mathbf{c}' = [\text{gnn}(\mathcal{M}), \text{counts}(\mathbf{f}'_{<j}), \text{counts}(\mathcal{F} - \mathbf{f}'_{<j}), \text{one-hot}(j)], \quad (4.2)$$

where $\text{gnn}(\cdot)$ specifies a neural encoding of the molecular graph (Section 4.7.5), $\text{counts}(\cdot)$ specifies a count-based encoding of the associated prefix (Section 4.7.5), and $\text{one-hot}(\cdot)$ specifies a one-hot encoding of the node’s depth (or equivalently, which element the predicted count is for). In our experiments, we use a fixed ordering of the chemical elements (Section 4.7.2), but optimizing or even learning the tree construction order could be carried out [155].

Formulae as differences. Following [84], we find it helpful to not only parameterize product formulae in terms of their element counts, but also in terms of the elements that they have lost, i.e., their *difference* from the precursor formula. On the input side, this is already covered by including in the context vector a count-based embedding of the prefix formula minus the product formula ($\text{counts}(\mathcal{F} - \mathbf{f}'_{<j})$). However, on the output side this is achieved by combining the probabilities of a “forward” and a “difference” network:

$$p(f'_j = a | \mathbf{f}'_{<j}, \mathcal{M}) = \alpha_a \sigma(\text{MLP}^F(\mathbf{c}'))_a + (1 - \alpha)_a \sigma(\text{MLP}^D(\mathbf{c}'))_{\mathcal{F}_j - a}, \quad (4.3)$$

where $\text{MLP}^F(\cdot)$ and $\text{MLP}^D(\cdot)$ specify multi-layer perceptrons (MLPs) for predicting the probability of observing a count of a and a loss of $\mathcal{F}_j - a$ atoms respectively; α is a variable (output from a third, unshown network) deciding how to weight these predictions; and $\sigma(\cdot)$ is the element-wise sigmoid function.

4.3.2 SCARF-Weave: Predicting intensities given product formulae

Given the product formulae outputs from SCARF-Thread, SCARF-Weave predicts corresponding intensities at each formula. This is a set-to-set problem, well suited for any equivariant set2set architecture [156, Section 3.1]. In our experiments, we use a Set Transformer [91], [92], which enables the model to consider all the formulae present in the mass spectrum (and their possible interactions) when predicting final intensities.

We choose to represent formula in the set similarly to the context vectors used in SCARF-Thread. For each input, we concatenate a vector embedding of the initial molecular graph with count-based embeddings of the product formula and its difference from the precursor formula (Figure 4.4). We again defer the particularities of the embedding functions to the Appendix (Section 4.7.5).

³While it is possible for fragments to fuse together, potentially taking the count of a chemical element over the number in the original precursor formula, we postpone the extension to modeling such rare events to future work.

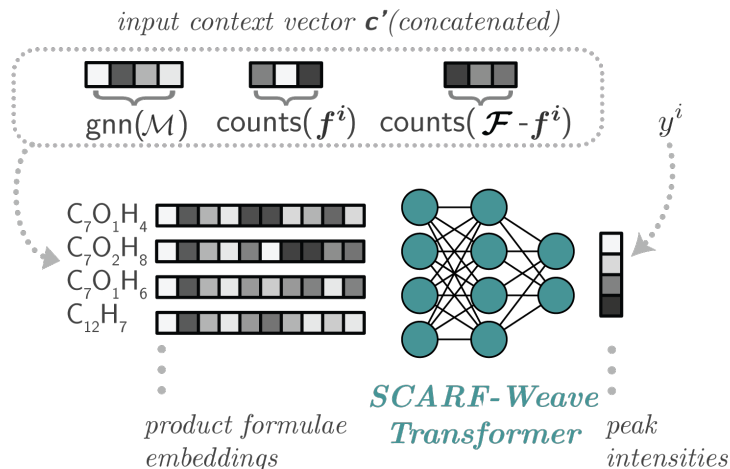


Figure 4.4: The SCARF-Weave network, which takes in the product formulae (e.g., predicted by SCARF-Thread) and predicts their intensities. We use a Set Transformer architecture [91], such that our model takes in the details of the other product formulae present when predicting intensities.

4.3.3 Training and inference

Provided with a dataset of molecules and formula-labeled mass spectra, we could train the two components of SCARF separately. However, in practice we find it beneficial to first train SCARF-Thread and then train SCARF-Weave on its outputs so that the distribution the latter model sees is the same at training and prediction time. SCARF-Weave is trained using a cosine loss (Section 9), as this most closely resembles the “retrieval” setting (Section 4.4.3).

SCARF-Thread is trained using the binary cross entropy losses associated with the multi-label classification tasks at each non-leaf node in the prefix tree. We use teacher forcing, i.e., we train on each level of the tree in parallel by conditioning on the ground-truth set of prefixes at each stage. In our experiments, when generating the set of product formulae from this model we always pick the top 300. Empirically, we find that this provides better performance than picking a variable number based on a likelihood threshold.

4.4 Experiments

We evaluate SCARF on spectra prediction (Section 4.4.2) and molecule identification in a retrieval task (Section 4.4.3).

4.4.1 Dataset

We train and validate SCARF on two libraries: a gold standard commercial tandem mass spectrometry dataset, NIST20 [63], as well as a more heterogeneous public dataset, NPLIB1, extracted from the GNPS database [28] by [57] and subsequently processed by [64]. We prepare both datasets by extracting and preprocessing spectra, as well as filtering to compounds

that (a) are under 1,500 Da (i.e., typically under 100 heavy atoms), (b) only contain predefined elements, and (c) are only charged with common positive-mode adduct types (Section 4.4.1).

Overall, NIST20 contains 35,129 total spectra with 24,403 unique structures, and 12,975 unique molecular formulae; NPLIB1 contains 10,709 spectra, 8,553 unique structures and 5,433 unique molecular formulae. Both datasets are evaluated using a structure-disjoint 90%/10% train/test split with 10% of training data held out for validation, such that all compounds in the test set are not seen in the train and validation sets.

Annotating spectra. We emphasize that SCARF can be trained with any product formula annotations, which can be labeled [63] or inferred with varied computational strategies [50]. Herein, we utilize the MAGMa algorithm [95]. In brief, for a given molecule-spectrum pair in the training dataset, the molecule is combinatorially fragmented at each atom up to a depth of 3 breakages to create sub-fragments. This creates a bank of possible molecular formulae, and each peak in the spectrum is assigned to its nearest possible formula within a mass difference of 20 parts-per-million.

4.4.2 Spectra prediction

Predicting product formulae (SCARF-Thread). SCARF-Thread is trained and used to reconstruct prefix trees and evaluated by its ability to recover the ground truth product formula set. The set of generated product formulae is rank-ordered by the probability of each product formula and filtered to the top k predicted product formulae. The fraction of ground truth formulae (22.29 peaks on average in NIST20) contained in the top k set is computed as *coverage*.

We compare coverage achieved by SCARF-Thread to several baselines: (i) CFM-ID [73], a fragmentation based approach (Section 4.7.3); (ii) a random baseline that samples product formulae from a uniform distribution; (iii) a frequency baseline, which ranks product formulae by the frequency the product formula candidate (or product formula difference) appears in the training set; (iv) an LSTM autoregressive neural network baseline (Section 4.7.3) that is trained to predict molecular formula vectors in sequence from highest to lowest intensity; and two model ablations, (v) SCARF-Thread-D and (vi) SCARF-Thread-F, which only make uni-directional elements difference or forward predictions of element counts respectively (i.e., α in Equation 4.3 is fixed to $\mathbf{0}$ for (v) and $\mathbf{1}$ for (vi)).

In general, SCARF-Thread starkly outperforms all baselines tested (Table 4.1). By generating 300 peaks, SCARF-Thread is able to cover on average 91% and 72% of the true formulae in the ground truth test set for NIST20 and NPLIB1 respectively. Our difference- and forward-only directional prediction ablations demonstrate the benefits of modeling both the atom counts for each element and the differences in counts from the original molecule.

Predicting mass spectra. We next evaluate the strength of SCARF-Weave for intensity prediction on the same test dataset. We compare against five baselines: a fragmentation-based approach, CFM-ID [73]; two NEIMS [84] binned prediction models (Section 4.7.3), using either feed forward network modules (FFNs), as in the original work, or graph neural

Table 4.1: **Model coverage (higher better) of true peak formulae as determined by MAGMa at various max formula cutoffs for the NIST20 and NPLIB1 datasets.** Best result for each column is in bold. Results are computed for a single test set; all re-trained models (i.e., Autoregressive and SCARF variants) are averaged across three random seeds.

| Dataset | NIST20 | | | | NPLIB1 | | | |
|----------------|---------------|---------------|----------------|-----------------|---------------|---------------|----------------|-----------------|
| | Coverage @ 10 | Coverage @ 30 | Coverage @ 300 | Coverage @ 1000 | Coverage @ 10 | Coverage @ 30 | Coverage @ 300 | Coverage @ 1000 |
| Random | 0.009 | 0.026 | 0.232 | 0.532 | 0.004 | 0.014 | 0.126 | 0.336 |
| Frequency | 0.173 | 0.275 | 0.659 | 0.830 | 0.090 | 0.151 | 0.466 | 0.688 |
| CFM-ID | 0.197 | 0.282 | – | – | 0.170 | 0.267 | – | – |
| Autoregressive | 0.204 | 0.262 | 0.309 | 0.317 | 0.072 | 0.082 | 0.095 | 0.099 |
| SCARF-Thread-D | 0.248 | 0.425 | 0.839 | 0.941 | 0.158 | 0.284 | 0.681 | 0.856 |
| SCARF-Thread-F | 0.249 | 0.476 | 0.855 | 0.943 | 0.155 | 0.306 | 0.708 | 0.859 |
| SCARF-Thread | 0.308 | 0.552 | 0.907 | 0.968 | 0.164 | 0.309 | 0.724 | 0.879 |

network modules (GNNs) as in SCARF-Weave and described by [150]; a retrained variant of 3DMolMS ([157]; Section 4.7.3), a binned spectrum predictor that utilizes a point cloud neural network over a single molecular conformer input generated by RDKit [158]; and FixedVocab, a formula prediction model that predicts intensities at a fixed library of formulae and formulae differences inspired by GRAFF-MS ([152]; Section 4.7.3).

To enable fair comparison across models, we predict test spectra at 15k bins (0.1 bin resolution between 0 and 1500) with a maximum of 100 peaks for each predicted molecule. With the exception of CFM-ID, all models are hyperparameter optimized (Section 4.7.5), retrained completely, and conditioned on the same covariate inputs as SCARF; such steps lead to large performance boosts to the prior NEIMS method in particular. We evaluate the quality of our predictions based upon four core criteria reflecting our original desiderata of accuracy, physical-sensibleness, and speed:

1. *Cosine sim.*: Cosine similarity between the ground truth and predicted spectra, indicating spectrum prediction accuracy.
2. *Coverage*: The fraction of ground truth spectrum peaks covered by the predicted spectrum.
3. *Valid*: The fraction of predicted peaks that can be explained by a subformula (that obeys basic ring-double bond equivalent heuristics [131]) of the predicted molecule.
4. *Time (s)*: The wall time it takes (using a single CPU and no batched calculations) to load the model and predict spectra for 100 randomly selected molecules.

SCARF is more accurate than all other approaches on NIST20, improving cosine similarity over a GNN binned prediction approach by over 0.02 points in NIST20 and 0.01 in NPLIB1 (Table 4.2). Further, our method is more physically grounded insofar as all predicted peaks are guaranteed to be valid subformulae, unlike the unconstrained binned

Table 4.2: **Spectra prediction in terms of cosine similarity, coverage (proportion of ground-truth peaks that are covered by the top 100 non-zero predictions), validity (the fraction of predicted peaks for which a chemically plausible explanation is possible), and time.** Best value in each column is typeset in bold (higher is better for all metrics but time). Results are averaged across 3 random seeds on a single data split for all retrainable models (i.e., not CFM-ID).

| Dataset | NIST20 | | | NPLIB1 | | | Time (s) |
|-------------|--------------|--------------|--------------|--------------|--------------|--------------|------------|
| | Cosine sim. | Coverage | Valid | Cosine sim. | Coverage | Valid | |
| CFM-ID | 0.412 | 0.278 | 1.000 | 0.377 | 0.235 | 1.000 | 1114.7 |
| 3DMolMS | 0.510 | 0.734 | 0.945 | 0.394 | 0.507 | 0.919 | 3.5 |
| FixedVocab | 0.704 | 0.788 | 0.997 | 0.568 | 0.563 | 0.998 | 5.5 |
| NEIMS (FFN) | 0.617 | 0.746 | 0.948 | 0.491 | 0.524 | 0.949 | 3.9 |
| NEIMS (GNN) | 0.694 | 0.780 | 0.947 | 0.521 | 0.547 | 0.943 | 4.9 |
| SCARF | 0.726 | 0.807 | 1.000 | 0.536 | 0.552 | 1.000 | 21.1 |

approaches, where nearly 5% of peak predictions cannot be explained by a valid molecular formula. Importantly, SCARF still operates 2 orders of magnitude faster than CFM-ID (Table 4.2).

The heterogeneity, reduced dataset size, and increased average molecular weight (Figure 4.8) of NPLIB1 leads to substantially worse absolute performance across all models. Interestingly, in this setting, the FixedVocab approach [152] performs better, perhaps because the strict priors of formula constraints are more helpful with fewer and more challenging training examples. We further stratify results by molecule size in Figure 4.7, showing that all models are generally more accurate on smaller compounds. We additionally validate that cosine similarity is not merely measuring a model’s ability to predict the parent mass peak by computing a modified cosine similarity with the original molecule’s mass masked (Table 4.9).

4.4.3 Retrieval

A key application for forward spectrum prediction is to use predicted spectra to determine the most plausible molecular structure assignment. We posit forward spectrum prediction models should be particularly helpful in differentiating structurally similar molecules and design a retrieval task to showcase such potential. For each test set molecule, we extract 49 potential “decoy” options based upon the most structurally similar *isomers* (i.e., compounds with the same precursor formula) within PubChem [128] as judged by Tanimoto similarity using Morgan fingerprints. While retrieval could be conducted on the entirety of PubChem or other similarly large molecular databases, we believe this subset retrieval setting is more practical and better mirrors a real-world setting (see Section 4.7.4 for justification). We predict the spectra for all molecules and rank them according to their similarity to the ground truth spectrum, computing the *accuracy* for retrieval. Herein, we specifically emphasize models and retrieval on the NIST20 dataset, as it is a much larger and higher quality dataset.

SCARF reaches a top 1 and top 5 retrieval accuracy in this task of 18.7% and 54.1%

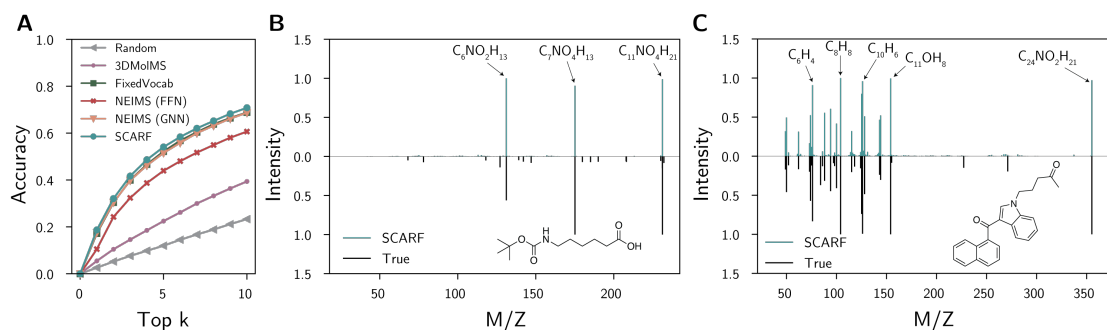


Figure 4.5: **SCARF enables more accurate retrieval of ground truth molecules within the NIST20 dataset.** **A.** Average retrieval accuracy of SCARF at various top k thresholds. Retrieval is conducted on the same test split, and retrieval accuracy is averaged across models trained for three separate random seeds. **B-C.** Example spectrum predictions made by SCARF (top) compared to the ground truth spectrum (bottom). Up to 5 predicted peaks with the highest intensity are annotated with their molecular formula explanation as predicted by SCARF. The full molecule is shown inset. Further examples are in the Appendix (Figure 4.6).

respectively, representing an improvement over the methods with the second best top 1 accuracy of 17.5% (NEIMS (GNN)) and top 5 accuracy of 52.2% (Fixed Vocab) (Figure 4.5A). We highlight two example predictions from SCARF (Figure 4.5B-C), with additional randomly sampled test set predictions shown in Figure 4.7.1. We repeat similar experiments within NPLIB1, but find that cosine similarity performance is uncorrelated with relative ranking performance; feed forward fingerprint based approaches are better at retrieval, despite relatively weak cosine similarity (Section 4.7.1). FixedVocab [152] performs especially well on NPLIB1, again likely due to the helpful biases imparted by constraining the formula vocabulary.

This result underscores previous observations regarding how database and model biases can skew retrieval results under certain settings [159]. That is, models may be more or less robust for certain classes of molecules, so the composition of these classes in the retrieval library may affect the retrieval accuracy accordingly. The observed discrepancy between cosine similarity and retrieval performance can further be explained by the dataset shift required for computing retrieval accuracy; cosine similarity is evaluated on “in-domain” data, whereas retrieval relies also on accuracy on unlabeled data that may be “out-of-domain.”

4.5 Related Work

Forward vs backward models. Computational tools to identify mass spectra are often divided into two categories: forward and backward models. Forward models, i.e., spectrum predictors, such as SCARF or the methods discussed in Section 4.2, operate in the causal direction and try to predict the spectrum given the molecule. Backward models start from the spectrum and predict features or even full molecule structures. Early backward models used heuristics, expert rules, and even neural networks [160]–[162]. Such approaches have

more recently been augmented with kernel methods and more modern, deep representation learning techniques [53], [56], [64], [83]. These models are complementary to spectrum predictors.

Mass spectra for proteomics. Although this paper has focused on small molecules, similar trends of deep representation learning for mass spectra are also emerging in the adjacent field of proteomics [163], with [164] recently proposing a benchmark challenge in this domain. While small molecule and protein spectra are similar, proteomics spectra tend to be more easily predicted as fragments are often formed at peptide bonds. We believe adapting SCARF to this task would be an interesting direction for future work.

Neural set generation. Our work is also related to methods for modeling sets and multisets. SCARF-Thread generates a set as output, which has been studied elsewhere in the context of n-gram generation [155], object detection [165], [166], and point cloud generation [167]. The product formulae sets we generate, however, are different to those considered in these other task; in our setting, each member of the set (i.e., individual product formula) represents a multiset of atom types (i.e., multiple carbons, multiple hydrogens, etc.) and is constrained physically by the precursor formula.

4.6 Conclusion

In this paper we introduced SCARF, an approach utilizing prefix tree data structures to efficiently decode mass spectra from molecules. By first predicting product formulae and then assigning these formulae intensities, we are able to combine the advantages of previous neural and fragment based approaches, providing fast and physically grounded predictions. We show how these resulting predictions are both more accurate in predicting experimentally-observed spectra and yield improvements in identifying a molecule’s structure from its respective mass spectrum, as tested on a widely used dataset.

In term of limitations, our model is data dependent, as indicated by the relative performance across the NIST20 and NPLIB1 datasets. SCARF is also highly reliant upon the quality of product formula label assignment. The current commercial status of mass spectrometry training data poses a barrier to entry, and identifying high quality public domain data is critical for future studies. A key contribution of this work is to retrain and optimize the hyperparameters of competing methods on a publicly available dataset under equivalent conditions to allow for future extensions. Directly training on a ranking-based loss or learning a model specific spectra distance function may be one way to improve upon our model’s performance in the retrieval setting, and we outline additional potential ideas more explicitly in Section 4.7.6.

Future directions will involve further developing SCARF for real world use cases such as unknown metabolite elucidation in clinical samples. Specific directions will include more carefully modeling covariates (e.g., collision energies and MS/MS instrument types), grounding product formulae in molecular graph substructures, and utilizing such models to augment inverse spectrum-to-molecule annotation tools.

4.7 Additional Results

4.7.1 Extended results

We benchmark models in terms of retrieval accuracy as described (Section 4.4.3) for both the NIST20 and NPLIB1 datasets (Table 4.3, 4.4). We recreate retrieval experiments using the full PubChem retrieval library in Table 4.5. We reproduce main text results with standard error values included in Tables 4.6, 4.7, and 4.8. We showcase additional spectra predictions from our model trained on NIST20 in Figure 4.6.

Table 4.3: **NIST20 spectra prediction retrieval top k accuracy for different values of k**. All values represent the mean across three separate random seeds \pm the standard error of the mean for a single test set.

| top k | 1 | 2 | 3 | 4 | 5 | 8 | 10 |
|-------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|
| Random | 0.026 \pm 0.000 | 0.052 \pm 0.001 | 0.076 \pm 0.001 | 0.098 \pm 0.001 | 0.120 \pm 0.000 | 0.189 \pm 0.001 | 0.233 \pm 0.002 |
| 3DMolMS | 0.055 \pm 0.002 | 0.105 \pm 0.000 | 0.146 \pm 0.002 | 0.185 \pm 0.004 | 0.225 \pm 0.004 | 0.332 \pm 0.003 | 0.394 \pm 0.004 |
| FixedVocab | 0.172 \pm 0.002 | 0.304 \pm 0.002 | 0.399 \pm 0.001 | 0.466 \pm 0.004 | 0.522 \pm 0.006 | 0.638 \pm 0.005 | 0.688 \pm 0.003 |
| NEIMS (FFN) | 0.105 \pm 0.002 | 0.243 \pm 0.006 | 0.324 \pm 0.006 | 0.387 \pm 0.006 | 0.440 \pm 0.007 | 0.549 \pm 0.005 | 0.607 \pm 0.002 |
| NEIMS (GNN) | 0.175 \pm 0.003 | 0.305 \pm 0.001 | 0.398 \pm 0.001 | 0.462 \pm 0.002 | 0.515 \pm 0.003 | 0.632 \pm 0.003 | 0.687 \pm 0.003 |
| SCARF | 0.187 \pm 0.004 | 0.321 \pm 0.006 | 0.417 \pm 0.004 | 0.486 \pm 0.004 | 0.541 \pm 0.005 | 0.652 \pm 0.004 | 0.708 \pm 0.005 |

Table 4.4: **NPLIB1 spectra prediction retrieval top k accuracy for different values of k**. All values represent the mean across three separate random seeds \pm the standard error of the mean for a single test set.

| top k | 1 | 2 | 3 | 4 | 5 | 8 | 10 |
|-------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|
| Random | 0.033 \pm 0.001 | 0.061 \pm 0.005 | 0.092 \pm 0.003 | 0.118 \pm 0.003 | 0.141 \pm 0.006 | 0.216 \pm 0.006 | 0.258 \pm 0.006 |
| 3DMolMS | 0.087 \pm 0.001 | 0.159 \pm 0.010 | 0.218 \pm 0.004 | 0.268 \pm 0.006 | 0.317 \pm 0.006 | 0.427 \pm 0.008 | 0.488 \pm 0.005 |
| FixedVocab | 0.193 \pm 0.003 | 0.314 \pm 0.004 | 0.390 \pm 0.003 | 0.448 \pm 0.005 | 0.492 \pm 0.001 | 0.587 \pm 0.005 | 0.635 \pm 0.006 |
| NEIMS (FFN) | 0.195 \pm 0.003 | 0.313 \pm 0.002 | 0.388 \pm 0.003 | 0.447 \pm 0.006 | 0.488 \pm 0.002 | 0.585 \pm 0.007 | 0.624 \pm 0.010 |
| NEIMS (GNN) | 0.174 \pm 0.007 | 0.285 \pm 0.004 | 0.362 \pm 0.002 | 0.422 \pm 0.001 | 0.471 \pm 0.002 | 0.586 \pm 0.007 | 0.640 \pm 0.005 |
| SCARF | 0.135 \pm 0.007 | 0.242 \pm 0.001 | 0.320 \pm 0.001 | 0.389 \pm 0.004 | 0.444 \pm 0.002 | 0.569 \pm 0.001 | 0.630 \pm 0.008 |

Table 4.5: **Retrieval accuracy on NIST20 for a single 500 molecule subset of the test set using a library of 49 decoys or all decoys contained in PubChem (“None”)**. Results were repeated for 3 random training seeds of the model and are shown \pm the standard error of the mean. The top value is shown in bold.

| PubChem limit | 50 | | | None | | |
|---------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|
| Top k | 1 | 2 | 3 | 1 | 2 | 3 |
| FixedVocab | 0.168 \pm 0.003 | 0.308 \pm 0.003 | 0.410 \pm 0.005 | 0.145 \pm 0.004 | 0.258 \pm 0.004 | 0.323 \pm 0.001 |
| NEIMS (FFN) | 0.102 \pm 0.002 | 0.237 \pm 0.003 | 0.315 \pm 0.004 | 0.087 \pm 0.003 | 0.183 \pm 0.008 | 0.236 \pm 0.007 |
| NEIMS (GNN) | 0.169 \pm 0.003 | 0.300 \pm 0.004 | 0.402 \pm 0.005 | 0.138 \pm 0.004 | 0.239 \pm 0.005 | 0.312 \pm 0.008 |
| SCARF | 0.204 \pm 0.009 | 0.326 \pm 0.005 | 0.432 \pm 0.009 | 0.167 \pm 0.003 | 0.258 \pm 0.001 | 0.336 \pm 0.003 |

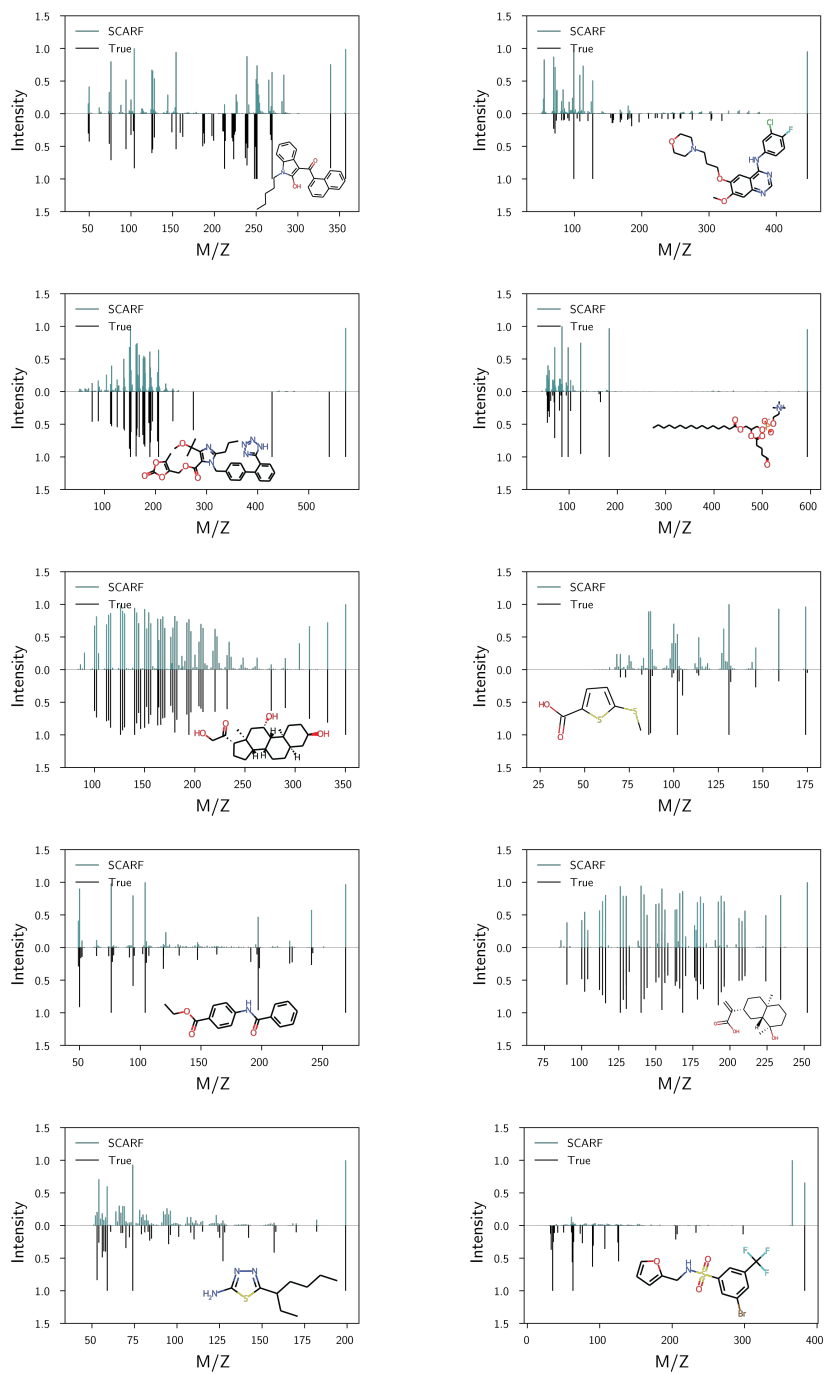


Figure 4.6: **Example spectra predictions from the NIST20 dataset for 10 randomly selected test molecules.** The ground truth spectra are shown underneath in black, with predictions above in teal. Molecules are shown inset.

Table 4.6: **Model coverage of true peak formulae as determined by MAGMa at various max formula cutoffs for the NPLIB1 dataset.** Results are calculated for a single held out test split, shown \pm the standard error of the mean across three random seeds for all models that were retrained. The best value (i.e., highest) is typeset in bold for each column.

| Coverage @ | 10 | 30 | 300 | 1000 |
|----------------|-------------------|-------------------------------------|-------------------------------------|-------------------------------------|
| Random | 0.004 | 0.014 | 0.126 | 0.336 |
| Frequency | 0.090 | 0.151 | 0.466 | 0.688 |
| CFM-ID | 0.170 | 0.267 | – | – |
| Autoregressive | 0.072 \pm 0.001 | 0.082 \pm 0.002 | 0.095 \pm 0.001 | 0.099 \pm 0.000 |
| SCARF-D | 0.158 \pm 0.001 | 0.284 \pm 0.003 | 0.681 \pm 0.002 | 0.856 \pm 0.002 |
| SCARF-F | 0.155 \pm 0.002 | 0.306 \pm 0.003 | 0.708 \pm 0.003 | 0.859 \pm 0.001 |
| SCARF | 0.164 \pm 0.009 | 0.309 \pm 0.014 | 0.724 \pm 0.013 | 0.879 \pm 0.004 |

4.7.2 Dataset preparation

NIST20 [63] is prepared by extracting all positive-mode experimental spectra collected in higher-energy collision-induced dissociation (HCD) mode (i.e., collected on Orbitrap mass spectrometers). Spectra are filtered, so that we keep only those for which the associated molecule (M) has (i) a mass under 1,500 Da, (ii) contains only elements from a predefined set (i.e., C, N, P, O, S, Si, I, H, Cl, F, Br, B, Se, Fe, Co, As, Na, and K), and (iii) is charged with common adduct types (i.e., $[M+H]^+$, $[M+Na]^+$, $[M+K]^+$, $[M-H_2O+H]^+$, $[M+NH_3+H]^+$, and $[M-2H_2O+H]^+$). Because non-standard empirical spectra databases [28] often do not include the measured collision energies, we pool all collision energies for each compound-adduct pairing to create a single spectrum. We refer the reader to Young et al. [86] for detailed instructions for purchasing and extracting the NIST20 dataset.

All spectrum intensities are square-root transformed to provide higher weighting to lower intensity peaks, normalized to a maximum intensity of 1 (i.e., through dividing by the maximum intensity), filtered to exclude any noise peaks with normalized intensity under 0.003, and subsetted to only the top 50 highest intensity peaks. All peaks are mass-shifted by the weight of the parent adduct (i.e., if the spectrum is $[M+H]^+$, the weight of a proton is subtracted from each child peak).

Product formulae assignments

Because the precursor ion and adduct species are known for the training dataset, we subtract the precursor adduct mass from every peak in the training set, and attempt to annotate each peak with a plausible product formula (i.e., a subset of the true precursor formula).

We opt to constrain the training product formulae to be subsets of contiguous heavy atoms of the parent molecule as derived with the MAGMa algorithm [95].

We note two important limitations of these heuristics. First, by using molecular substructures to annotate product formulae, our model is less prone to correctly identifying

Table 4.7: **Model coverage of true peak formulae as determined by MAGMa at various max formula cutoffs for the NIST20 dataset.** Results are calculated for a single held out test split, shown \pm the standard error of the mean across three random seeds for all models that were retrained. The best value (i.e., highest) is typeset in bold for each column.

| Coverage @ | 10 | 30 | 300 | 1000 |
|----------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|
| Random | 0.009 | 0.026 | 0.232 | 0.532 |
| Frequency | 0.173 | 0.275 | 0.659 | 0.830 |
| CFM-ID | 0.197 | 0.282 | – | – |
| Autoregressive | 0.204 \pm 0.001 | 0.262 \pm 0.002 | 0.309 \pm 0.005 | 0.317 \pm 0.006 |
| SCARF-D | 0.248 \pm 0.001 | 0.425 \pm 0.002 | 0.839 \pm 0.002 | 0.941 \pm 0.001 |
| SCARF-F | 0.249 \pm 0.001 | 0.476 \pm 0.002 | 0.855 \pm 0.000 | 0.943 \pm 0.001 |
| SCARF | 0.308 \pm 0.002 | 0.552 \pm 0.001 | 0.907 \pm 0.002 | 0.968 \pm 0.001 |

complex rearrangements. Second, it is also possible for adduct switching to occur. Namely, if the precursor ion has a sodium adduct ($[M+Na]^+$), some of the product formulae may actually switch and acquire a hydrogen adduct instead. We assume no adduct switching in our formulation, instead focusing on the novelty of the prefix tree decoding approach, as these represent data labeling challenges, rather than modeling challenges.

In addition, any predictive models of product formulae distributions will more closely predict spectra that would be produced on instrumentation similar to the training sets utilized [148], [149]. Given this, we encourage users of such models to treat these predictions as putative, rather than experimentally valid.

Dataset statistics

To probe the composition of our two primary datasets, we investigate both the molecular weight and chemical classes contained in the NPLIB1 and NIST20 datasets. We find that the average molecular weight is higher for NPLIB1 (Figure 4.8A), consistent with the increased complexity of natural product molecules. We additionally compute chemical classes of the compounds using NPClassifier [101] to identify the types of compounds present in both datasets (Figure 4.8B-C). NPLIB1 is enriched for steroids, coumarins, and various complex alkaloid natural products. On the other hand, NIST20 is enriched for small peptides, nicotinic acid alkaloids, and pseudoalkaloids, among others. While these descriptions are helpful to identify the dataset composition, chemical compound classification is itself a learned classification and should be interpreted cautiously.

4.7.3 Baselines

We further describe select baseline models.

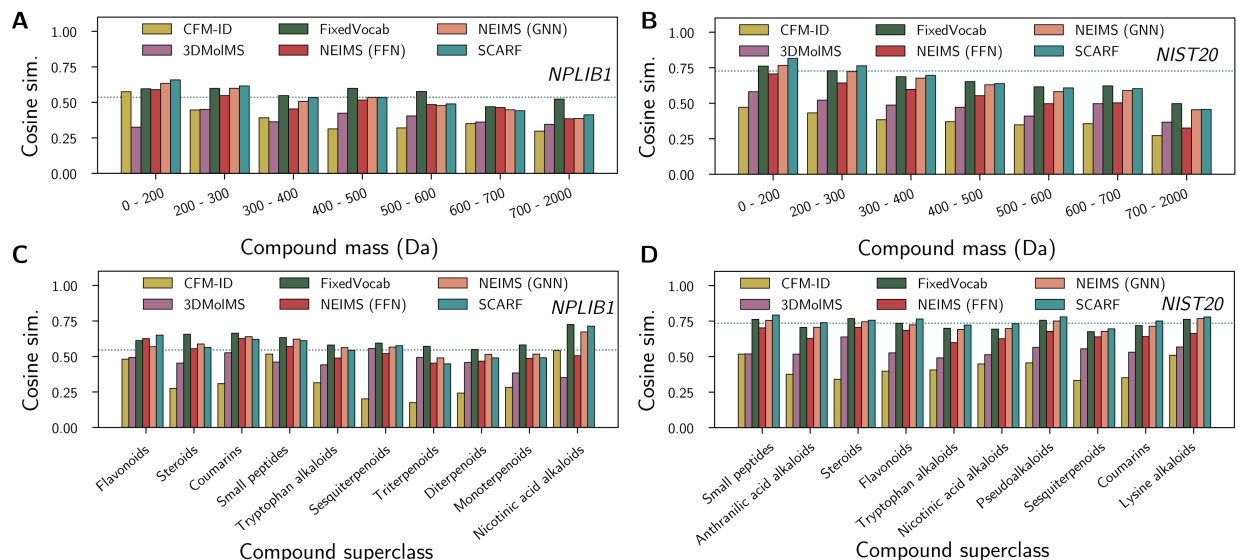


Figure 4.7: **Cosine similarity of predicted spectra stratified by properties.** Results are stratified across molecular weight for both NPLIB1 (A) and NIST20 (B). We further stratify results across putative chemical classes of input molecules using NPClassifier [101] for both NPLIB1 (C) and NIST20 (D). The dotted line indicates the average predictive cosine similarity of SCARF across all examples and averaged over three random splits.

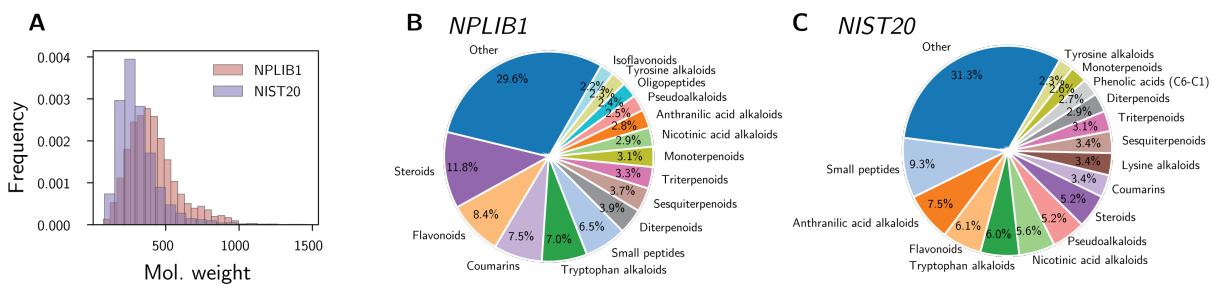


Figure 4.8: **Spectra dataset molecule characterizations.** A. Distribution of the molecular weight of compounds across NPLIB1 and NIST20. B-C. Chemical classes contained in NPLIB1 (B) and NIST20 (C) with the top 15 classes shown and all others grouped in 'Other'. Chemical classes are computed using NPClassifier [101].

Table 4.8: **Spectra prediction in terms of cosine similarity, coverage (proportion of ground-truth peaks that are covered by the top 100 non-zero predictions), validity (the fraction of predicted peaks for which a chemically plausible explanation is possible), and time.** Best value in each column is typeset in bold (higher is better for all metrics but time). Values are shown \pm the standard error of the mean computed across three random seeds on a single test set for all models that could be retrained (i.e., not CFM-ID).

| Dataset | NIST20 | | | NPLIB1 | | | Time (s) |
|-------------|-------------------------------------|-------------------------------------|------------------------------------|-------------------------------------|-------------------------------------|------------------------------------|------------|
| | Cosine sim. | Coverage | Valid | Cosine sim. | Coverage | Valid | |
| CFM-ID | 0.412 \pm 0.000 | 0.278 \pm 0.000 | 1.00 \pm 0.000 | 0.377 \pm 0.000 | 0.235 \pm 0.000 | 1.00 \pm 0.000 | 1114.7 |
| 3DMolMS | 0.510 \pm 0.000 | 0.734 \pm 0.001 | 0.94 \pm 0.001 | 0.394 \pm 0.002 | 0.507 \pm 0.001 | 0.92 \pm 0.000 | 3.5 |
| FixedVocab | 0.704 \pm 0.000 | 0.788 \pm 0.000 | 1.00 \pm 0.000 | 0.568 \pm 0.002 | 0.563 \pm 0.001 | 1.00 \pm 0.000 | 5.5 |
| NEIMS (FFN) | 0.617 \pm 0.000 | 0.746 \pm 0.001 | 0.95 \pm 0.001 | 0.491 \pm 0.002 | 0.524 \pm 0.001 | 0.95 \pm 0.000 | 3.9 |
| NEIMS (GNN) | 0.694 \pm 0.000 | 0.780 \pm 0.000 | 0.95 \pm 0.001 | 0.521 \pm 0.002 | 0.547 \pm 0.003 | 0.94 \pm 0.001 | 4.9 |
| SCARF | 0.726 \pm 0.001 | 0.807 \pm 0.000 | 1.00 \pm 0.000 | 0.536 \pm 0.007 | 0.552 \pm 0.008 | 1.00 \pm 0.000 | 21.1 |

Table 4.9: **Spectra prediction accuracy comparing inclusion (Cosine sim.) and exclusion (Cosine sim. (no MS1)) of the precursor mass.** For all compounds, the peak at the mass of the input compound is masked in the prediction and ground truth to compute Cosine sim. (no MS1). All results represent an average on a single test set across three random seeds.

| Dataset | NIST20 | | NPLIB1 | |
|-------------|--------------|----------------------|--------------|----------------------|
| | Cosine sim. | Cosine sim. (no MS1) | Cosine sim. | Cosine sim. (no MS1) |
| CFM-ID | 0.412 | 0.289 | 0.377 | 0.326 |
| 3DMolMS | 0.510 | 0.517 | 0.394 | 0.390 |
| FixedVocab | 0.704 | 0.637 | 0.568 | 0.505 |
| NEIMS (FFN) | 0.617 | 0.557 | 0.491 | 0.454 |
| NEIMS (GNN) | 0.694 | 0.620 | 0.521 | 0.477 |
| SCARF | 0.726 | 0.663 | 0.536 | 0.466 |

CFM-ID baseline

CFM-ID [73] is a long standing and important approach to fragmentation prediction. Because CFM-ID is fit using a time intensive EM training approach on an analogous dataset, we utilize the pretrained Docker implementation provided by the authors in line with [152]. CFM-ID has two options for predicting molecules in either positive or negative adduct mode with hydrogen adducts (i.e., $[M+H]^+$ or $[M-H]^-$). To directly compare to our method, we predict spectra in positive mode and remove hydrogens from all predicted peaks, as all training peaks are shifted by removing their adducts.

CFM-ID also produces predictions at three collision energies (i.e., low, medium, or high fragmentation). Because we opt not to include these, we merge these predictions and re-normalize the result to a maximum of 1.

Table 4.10: NIST20 retrieval accuracy averaged across three random seeds of model training stratified by the weight of the target molecule.

| Dataset | NIST20 | | | | | | |
|------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Molecular weight | 0 - 200 | 200 - 300 | 300 - 400 | 400 - 500 | 500 - 600 | 600 - 700 | 700 - 2000 |
| Num. compounds | 624 | 1358 | 900 | 387 | 129 | 56 | 89 |
| Random | 0.024 | 0.021 | 0.027 | 0.025 | 0.031 | 0.048 | 0.101 |
| 3DMolMS | 0.039 | 0.041 | 0.057 | 0.077 | 0.070 | 0.125 | 0.199 |
| FixedVocab | 0.143 | 0.184 | 0.168 | 0.172 | 0.196 | 0.220 | 0.161 |
| NEIMS (FFN) | 0.110 | 0.122 | 0.092 | 0.078 | 0.111 | 0.083 | 0.064 |
| NEIMS (GNN) | 0.155 | 0.192 | 0.164 | 0.182 | 0.147 | 0.214 | 0.161 |
| SCARF | 0.191 | 0.211 | 0.165 | 0.163 | 0.168 | 0.161 | 0.169 |

Autoregressive baseline

When considering the task of generating spectrum formulae candidates, we compare SCARF-Thread to an autoregressive recurrent neural network baseline, which is based around a long short term memory (LSTM) module [168].

The LSTM generates formulae consecutively from a single concatenated encoding of the input molecule and input full formula. At each step in the recurrent process, a one-hot encoding of the previous predicted element count is concatenated to a one-hot encoding of the element type being predicted in the current step. By embedding this information into the network, we can avoid predicting element types that do not appear in the parent molecule’s molecular formula. If the parent molecular formula has 5 element types, each autoregressively predicted formula requires generating only 5 element type counts; this eliminates the need for a stop condition between each formula. Formulae are generated autoregressively, from highest to lowest intensity. When predicting the counts of the next element type, we employ the same difference and forward count prediction strategy as used in SCARF for fair comparison. The model is trained with a cross entropy loss and full hyperparameters are listed in Table 4.12.

NEIMS baseline

NEIMS [84] is a highly efficient binned spectrum prediction approach, originally developed for gas chromatography-mass spectra (GC-MS). To enable a fair comparison, we optimize its hyperparameters on our dataset and add higher resolution bins. Furthermore, we also train a graph neural network-based version “NEIMS (GNN)” (in addition to the network that more closely matches [84]’s original model and operates on the molecular fingerprint, “NEIMS (FFN)”). The adduct type is either concatenated to all atom features—as for NEIMS (GNN)—or to the fingerprint vector—as for NEIMS (FFN).

3DMolMS baseline

3DMolMS [157] is a binned spectrum prediction approach developed simultaneously to this work. Unlike the other binned NEIMS approach, 3DMolMS utilizes a point-based deep

neural network model operating on the point cloud of an input molecule. To project a 2D molecule or SMILES string into 3D space, a single 3D conformer is first generated using RDKit [158]. After several 3D convolutions, the atom-wise representations are pooled, covariates corresponding to the settings of the machine and experiment are concatenated, and the result is projected into a fixed length binned spectrum.

To enable a fair comparison, we copy the 3DMolMS architecture into our modeling framework with minor tweaks to the network. Rather than use variable sized hidden layers, we fix hidden layer sizes to a single value across convolutions. In addition, we only use the covariate of the adduct type for consistency with our model, excluding collision energy and instrument type. We hyperparameter optimize the model independently. We find that the performance of this model is substantially lower than the NEIMS baseline, likely due to the additional use of the “difference” prediction module in the NEIMS approach that allows the network to predict intensities at both fragments and neutral losses.

FixedVocab baseline

Concurrently to our work, Murphy, Jegelka, Fraenkel, *et al.* introduce an alternative formula prediction strategy for mass spectrum prediction, a model they term GRAFF-MS [152]. Unlike SCARF, GRAFF-MS utilizes a fixed vocabulary of molecular formulae and molecular formula differences, predicting intensities at each such value without learning to encode the formulae. These formulae and formula differences are selected in a greedy fashion based upon their frequency in the training set.

Because no code was released for this approach at the time this work was conducted, we reimplement a variant of their approach that emphasizes the use of a fixed vocabulary of formulae and differences. For methodological consistency, we utilize equivalent formula annotations as used by SCARF (i.e., one annotation per peak), do not model collision energies or instrument types, and utilize the same graph encoder as SCARF for encoding each molecule. We treat the number of fragment and difference formulae as a tunable hyperparameter (which is optimized along with the rest of the hyperparameters – see Section 4.7.5). We mask all invalid formulae and differences and utilize a cosine similarity loss with the original spectrum to train the model. To convert predicted formula and difference intensities into a binned spectrum, each formula-intensity pairing is projected into its respective binned position using a scatter max calculation. We note that because alternative adducts and isotopes are not labeled in our preprocessing step, we do not predict isotopic or adduct variants for each fragment.

Given the differences between codebases, it is possible that the performance of our reimplement does not exactly match the original implementation, and we instead refer to it as “FixedVocab” rather than GRAFF-MS in table presentations. An earlier version of our work understated the FixedVocab model’s performance due to an implementation decision along these lines (specifically, not including the “0” neutral loss as a predicted vocabulary entry). This has since been rectified, increasing the accuracy of the FixedVocab model.

4.7.4 Retrieval subsets vs. PubChem

We restrict our retrieval experiments in Section 4.4.3 to only the top 49 decoys per test case for two reasons. First, from a practical perspective, running the forward model on every isomer match in PubChem (approx. thousands each, >1,000,000 for only 500 test cases) makes benchmarking across all considered models substantially more challenging both for this work and also future work. Second, we also believe that this top 50 challenge represents a more realistic setting. In practice, retrieval compound libraries will often be carefully crafted and designed to contain molecules similar to the unknown molecule rather than all possible isomers (using either prior knowledge or “backward” models such as CSI:FingerID [53] and MIST [64]). We conduct a side-by-side analysis on a small 500 molecule subset of the test set comparing the setting described above (with 49 decoys) to a setting with no limit on the number of decoys. The results are shown in Table 4.5, showing that SCARF still performs well in this setting.

4.7.5 Model details

Here, we describe details of our model’s training setup, architecture, and hyperparameters that were omitted from the main text. Definitive details can also be found in the code at <https://github.com/samgoldman97/ms-pred>.

Training

We train each of our models on a single RTX A5000 NVIDIA GPU (CUDA Version 11.6), making use of the Torch Lightning [140] library to manage the training. SCARF-Thread and SCARF-Weave take on the order of 1.5 and 2.5 hours of wall time to train respectively.

Molecule encoding

Within both SCARF-Thread and SCARF-Weave, a key component is an encoding of the molecular graph using a message passing graph neural network, $\text{gnn}(\mathcal{M})$. Such graph neural network models are now well described [169]–[171], so we will skip a detailed explanation of them here. In our experiments, we use gated graph sequence neural networks [172]. We made use of the implementation of this network in the DGL library [173] and use as atom features those shown in Table 4.11 (which are computed using RDKit [158] or DGL [173]).

Molecular formulae representations

When forming representations of formulae (including formulae prefixes) we use a count-based encoder, $\text{counts}(\mathbf{f})$. This encoder takes in as input the counts of all individual elements in the formula (which also can be “undefined” for counts of elements not yet specified – indicated as ‘*’ in Figure 4.3B) and returns a vector representation in \mathbb{R}^d . The encoder is based upon the Fourier feature mapping proposed by [139], but using only sin basis functions (to reduce the number of parameters required by our networks). [139] has shown that such features perform better than encoding integers directly; furthermore, compared to learned

Table 4.11: **Graph neural network (GNN) atom features.**

| Name | Description |
|-------------------------|---|
| Element type | one-hot encoding of the element type |
| Degree | one-hot encoding of number of bonds atom is associated with |
| Hybridization type | one-hot encoding of the hybridization (SP, SP2, SP3, SP3D, SP3D2) |
| Charge | one-hot encoding of atom’s formal charge (from -2 to 3) |
| Ring-system | binary flag indicating whether atom is part of a ring |
| Atom mass | atom’s mass as a <code>float</code> |
| Chiral tag | atom’s chiral tag as one-hot encoding |
| Adduct type | one-hot encoding of the adduct ion |
| Random walk embed steps | positional encodings of the nodes computed using DGL |

representations, using Fourier features enables us (at least in principle) to deal with counts at test time that have not been seen during training.

To be precise, each possible count, $v \in \mathbb{N}_0$, is encoded by our counts-based encoder into the vector:

$$\text{abs} \left(\left[\sin \left(\frac{2\pi v}{T_1} \right), \sin \left(\frac{2\pi v}{T_2} \right), \sin \left(\frac{2\pi v}{T_3} \right), \dots \right] \right),$$

where the periods (T_1 , T_2 , etc.) are set at increasing powers of two that enable us to discriminate all possible element counts given in the input, and $\text{abs}(\cdot)$ is the absolute value function such that we get positive embeddings. For the “undefined” count we learn a separate encoding of the same dimensionality.

Further details of SCARF-Thread

Pseudo-code for the SCARF-Thread model is shown in Algorithm 9. Note that the second for loop (on the line marked †) does not depend on previous iterations of the loop, so that in practice we perform this computation in parallel. At training time we use teacher forcing (Section 4.3.3), meaning the first for loop (marked †) is only run sequentially at inference time.

The function `scarf-thread-net(\cdot)` represents the network shown in Figure 4.3B and generates the set of subsequent valid element counts given a prefix (i.e., the child nodes of a given prefix node). As discussed in the main text, we treat this as a multi-label binary classification task and predict the binary label for each possible count using forward and difference MLPs (Equation 4.3). We fix a maximum possible element count (i.e., the number of possible classes in this classification problem), $N = 160$. We do not allow product formulae to have more of a given element than is present in the precursor formula, \mathcal{F} , and we achieve this by setting the probability of these classes to zero.

Further details of SCARF-Weave

As discussed in the main text, SCARF-Weave is based off [91]’s Set Transformer. After forming the input encoding using the molecule and count-based encoder (Section 4.7.5 & Section 4.7.5), we further refine this embedding using an MLP (multi-layer perceptron)

Algorithm 4.7.1: Pseudo-code for SCARF-Thread, which generates prefix trees from a root node autoregressively, one level at a time.

Data: Input molecule, \mathcal{M} , with corresponding input formula, \mathcal{F} .
Result: Set of product formulae, $\rho_e = \{\mathbf{f}^i\}_{i=1}^n$.

```
1  $\mathbf{h}_{\mathcal{M}} \leftarrow \text{gnn}(\mathcal{M})$  ; ▷ Form embedding of precursor molecule.
2  $\rho_0 \leftarrow \{*\}$  ; ▷ Store the set of initial prefixes which is just the undefined formula, *.
3 † for  $j \in [1, \dots, e]$  do ▷ Loop over all possible elements.
4    $\rho_j \leftarrow \{\}$  ; ▷ Create the set of prefixes the next time around.
5    $\mathbf{h}_j \leftarrow \text{one-hot}(j)$  ; ▷ Encoding of which element we are predicting the count of.
6 ‡ for  $\mathbf{f}'_{<j} \in \rho_{j-1}$  do (in parallel) ▷ Loop over all current prefixes.
7    $\mathbf{c}' = [\mathbf{h}_{\mathcal{M}}, \text{counts}(\mathbf{f}'_{<j}), \text{counts}(\mathcal{F} - \mathbf{f}'_{<j}), \mathbf{h}_j]$  ; ▷ Create context vector, Equation 4.2
8    $\{f_j^{i'}\}_{i'=1}^{n'} \leftarrow \text{scarf-thread-net}(\mathbf{c}', \mathcal{F})$  ; ▷ Predict the set of valid next element counts under this
   prefix.
9    $\rho_j \leftarrow \rho_j \cup \text{create-new-prefixes}(\mathbf{f}'_{<j}, \{f_j^{i'}\}_{i'=1}^{n'})$  ; ▷ Create new prefixes for the next element.
return  $\rho_e$ 
```

network. The output of this is passed into a series of l_3 Transformer [92] layers (Section 4.7.5 defines the exact number used in the experiments) with 8 attention heads each.

We use a cosine distance loss to train the parameters of SCARF-Weave. This loss is also used for the FFN and GNN baselines (Table 4.2). To ensure consistency with the baselines, we first project the output of SCARF-Weave into a binned histogram representation (Section 4.7.5 defines the number of bins used); for each bin we take the max intensity across all applicable formulae. Given a predicted binned spectra, $\hat{\mathbf{s}}$, and the ground-truth binned spectra, \mathbf{s} , the cosine distance is defined as the negative of the cosine similarity (computed using PyTorch’s `torch.cosine_similarity` function [174]):

$$\text{cos-sim}(\hat{\mathbf{s}}, \mathbf{s}) = \frac{\hat{\mathbf{s}} \cdot \mathbf{s}}{\max(\|\hat{\mathbf{s}}\|_2 \|\mathbf{s}\|_2, \epsilon)}, \quad (4.4)$$

where $\epsilon = 1 \times 10^{-8}$ is used to ensure numerical stability.

Hyperparameters

To enable fair comparison across models, hyperparameters were tuned for SCARF, the FFN binned prediction baseline, and the GNN binned prediction baseline. Parameters were tuned using RayTune [124] with Optuna [123] and an ASHAScheduler. Each model was allotted 50 different hyperoptimization trials for fitting. Models were hyperparameter optimized on a smaller 10,000 spectra subset of NIST20. Parameters are detailed in Table 4.12.

4.7.6 Limitations and future work

We outline several potential directions for future work to address limitations of this work.

1. *Improving the gold standard training annotation pre-processing.* Because SCARF is flexible in that it can match distributions of formula assignments, a key step to improving and building upon this approach is to develop more robust assignments of

formula to training spectra. This includes adding complexity and removing potential assumptions, such as allowing annotations to account for rearrangement, elimination, or charge transfer. A second goal is to identify potentially low quality training spectra, such as ones that emerge from mixtures, and remove these from the inputs. Another potential way to handle such cases would be to model each spectrum peak as an *ensemble* of potential equivalent-mass formulae, which would be particularly helpful in relating SCARF to inverse models such as MIST [64] in which the structure of the molecule and formula identity of each peak cannot be known *a priori*.

2. *Incorporating other model covariates.* Incorporating collision energy features explicitly into the model, as well as negative-ion mode inputs, will increase its usability. This could be enabled by aggregating public data containing these annotations.
3. *Featurizing molecule inputs using different or more powerful molecular encoders.* Recent and simultaneous work to this used a pretrained graph encoder as part of a binned spectrum prediction approach, MassFormer [86]. It is possible to include more powerful molecule or formula encoders into SCARF.
4. *Consideration of interpretability by subgraph attribution and combination with ICEBERG, which appears in Chapter 5.* Following this initial work, we developed a second model, ICEBERG [145], that uses a similar two step modeling approach, but instead encodes fragments, not formula. This increases accuracy and robustness, especially for retrieval, but substantially slows the model. In comparison to ICEBERG, SCARF still has several benefits including speed, the lack of required substructure labeling, and ability to capture potential skeletal rearrangements of molecules (i.e., discontinuities in structure that may not be possible to model by only breaking bonds). An open question and exciting opportunity in the future is to combine these two levels of abstraction and make formula-level predictions with graph-level attribution or featurization.
5. *Retrieval-specific loss functions to enhance retrieval performance.* A significant finding of this work was the noise associated with the retrieval task and lack of correlation with spectrum prediction performance as measured by cosine similarity. Future work may consider how to more directly define loss functions that reflect the task of retrieval.

Table 4.12: **Model and baseline hyperparameters.**

| Model | Parameter | Grid | Value |
|----------------|--------------------------|---------------------------------|-------------------|
| Autoregressive | learning rate | $[1e-4, 1e-3]$ | 0.0009 |
| | learning rate scheduler | - | StepDecay (5,000) |
| | learning rate decay | $[0.7, 1.0]$ | 0.85 |
| | dropout | $\{0.0, 0.1, 0.2, 0.3\}$ | 0.2 |
| | hidden size, d | $\{128, 256, 512\}$ | 512 |
| | gnn layers | $[1, 6]$ | 1 |
| | rnn layers | $[1, 3]$ | 3 |
| | batch size | $\{8, 16, 32, 64\}$ | 64 |
| | weight decay | $\{0, 1e-6, 1e-7\}$ | $1e-6$ |
| | use differences (Eq.4.3) | $\{\text{True}, \text{False}\}$ | True |

Table 4.12 – continued from previous page

| Model | Parameter | Grid | Value |
|-------------|--------------------------------------|-----------------------|-------------------|
| | conv type | - | GatedGraphConv |
| | random walk embed steps (Table 4.11) | [0,20] | 20 |
| | graph pooling | {mean, attention} | mean |
| NEIMS (FFN) | learning rate | $[1e-4, 1e-3]$ | 0.00087 |
| | learning rate scheduler | - | StepDecay (5,000) |
| | learning rate decay | [0.7, 1.0] | 0.722 |
| | dropout | {0.0, 0.1, 0.2, 0.3} | 0.0 |
| | hidden size, d | {64, 128, 256, 512} | 512 |
| | layers, l | {1, 2, 3} | 2 |
| | batch size | {16, 32, 64, 128} | 128 |
| | weight decay | {0, $1e-6$, $1e-7$ } | 0 |
| | use differences (Eq.4.3) | {True, False} | True |
| | num bins (Section 9) | - | 15,000 |
| NEIMS (GNN) | learning rate | $[1e-4, 1e-3]$ | 0.00052 |
| | learning rate scheduler | - | StepDecay (5,000) |
| | learning rate decay | [0.7, 1.0] | 0.767 |
| | dropout | {0.0, 0.1, 0.2, 0.3} | 0.0 |
| | hidden size, d | {64, 128, 256, 512} | 512 |
| | layers, l | [1, 6] | 4 |
| | batch size | {16, 32, 64} | 64 |
| | weight decay | {0, $1e-6$, $1e-7$ } | $1e-7$ |
| | use differences (Eq.4.3) | {True, False} | True |
| | num bins (Section 9) | - | 15,000 |
| | conv type | - | GatedGraphConv |
| | random walk embed steps (Table 4.11) | [0,20] | 19 |
| | graph pooling | {mean, attention} | mean |
| 3DMolMS | learning rate | $[1e-4, 1e-3]$ | 0.00074 |
| | learning rate scheduler | - | StepDecay (5,000) |
| | learning rate decay | [0.7, 1.0] | 0.86 |
| | dropout | {0.0, 0.1, 0.2, 0.3} | 0.3 |
| | hidden size, d | {64, 128, 256, 512} | 256 |
| | layers, l | [1, 6] | 2 |
| | top layers | [1, 3] | 2 |
| | neighbors, k | [3, 6] | 5 |
| | batch size | {16, 32, 64} | 16 |
| | weight decay | {0, $1e-6$, $1e-7$ } | $1e-6$ |
| | num bins (Section 9) | - | 15,000 |
| FixedVocab | learning rate | $[1e-4, 1e-3]$ | 0.00018 |
| | learning rate scheduler | - | StepDecay (5,000) |
| | learning rate decay | [0.7, 1.0] | 0.92 |
| | dropout | {0.0, 0.1, 0.2, 0.3} | 0.3 |
| | hidden size, d | {64, 128, 256, 512} | 512 |
| | layers, l | [1, 6] | 6 |
| | batch size | {16, 32, 64} | 64 |
| | weight decay | {0, $1e-6$, $1e-7$ } | $1e-6$ |
| | num bins (Section 9) | - | 15,000 |
| | conv type | - | GatedGraphConv |
| | random walk embed steps (Table 4.11) | [0,20] | 11 |
| | graph pooling | {mean, attention} | mean |

Table 4.12 – continued from previous page

| Model | Parameter | Grid | Value |
|--------------|---------------------------------------|-----------------------------------|-------------------|
| | formula library size | {1000, 5000, 10000, 25000, 50000} | 5000 |
| SCARF-Thread | learning rate | $[1e-4, 1e-3]$ | 0.000577 |
| | learning rate scheduler | - | StepDecay (5,000) |
| | learning rate decay | [0.7, 1.0] | 0.894 |
| | dropout | {0.0, 0.1, 0.2, 0.3} | 0.3 |
| | hidden size, d | {128, 256, 512} | 512 |
| | mlp layers, l_1 | [1, 3] | 2 |
| | gnn layers, l_2 (Section 4.7.5) | [1, 6] | 4 |
| | batch size | {8, 16, 32, 64} | 16 |
| | weight decay | {0, $1e-6, 1e-7$ } | $1e-6$ |
| | use differences (Eq.4.3) | {True, False} | True |
| | conv type | - | GatedGraphConv |
| | random walk embed steps (Table 4.11) | [0,20] | 20 |
| | graph pooling | {mean, attention} | mean |
| SCARF-Weave | learning rate | $[1e-4, 1e-3]$ | 0.00031 |
| | learning rate scheduler | - | StepDecay (5,000) |
| | learning rate decay | [0.7, 1.0] | 0.962 |
| | dropout | {0.0, 0.1, 0.2, 0.3} | 0.2 |
| | hidden size, d | {128, 256, 512} | 512 |
| | mlp layers, l_1 (Section 9) | [1, 3] | 2 |
| | gnn layers, l_2 (Section 4.7.5) | [1, 6] | 3 |
| | transformer layers, l_3 (Section 9) | [0, 3] | 2 |
| | batch size | {4, 8, 16, 32, 64} | 32 |
| | weight decay | {0, $1e-6, 1e-7$ } | 0 |
| | num bins (Section 9) | - | 15,000 |
| | conv type | - | GatedGraphConv |
| | random walk embed steps (Table 4.11) | [0,20] | 7 |
| | graph pooling | {mean, attention} | attention |

Chapter 5

Generating Mass Spectra as Substructure Sets by Breaking Bonds

This work has previously been made available as a preprint S. Goldman, J. Li, and C. W. Coley, “Generating Molecular Fragmentation Graphs with Autoregressive Neural Networks,” *arXiv preprint arXiv:2304.13136*, 2023. It is undergoing revisions and reproduced here with minor changes to the initially released preprint. I led this work’s conceptualization, implementation, and writing, with input and support from co-authors. Co-author Janet Li was essential in initializing this project, as she worked on a variant of this idea during an undergraduate senior thesis that I mentored and supervised.

5.1 Introduction

Identifying unknown molecules in complex metabolomic or environmental samples is of critical importance to biologists [67], forensic scientists [175], and ecologists alike [18]. Tandem mass spectrometry, MS/MS, is the standard analytical chemistry method for analyzing such samples, favored for its speed and sensitivity [68]. In brief, MS/MS metabolomics experiments isolate, ionize, and fragment small molecules, resulting in a characteristic spectrum for each where peaks correspond to molecular sub-fragments (Figure 5.1A). Importantly, these experiments are high throughput, leading to thousands of detected spectra per single experiment for complex samples such as human serum.

The most straightforward way to identify an unknown molecule from its fragmentation spectrum is to compare the spectrum to a library of known standards [29]. However, spectral libraries only contain on the order of 10^4 compounds—a drop in the bucket compared to the vast size of biologically-relevant chemical space, oft cited as large as 10^{60} [176]. Of the many tandem spectra deposited into a large community library, 87% still cannot be annotated [29]. The accurate prediction of mass spectra from molecular structures would enable these libraries to be augmented with hypothetical compounds and significantly advance the utility of mass spectrometry for structural elucidation. This paradigm of comparing unknown spectra to putative spectra is well established in the adjacent field of proteomics due to the ease of predicting protein fragmentations [177].

Because tandem mass spectrometry experiments physically break covalent bonds in a

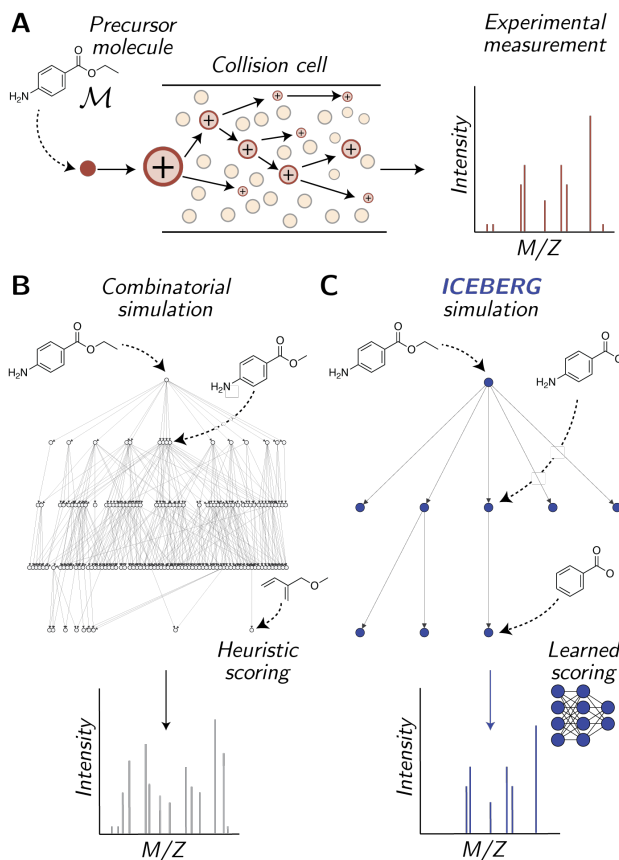


Figure 5.1: **ICEBERG** enables the prediction of tandem mass spectra by efficiently navigating the space of possible fragmentation events. **A**. Example experimental mass spectrum. An input molecule, benzocaine, is depicted entering a mass spectrometer collision cell and fragmenting. The observation of the resulting charged fragments results in a characteristic spectrum. **B**. A combinatorial mass spectrum simulation. The root molecule, benzocaine, is iteratively fragmented by removing atoms or breaking bonds, resulting in a large fragmentation tree. Heuristic rules score nodes in the tree to predict intensities. **C**. ICEBERG spectrum simulation. ICEBERG learns to generate only the most relevant substructures. After generating fragments, a neural network module scores the resulting fragments to predict intensities.

process known as “collision-induced-dissociation” (CID) to create fragments, simulating such fragmentation events computationally is a natural strategy for prediction. Tools from the last decade including MetFrag [71], MAGMa [95], and CFM-ID [51], [73] use fragmentation rules (based on removing atoms or bonds) and local scoring methods to (a) enumerate molecular fragmentation trees and (b) estimate the intensity at each node in the tree with a mix of heuristic rules and statistical learning (Figure 5.1B).

However, these combinatorial methods are computationally demanding and often make inaccurate predictions by *overestimating* the possible fragments (Figure 5.1B, bottom). We recently found CFM-ID to be far less accurate than black-box neural networks [130], an observation separately confirmed by Murphy, Jegelka, Fraenkel, *et al.* [152]. Further, current learned fragmentation models are not easily adapted or scaled to new datasets; Murphy, Jegelka, Fraenkel, *et al.* estimate it would take the leading fragmentation approach, CFM-ID [73], approximately three months on a 64-core machine to train on an approximately 300,000 spectrum dataset.

Alternative strategies that utilize black box neural networks to predict MS/MS spectra have been attempted. They encode an input molecule (i.e., as a fingerprint, graph, or 3D structure) and predict either a 1D binned representation of the spectrum [84], [86], [150], [157], or a set of output formulae corresponding to peaks in the spectrum [130], [152], [153]. While we have demonstrated that predicting molecular formulae provides a fast, accurate, and interpretable alternative to binned representation approaches [130], the improved accuracy surprisingly did not directly translate to better database retrieval for complex natural product molecules contained within the Global Natural Products Social (GNPS) database [28]. We hypothesized that combining the flexibility of neural networks to learn from experimental MS/MS data in reference libraries with the structural-bias of combinatorial fragmentation approaches could lead to increased prediction performance on complex natural product molecules.

Herein, we introduce a hybrid strategy for simulating molecular fragmentation graphs using neural networks, *Inferring Collision-induced-dissociation by Estimating Breakage Events and Reconstructing their Graphs* (ICEBERG). ICEBERG is a two-part model that simulates probable breakage events (Generate) and scores the resulting fragments using a Transformer architecture (Score) (Figure 5.1C; details in Figure 5.2). Our core computational contribution is to leverage previous exhaustive cheminformatics methods for the same task, specifically MAGMa [95], in order to build a training dataset, from which our model learns to make fast estimates prioritizing only likely bond breakages. In doing so, we lift MAGMa and previous bond-breaking approaches into a neural network space with demonstrable benefits in performance.

We evaluate ICEBERG on two datasets: NPLIB1 (GNPS data [28] as used to train the CANOPUS model [57]) and NIST20 [63], which test the model’s ability to predict both complex natural products and small organic standard molecules, respectively. We find that ICEBERG increases cosine similarity of predicted spectra by over 0.05, a 10% improvement over a recent state of the art method on NPLIB1 data. When used to identify molecules in retrospective retrieval studies, ICEBERG leads to a 46% improvement in top 1 retrieval accuracy on a challenging natural product dataset compared to the next best model tested. ICEBERG is fully open-sourced with pretrained weights alongside other existing prediction baseline methods available on GitHub at <https://github.com/samgoldman97/ms-pred>.

5.2 Methods

5.2.1 Datasets

We train our models on the two datasets, NIST20 [63] as generated by the National Institute of Standards and NPLIB1 extracted from the GNPS database [28] and prepared previously by Dührkop, Nothias, Fleischauer, *et al.* [57] and Goldman, Wohlwend, Stražar, *et al.* [64] For each spectrum in the dataset, we first merge all scans at various collision energies, combine peaks that are within 10^{-4} m/z tolerance from each other, renormalize the resulting spectrum by dividing by the maximum observed intensity, and take the square-root of each intensity. We subset the resulting spectrum to keep the top 50 peaks with intensity above 0.003. This normalization process is identical to our previous work [130] and emphasizes (a) removing peaks that are likely noise and (b) combining various collision energies. We refer the reader to [130] for exact details on dataset extraction.

To further normalize the dataset, for each spectrum, we subtract the mass of the adduct ion from each resulting MS2 peak. Concretely, the precursor molecule is ionized with an adduct ion, for instance, H^+ . In this case, the mass of each peak in the spectrum is shifted by the mass of H^+ before proceeding further. In doing so, we normalize against different ionizations. While adduct switching is possible, we note that this is a rarer phenomenon and can be easily interchanged at the data preprocessing step. We make the simplifying assumption that all peaks are singly charged and use mass and m/z interchangeably. Ultimately, each spectrum \mathcal{Y} can be considered a set of mass, intensity tuples, $\mathcal{Y} = \{(m_0, y_0), (m_1, y_1), \dots (m_{|\mathcal{Y}|}, y_{|\mathcal{Y}|})\}$.

5.2.2 Canonical DAG construction

We build a custom re-implementation of the MAGMa algorithm [95] to help create explanatory directed acyclic graphs (DAGs) for each normalized and adduct-shifted spectrum.

Given an input molecule \mathcal{M} , MAGMa iteratively breaks each molecule by removing atoms. Each time an atom is removed, multiple fragments may form, from which we keep all fragments of > 2 heavy (non-hydrogen) atoms. To prevent combinatorial explosion of DAG nodes, we use a Weisfeiler-Lehman isomorphism test [178] to generate a unique hash ID of each generated fragment and reject new fragments with hash IDs already observed. When conducting this test, to remain insensitive to *how* this fragment originated, we hash only the atom identities and bonds in the fragment graph, *not the number of hydrogen atoms*. For instance, consider an ethane fragment in which the terminal carbon was originally double-bonded to a single neighboring atom in the precursor molecule compared to an ethane fragment in which the terminal carbon was single-bonded to two adjacent atoms in the original precursor—our approach applies the same hash ID to both fragments. The molecular formula and hydrogen status for the fragment is randomly selected from the fragments that required the *minimal* number of atom removals. Each fragment corresponds to multiple potential m/z observations due to the allowance for hydrogen shifts equal to the number of broken bonds.

After creating the fragmentation graph for \mathcal{M} , a subset of the fragments are selected to explain each peak in \mathcal{Y} , using the minimum mass differences of under 20 parts-per-million as

the primary filter and the minimal MAGMa heuristic score as a secondary filter. We include nodes along all paths back to the root molecule for each selected fragment. To prune the DAG to select only the most likely paths to each fragment, we design a greedy heuristic. Starting from the lowest level of the DAG, we iteratively select the parent nodes for inclusion into the final DAG that “cover” the highest number of peak-explaining nodes. Finally, the “neutral loss” fragments are added into the DAG, as they provide useful training signals for ICEBERG Generate to learn when to stop fragmenting each molecule.

5.2.3 Model details

DAG generation prediction Using the ground truth DAG as described above, we train a neural network, ICEBERG Generate, to reconstruct the DAG from an input molecule and adduct type. Concretely, our model learns to predict for each fragment, $\mathcal{S}^{(i)}$, the probability that it will fragment at the j^{th} atom:

$$p\left(\mathcal{F}[\mathcal{S}_j^{(i)}]|\mathcal{S}^{(i)}, \mathcal{M}, C\right) = g_{\theta}^{\text{Generate}}(\mathcal{M}, \mathcal{S}^{(i)}, C)_j \quad (5.1)$$

To make this atom-wise prediction, we encode information about the root molecule, fragment molecule, their difference, their respective molecular formulae, the adduct, and the number of bonds that were broken between the root molecule and fragment. To embed the root molecule, we utilize a gated graph neural network [172], $\text{GNN}(\mathcal{M})$, where either average or weighted summations are used to pool embeddings across atoms (specified by a hyperparameter). We utilize the same network to learn representations of the fragment, $\text{GNN}(\mathcal{S}^{(i)})$ and define $\text{GNN}(\mathcal{S}^{(i)})_j$ as the graph neural network-derived embedding of fragment i at the j^{th} atom prior to pooling operation. For all graph neural networks, a one-hot encoding of the adduct type is also added as atom-wise features alongside the bond types and atom types. We define the molecular formula f for each DAG fragment and specify an encoding, Enc , using the Fourier feature scheme defined in [130]. We encode the root and i^{th} node of the fragmentation DAG as $\text{Enc}(f_0)$ and $\text{Enc}(f_i)$, respectively. Lastly, we define a one hot vector for the number of bonds broken, b .

All the encodings described above are concatenated together and a shallow multilayer perceptron (MLP) ending with a sigmoid function is utilized to predict binary probabilities of fragmentation at each atom.

$$p\left(\mathcal{F}[\mathcal{S}_j^{(i)}]|\mathcal{S}^{(i)}, \mathcal{M}, C\right) = \text{MLP}\left([\text{GNN}(\mathcal{M}), \text{GNN}(\mathcal{M}) - \text{GNN}(\mathcal{S}^{(i)}), \text{GNN}(\mathcal{S}^{(i)})_j, \text{Onehot}(b), \text{Enc}(f_i), \text{Enc}(f_0 - f_i)]\right) \quad (5.2)$$

The model is trained to maximize the probability of generating the DAG by minimizing the binary cross entropy loss over each atom for every fragment in an observed spectrum.

DAG intensity prediction The trained Generate module is used to generate DAGs for each input molecule in the training set. In this generation step, molecules are iteratively fragmented beginning with the root \mathcal{M} and the probability of each fragment is computed autoregressively. We define the node indices for an ordering from each fragment $\mathcal{S}^{(i)}$ back to

the root node through its highest likelihood path $\pi[i]$, where $\pi[i, j]$ defines the j^{th} node on this factorization path.

$$p(\mathcal{S}^{(i)}|\mathcal{M}, C) = p(\mathcal{S}^{(i)}|\mathcal{S}^{(\pi(i,1))}, \mathcal{M}, C) \prod_{j=1}^{|\pi[i]|} p(\mathcal{S}^{(\pi[i,j])}|\mathcal{S}^{(\pi[i,j+1])}, \mathcal{M}, C) \quad (5.3)$$

At each step, we maintain only the top 100 most likely fragments in the DAG as a practical consideration until reaching the maximum possible fragmentation depth. To further reduce complexity in the inference step, we maintain the highest scoring isomer from the DAG. This resulting set of fragments is featurized and passed to a Set Transformer module to generate output values at each fragment. Following the notation from the generative model, we featurize each individual fragment with a shallow MLP to generate hidden representations, h_i :

$$h_i = \text{MLP}\left(\left[\text{GNN}(\mathcal{M}), \text{GNN}(\mathcal{M}) - \text{GNN}(\mathcal{S}^{(i)}), \text{GNN}(\mathcal{S}^{(i)}), \text{Onehot}(b), \text{Enc}(f_i), \text{Enc}(f_0 - f_i)\right]\right) \quad (5.4)$$

These are subsequently jointly embedded with a Transformer module and used to predict unnormalized intensity weights at each possible hydrogen shift δ alongside an attention weight α to determine how heavily to weight each prediction for its specified hydrogen shift. To compute the attention weight, we take a softmax over all prediction indices that fall into the same intensity bin (0.1 resolution), $M(i, \delta)$:

$$\hat{y}_\delta^{(i)} = \text{MLP}_{inten}\left(\text{Transformer}(h_0, h_1, h_2, \dots, h_{|\mathcal{T}|})_i\right)_\delta, \quad (5.5)$$

$$\alpha_\delta^{(i)} = \text{Softmax}_{k \in M(i, \delta)}\left(\text{MLP}_{attn}\left(\text{Transformer}(h_0, h_1, h_2, \dots, h_{|\mathcal{T}|})_k\right)\right)_{i, \delta} \quad (5.6)$$

The final intensity prediction for the bin at mass m is then a then a weighted sum over all predictions that fall within this mass bin followed by a sigmoid activation function:

$$\hat{y}_m = \sigma\left(\sum_i \sum_\delta \alpha_\delta^{(i)} \hat{y}_\delta^{(i)} \mathcal{I}[M(i, \delta) = m]\right) \quad (5.7)$$

The model is trained to maximize the cosine similarity between the predicted spectrum and ground truth spectrum.

Model training All models are implemented and trained using Pytorch Lightning [140], the Adam optimizer [141], and the DGL library [173]. Ray [124] is used to complete hyperparameter optimizations over all models and baselines. Models are trained on a single RTX A5000 NVIDIA GPU (CUDA Version 11.6) in under 3 hours for each module. A complete list of hyperparameters and their definition can be found in the Supporting Information.

5.2.4 Baselines

A key component of this work is to extend our robust comparison to previous and contemporary methods [130]. To conduct this benchmarking and specifically emphasize methodological differences rather than data differences, we port and modify code from a number of previous methods including NEIMS [84], NEIMS with a graph neural network [150] 3DMolMS [157], SCARF [130] (Chapter 4), a modified version of GRAFF-MS we refer to as FixedVocab [152], MassFormer [86], and CFM-ID [51] into a single GitHub repository.

Following our previous approach in [130], we emphasize conditioning on the same experimental settings of adduct type *across methods* for fair comparison, excluding collision energies and instrument types as extensions for future work. We rigorously hyperparameter optimize each method for our data regime and train each model on data splits with the exception of CFM-ID for which retraining is not feasible [86], [130], [152].

All model predictions are transformed into binned representations for fair evaluation at a bin resolution of 0.1 from mass 0 to 1,500 Da. Further details are included in Section 5.5.2.

5.3 Results

5.3.1 ICEBERG is trained as a two-stage generative and scoring model

Learning to generate likely substructures. ICEBERG simulates a mass spectrum by generating the substructure fragments from an initial molecule that are most likely to be generated by collision induced dissociation and subsequently measured in the mass spectrometer. We define an input molecule \mathcal{M} (benzocaine example shown in Figure 5.2A) and its observed spectrum \mathcal{Y} , which is a set of intensities at various mass-to-charge values (m/z), termed peaks. Each peak represents one or more observed molecular fragment.

A core question is then *how to generate the set of potential fragments*. These fragments can be sampled from the many possible substructure options, $\mathcal{S}^{(i)} \in (N^{(i)}, E^{(i)}) \subseteq \mathcal{M}$, where the set of nodes and edges in substructures are subsets of the atoms and bonds in the original molecule, $\mathcal{M} \in (N, E)$. Most often, this sampling is accomplished by iteratively and exhaustively removing edges or atoms from the molecular graph, creating a fragmentation graph $\mathcal{T} \in (\mathcal{S}, \mathcal{E})$, where all the nodes in this graph are themselves substructures of the original molecule $\mathcal{S} = \{\mathcal{S}^{(0)}, \mathcal{S}^{(1)}, \dots, \mathcal{S}^{(|\mathcal{T}|)}\}$ ([71], [73], [95]) (Figure 5.1b). However, such a combinatorial approach leads to *thousands* of molecular fragments, making this procedure slow and complicating the second step of estimating intensity values for all enumerated fragments.

We eschew combinatorial generation and instead leverage a graph neural network to parameterize breakage events of the molecule, defining the Generate module of ICEBERG (Figure 5.2A,B). Generate predicts the fragmentation graph iteratively, beginning with just the root of the graph $\mathcal{S}^{(0)} = \mathcal{M}$, borrowing ideas from autoregressive tree generation [179], [180]. At each step in iterative expansion, the model $g_{\theta}^{\text{Generate}}$ assigns a probability of fragmentation to each atom j in the current substructure fragment $\mathcal{S}^{(i)}$, $p(F[\mathcal{S}_j^{(i)}])$. Learned atom embeddings are concatenated alongside embeddings of the root molecule and and a

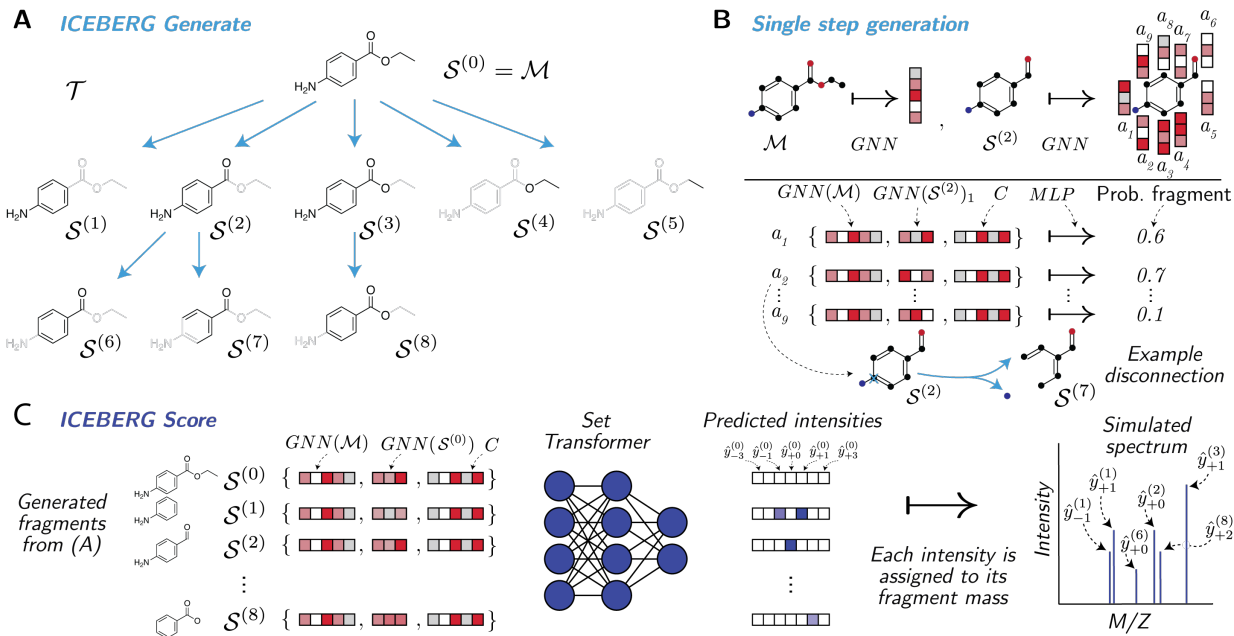


Figure 5.2: **Overview of ICEBERG.** **A.** The target fragmentation directed acyclic graph (DAG) for an example molecule \mathcal{M} , benzocaine. Fragments are colored in black with missing substructures in gray. **B.** Example illustration for the generative process at a single step in the DAG generation predicting subfragments of $\mathcal{S}^{(2)}$. The root molecule \mathcal{M} , fragment of interest $\mathcal{S}^{(2)}$, and context vector C are encoded and used to predict fragment probabilities at each atom of the fragment of interest. A sample disconnection is shown at atom a_2 , resulting in fragment $\mathcal{S}^{(7)}$. **C.** ICEBERG Score module. Fragments generated from **A** are encoded alongside the root molecule. A Set Transformer module predicts intensities for each fragment, allowing mass changes corresponding to the loss or gain of hydrogen atoms, resulting in the final predicted mass spectrum.

context vector C containing metadata such as the ionization adduct type in order to make this prediction. An illustrative example can be seen for fragment $\mathcal{S}^{(2)}$ in Figure 5.2B. Atom a_2 has the highest predicted probability, so this atom is then removed from the graph, leading to the subsequent child node $\mathcal{S}^{(7)}$ (Figure 5.2B). Importantly, the number of child fragments is determined by how many disjoint molecular graphs form upon removal of the j^{th} node from the molecular graph; in this example, fragments $\mathcal{S}^{(1)}$ and $\mathcal{S}^{(4)}$ originate from the same fragmentation event of $\mathcal{S}^{(0)}$ (Figure 5.2A).

In this way, ICEBERG predicts breakages at the level of each *atom*, following the convention of MAGMa [95] rather than each *bond* as is the convention with CFM-ID [73]. We strategically use this abstraction to ensure that all fragmentation events lead to changes in heavy-atom composition. We acknowledge that this formulation does not currently allow for the prediction of skeletal rearrangements and recombinations, which might further improve the model’s ability to explain fragmentation spectra. We refer the reader to Section 5.2.3 for a full description of the model $g_{\theta}^{\text{Generate}}(\mathcal{M}, \mathcal{S}^{(i)}, C)_j$, graph neural network architectures, and context vector inputs.

While this defines a neural network for generation, we must also specify an algorithm for how to *train* this network. Spectral library datasets contain only molecule and spectrum pairs, but not the directed acyclic graph (DAG) \mathcal{T} of the molecule’s substructures that generated the spectrum. We infer an explanatory substructure identity of each peak for model training by leveraging previous combinatorial enumeration methods, specifically MAGMa [95]. For each training molecule and spectrum pair, $(\mathcal{M}, \mathcal{Y})$, we modify MAGMa to enumerate all substructures of \mathcal{M} up to a depth of 3 sequential fragmentation events. We filter enumerated structures to include only those with m/z values appearing in the final spectrum, thereby defining a dataset suitable for training ICEBERG Generate (see Section 5.2.2). As a result, each paired example $(\mathcal{M}, \mathcal{Y})$, in the training dataset is labeled with an estimated fragmentation DAG. Generate learns from these DAGs to generate only the most relevant and probable substructures for a molecule of interest (see Section 5.2.3).

Predicting substructure intensities. After generating a set of potential substructure fragments, we employ a second module, ICEBERG Score, to predict their intensities (Figure 5.2C). Importantly, this design decision enables our models to consider two important physical phenomena: (i) neutral losses and (ii) mass shifts due to hydrogen rearrangements and isotope effects.

Because we elect to fragment molecules at the level of atoms (see Section 5.2.3), multiple substructures can result from a single fragmentation event. In physical experiments, not all of these substructure fragments will be observed; when fragmentation events occur in the collision cell, one fragment often retains the charge of the parent while the other is uncharged and therefore undetected, termed a “neutral loss”. By deferring prediction of intensities to a second module, Generate needs not predict or track whether structures are ionized, greatly reducing the complexity of the fragmentation DAG.

In addition to the occurrence of neutral losses, molecules often undergo complex rearrangements in the collision cell, leading to bond order promotions or reductions (e.g., spurious formation of double bonds when a single bond breaks to maintain valence), the most classic of which is the McLafferty rearrangement [148], [181]. While other approaches attempt to model and estimate where these rearrangements occur using hand-crafted rules [73], we instead adopt the framework of Ridder, Hooft, and Verhoeven [95] to consider hydrogen tolerances. That is, for each generated molecular substructure $\mathcal{S}^{(i)}$ we consider the possibility that this fragment is observed not only at its mass, but also at masses shifted by discrete hydrogen masses, $\pm\delta\text{H}$. This design choice also simplifies Generate by deferring specification of hydrogen counts to the second model. In addition to accounting for a mass shift of 1 hydrogen, such flexibility also allows the model to predict the common $M+1$ isotopes for carbon- and nitrogen- containing compounds.

Mathematically, we define a neural network, $g_{\theta}^{\text{Score}}$ that predicts multiple intensities for each fragment $\hat{y}_{\delta}^{(i)}$ corresponding to different hydrogen shifts, δ :

$$\hat{y}_{\delta}^{(i)} = g_{\theta}^{\text{Score}}(\mathcal{M}, \mathcal{S}^{(i)}, \mathcal{T}, C)_{\delta} \quad (5.8)$$

In practice, we predict up to 13 intensities at each fragment (i.e., $\{+0\text{H}, \pm 1\text{H}, \dots, \pm 6\text{H}\}$). For each individual subfragment, the tolerance is further restricted to the number of bonds broken, most often less than 6. We then take the masses of all fragments, perturb them by

the corresponding hydrogen or isotope shifts, and aggregate them into a set of unique m/z peaks by summing the intensities of perturbed fragments with the same m/z value.

To consider all fragments simultaneously in a permutation-invariant manner, g_{θ}^{Score} is parameterized as a Set Transformer network [91], [92]. We train this second module to maximize the cosine similarity between the ground truth spectrum and the predicted spectrum after converting the set of substructures and intensities to m/z peaks.

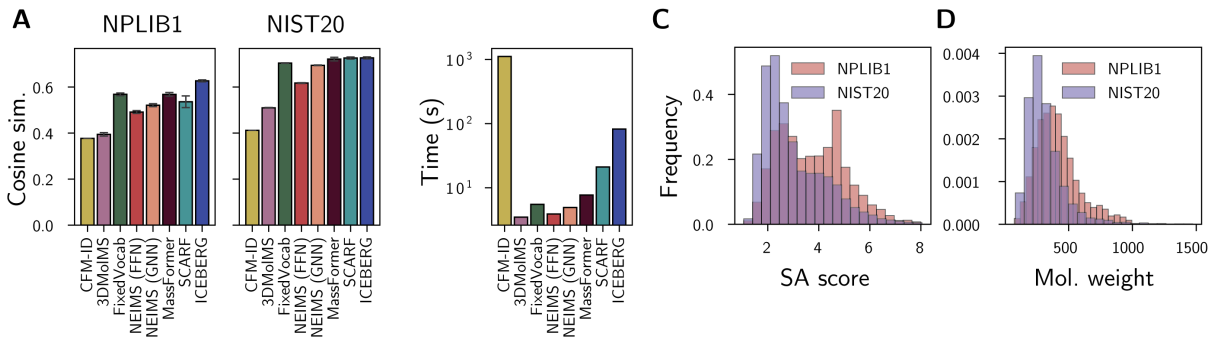


Figure 5.3: **ICEBERG predictions are highly accurate.** **A.** Cosine similarities to true spectra on NPLIB1 (left) and NIST20 respectively (right) for CFM-ID [73], 3DMolIMS [157], FixedVocab [152], NEIMS (FFN) [84], NEIMS (GNN) [150], MassFormer [86], SCARF [130], and ICEBERG. Error bars represent 95% confidence intervals using the standard error of the mean across three random seeds on a single test set split. **B.** Time required to predict spectra for 100 molecules randomly sampled from NIST20 on a single CPU, including the time to load models into memory. **C,D.** Comparison of NPLIB1 and NIST20 molecules in terms of synthetic accessibility (SA) score [182] and molecular weight (Mol. weight).

At test time, we generate the top 100 most likely fragments from ICEBERG Generate and predict intensities for these fragments and their possible hydrogen shifts using ICEBERG Score. We find this tree size allows our model to consider sufficiently many potential fragments while maintaining a speed advantage over previous fragmentation approaches.

5.3.2 ICEBERG enables highly accurate spectrum prediction

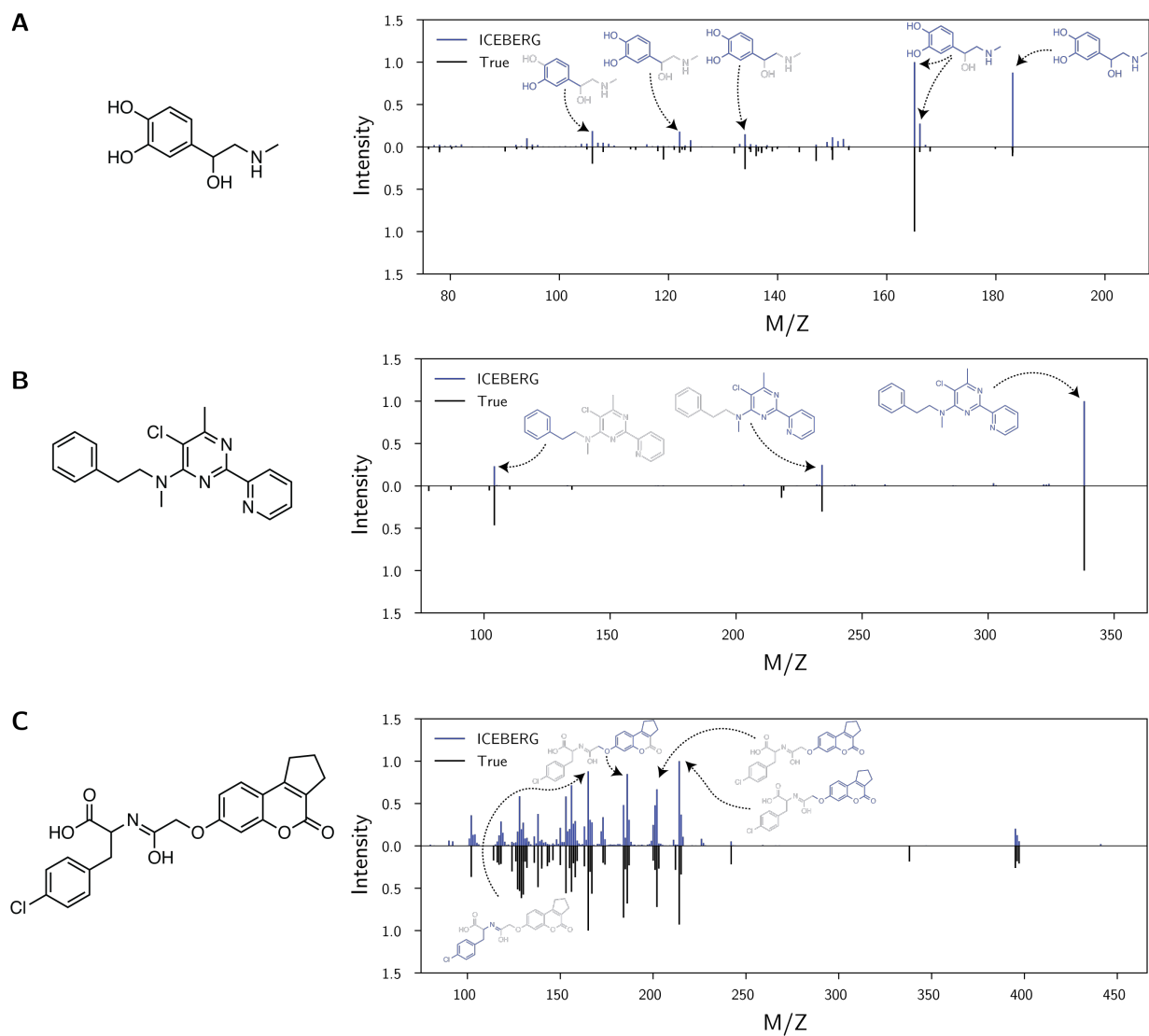


Figure 5.4: **Examples of predicted spectra from ICEBERG.** Predictions are shown as generated by ICEBERG trained on NPLIB1 for select test set examples GNPS: CCMSLIB00003137969 (A), MoNA: 001659 (B), and GNPS: CCMSLIB00000080524 (C). The input molecular structures are shown (left); fragmentation spectra are plotted (right) with predictions (top, blue) and ground truth spectra (bottom, black). Molecular fragments are shown inset. Spectra are plotted with m/z shifted by the mass of the precursor adduct. All examples shown were not included in the model training set.

We evaluate ICEBERG on its ability to accurately simulate positive ion mode mass spectra for both natural product like molecules and smaller organic molecules under 1,500 Da. Using the data cleaning pipeline from [130], we compile a public natural products dataset NPLIB1 with 10,709 spectra (8,533 unique structures) [28], [57], [64] as well as a gold standard chemical library NIST20 with 35,129 spectra (24,403 unique structures) [63]. We note that NPLIB1 was previously named ‘CANOPUS’, renamed here to disambiguate the data from the tool CANOPUS [57]. Both datasets are split into structurally disjoint 90%/10% train-test splits, with 10% of the training data reserved for model validation (see Section 5.2.1). We compare ICEBERG against an expansive suite of contemporary and competitive methods [73], [84], [86], [130], [150], [152], [157]; all methods excluding CFM-ID are reimplemented, hyperparameter optimized, and trained on equivalent data splits, further described in the Supporting Information.

To measure performance, we calculate the average cosine similarity between each predicted spectrum and the true spectrum, as cosine similarity is widely used to cluster mass spectra in molecular networking [69]. We find that ICEBERG outperforms the next method MassFormer on the natural product focused dataset (Figure 5.3A; Table 5.3). ICEBERG achieves an average cosine similarity of 0.627, compared to MassFormer with cosine similarity of 0.568—a 10% improvement.

Surprisingly, however, this boost in performance extends only marginally to the gold standard dataset, NIST20. ICEBERG, while still outperforming binned spectrum prediction approaches (i.e., NEIMS [84]) on this dataset, is nearly equivalent to SCARF (0.727 v. 0.726) [130]. Still, our model performs substantially better than CFM-ID and uses only a fraction of the computational resources (Figure 5.3B). Unlike previous physically inspired models, because ICEBERG only samples the most relevant fragments from chemical space, it requires just over 1 CPU second per spectrum.

We hypothesize that the discrepancy in performance improvement between NPLIB1 and NIST20 may be partially explained by differences in the chemical spaces they cover. Many molecules within NPLIB1 are natural products with more complicated chemical scaffolds. To characterize this, we analyzed the distributions for both the synthetic accessibility (SA) score [182], [183] (Figure 5.3C) and molecular weight (Figure 5.3D), both proxies for molecular complexity. In concordance with our hypothesis, we find that SA scores and molecular weight are substantially higher on NPLIB1 than NIST20: NPLIB1 has an average SA score of 3.75, compared to 3.01 for NIST20; the datasets have average molecular weights of 413 Da and 317 Da respectively.

5.3.3 Model predictions are interpretable

In addition to accurate predictions, a key benefit of simulating fragmentation events is that predictions are interpretable, even for highly complex molecules. Each predicted peak from ICEBERG is directly attributed to a fragment of the predicted molecule.

By inspecting certain patterns and examples, we find expected broken bonds. Weaker bonds such as carbon-oxygen and carbon-nitrogen bonds tend to more reliably break, compared to carbon-carbon bonds and more complex ring breakages (Figure 5.4A). A second strength of using fragmentation-based models can be seen in Figure 5.4B, where despite the heteroaromatic ring structures, our model is still able to correctly predict peak intensities

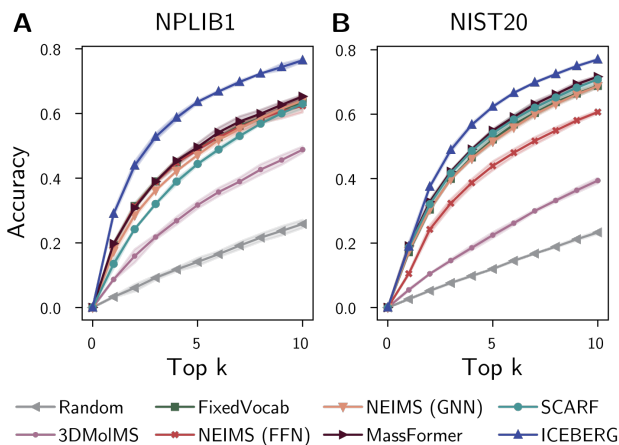


Figure 5.5: **ICEBERG enables improved spectrum retrieval over other methods on both NPLIB1 (A) and NIST20 (B) compared to other spectrum prediction models.** Top k retrieval accuracy is computed by ranking a list of 49 additional candidates by putative cosine similarity as predicted by the model and determining the fraction of times that the true molecule is within the first k entries. A 95% confidence interval across three random model training seeds is shaded.

by predicting a small number of carbon-nitrogen breakages.

Further alignment can be seen within the intensity prediction module. Because ICEBERG predicts multiple intensities for each substructure corresponding to hydrogen shifts, 2 peaks can be present when a single bond breaks. In fragmentation example of Figure 5.4A, the most intense peak is estimated at the mass shift of -1H from the original fragment, indicating that ICEBERG correctly recognizes the hydroxyl group will likely leave as neutral H_2O and result in a hydrogen rearrangement.

5.3.4 Fragmentation simulations lead to improved structural elucidation

We next demonstrate that ICEBERG improves the structural elucidation of unknown molecules using reference libraries of model-predicted spectra. We design a retrospective evaluation using our labeled data to resemble the prospective task of spectrum lookup within libraries. For each test spectrum, we extract up to 49 “decoy” isomers from PubChem [128] with the highest Tanimoto similarity to the true molecular structure. The consideration of up to 50 isomers mimics the realistic elucidation setting, as an unknown spectrum can yield clues regarding certain properties of its source molecule (e.g., computed using MIST [64], CSI:FingerID [53], or molecular networking [69]), which narrows the chemical space of possible molecules to a smaller, more relevant set. We predict the fragmentation spectrum for each isomer and, for each model, we rank these possible matches by their spectral similarity to the spectrum of interest and compute how often the true molecule is found within the top k ranked isomers for different values of k .

We find that ICEBERG improves upon the next best model by a margin of 9% accuracy (a nearly 46% relative improvement) in *top 1* retrieval accuracy for the NPLIB1 dataset (Figure 5.5A; Table 5.5). Previous models with high spectrum prediction accuracies have

Table 5.1: **Comparing the accuracy of spectrum prediction on NIST20 using random (easier) or scaffold (harder) split.**

| NIST20 | Cosine sim. | |
|-------------|--------------|----------------|
| | Random split | Scaffold split |
| CFM-ID | 0.412 | 0.411 |
| 3DMolMS | 0.510 | 0.466 |
| FixedVocab | 0.704 | 0.658 |
| NEIMS (FFN) | 0.617 | 0.546 |
| NEIMS (GNN) | 0.694 | 0.643 |
| MassFormer | 0.721 | 0.682 |
| SCARF | 0.726 | 0.669 |
| ICEBERG | 0.727 | 0.699 |

struggled on this task due to their poor ability to differentiate structurally similar isomers [130]. Our structure-based model appears to excel in retrieval and may have out-of-domain robustness beneficial to this task.

While the effect for top 1 retrieval accuracy on the NIST20 dataset is not pronounced (i.e., tied with MassFormer), ICEBERG outperforms the next best model by an absolute margin of over 5%, a 7.5% relative improvement at top 10 accuracy (Figure 5.5B, Table 5.4). These results underscore the real world utility of ICEBERG to identify unknown molecules of interest.

5.3.5 Challenging test splits better explain retrieval performance

The strong performance on the retrieval task, particularly for increasing values of k on NIST20, suggests that ICEBERG is able to generalize well to decoys not appearing in the training set and to account for how structural changes should affect fragmentation patterns. While encouraging, we observed only minor increases in cosine similarity accuracy when predicting spectra using NIST20 (Figure 5.3, Table 5.2).

To try to explain this apparent discrepancy, we reevaluate prediction accuracy on a more challenging dataset split. We retrain all models on the NIST20 utilizing a Murcko scaffold split of the data [42] with smaller scaffold clusters (i.e., more unique compounds) placed in the test set. This split enforces that molecules in the test set will be more distant and less similar to the training set, probing the ability of each model to generalize in a more stringent setting than our previous random split.

In the more strict scaffold split evaluation, the improved accuracy of ICEBERG over existing models is striking (Table 5.1). We find that ICEBERG outperforms MassFormer and SCARF by 0.017 and 0.03– 2% and 4% improvements respectively. These results suggest that, particularly for standard libraries with more homogeneous molecules, more challenging scaffold split evaluations may yield performance metrics that better correlate with performance on the structural elucidation problem (retrieval).

5.4 Discussion

We have proposed a physically-grounded mass spectrum prediction strategy we term ICEBERG. From a computational perspective, this integration of neural networks into fragmentation prediction is enabled by (a) bootstrapping MAGMa to construct fragmentation trees on which our model is trained, (b) posing the tree generation step as a sequential prediction over atoms, and (c) predicting multiple intensities at each generated fragment with a second module in order to account for hydrogen rearrangements and isotopic peaks. By learning to generate fragmentation events, ICEBERG is able to accurately predict mass spectra, yielding especially strong improvements for natural product molecules under evaluation settings of both spectrum prediction and retrieval.

ICEBERG establishes new state of the art performance for these tasks, yet there are some caveats we wish to highlight. First, while we learn to generate molecular substructures to explain each peak, there are no guarantees that they are the correct physical explanations given the number of potential equivalent-mass atom and bond rearrangements that could occur. Second, while we achieve increased accuracy, this comes at a higher computational cost of roughly 1 CPU second per molecule, nearly an order of magnitude more than other neural approaches like SCARF [130]. Future work will consider more explicitly how to synergize fragment- and formula- prediction approaches to achieve higher accuracy and speed. In addition to model architecture modifications, we anticipate model accuracy improvements from modeling other covariates such as collision energy, instrument type, and even jointly modeling MS/MS with other analytical chemistry measurements such as FTIR [184].

The discovery of unknown metabolites and molecules is rapidly expanding our knowledge of potential medical targets [13], the effects of environmental toxins [19], and the diversity of biosynthetically accessible chemical space [185]. We envision exciting possibilities to apply our new model to expand the discovery of novel chemical matter from complex mixtures.

5.5 Additional Results

We include a complete set of all model results and metrics (i.e., Cosine similarity, coverage, validity, and time for 100 spectrum predictions) in Tables 5.2 and 5.3. Interestingly, while cosine similarity is higher for ICEBERG on nearly all evaluations, the fraction of peaks in the ground truth spectrum explained by the predicted spectrum, *coverage*, does not increase. We posit that the more strict and rigid fragment-grounding causes our model to miss certain lower intensity peaks. In doing so, however, the model is able to maintain higher accuracy.

We note that not all predicted fragments are valid, as a very small fraction of predicted fragments fail the RDBE molecular formula test [131] when hydrogen shifts are considered.

Full results for retrieval are shown for NIST20 and NPLIB1 and in Tables 5.4 and 5.5 respectively.

5.5.1 Graph neural network details

ICEBERG relies upon graph neural network embeddings of the molecule $\text{GNN}(\mathcal{M})$. Given the widespread use and descriptions of such models, we refer the reader to Li, Tarlow,

Table 5.2: **Spectra prediction accuracy on NIST20, and NIST20 (scaffold split)**. Results are shown for CFM-ID [73], 3DMolMS [157], FixedVocab [152], NEIMS (FFN) [84], NEIMS (GNN) [150], MassFormer [86], SCARF [130], and ICEBERG. Cosine similarity is calculated at a bin resolution of 0.1 m/z; Coverage indicates the fraction of peaks in the ground truth spectrum explained by the predicted spectrum on average; Valid indicates the fraction of predicted peaks that can be explained as a subset of the precursor molecule’s molecular formula (filtered to formulae with ring-double bond equivalents of greater or equal to zero); Time (s) indicates the number of seconds required to predict 100 random spectra from NIST20 on a single CPU including the time to load the model. Values are shown \pm 95% confidence intervals of the mean across three random seeds for the random split. Scaffold split experiments are only repeated once and shown without confidence intervals. The best value in each column is typeset in bold.

| Dataset | NIST20 (Random split) | | | NIST20 (Scaffold split) | | | Time (s) |
|-------------|-------------------------------------|-------------------------------------|------------------------------------|-------------------------|--------------|-------------|------------|
| | Cosine sim. | Coverage | Valid | Cosine sim. | Coverage | Valid | |
| CFM-ID | 0.412 \pm 0.000 | 0.278 \pm 0.000 | 1.00 \pm 0.000 | 0.411 | 0.272 | 1.00 | 1114.7 |
| 3DMolMS | 0.510 \pm 0.001 | 0.734 \pm 0.002 | 0.94 \pm 0.002 | 0.466 | 0.707 | 0.97 | 3.5 |
| FixedVocab | 0.704 \pm 0.001 | 0.788 \pm 0.001 | 1.00 \pm 0.000 | 0.658 | 0.785 | 1.00 | 5.5 |
| NEIMS (FFN) | 0.617 \pm 0.001 | 0.746 \pm 0.002 | 0.95 \pm 0.001 | 0.546 | 0.717 | 0.96 | 3.9 |
| NEIMS (GNN) | 0.694 \pm 0.001 | 0.780 \pm 0.001 | 0.95 \pm 0.001 | 0.643 | 0.766 | 0.97 | 4.9 |
| MassFormer | 0.721 \pm 0.004 | 0.790 \pm 0.006 | 0.93 \pm 0.006 | 0.682 | 0.789 | 0.96 | 7.7 |
| SCARF | 0.726 \pm 0.002 | 0.807 \pm 0.001 | 1.00 \pm 0.000 | 0.669 | 0.807 | 1.00 | 21.1 |
| ICEBERG | 0.727 \pm 0.002 | 0.754 \pm 0.002 | 1.00 \pm 0.000 | 0.699 | 0.771 | 1.00 | 82.2 |

Brockschmidt, *et al.* [172] for a description of the gated graph neural networks we employ. We utilize the DGL library [173] to implement and featurize molecular graphs. Because our fragmentation method relies upon iteratively removing atoms, we often have molecular fragments with incomplete valence shells that would not be parsed by RDKit [158]. As such, we opt for a more minimal set of atom and bond features described in Table 5.7 .

5.5.2 Additional baseline details

We compare ICEBERG to three classes of spectrum prediction models: *binned prediction* models that encode a molecule and output a single fixed length discretized spectrum vector, *fragmentation prediction* models such as ICEBERG that directly learn to fragment and break the input molecule, and *formula prediction* models that predict molecular formula peaks and corresponding intensities. We briefly describe each model used below, emphasizing the novelty of their contribution.

For all baselines that utilize 2D graph neural networks that are not pretrained (i.e., NEIMS (GNN)), SCARF, and FixedVocab), we utilize a common graph neural network base architecture consisting of gated graph neural network message passing steps as in ICEBERG (Section 5.5.1) [172]. Baseline GNNs are provided with a wide range of atom features including element type, the degree of each atom, hybridization type, formal charge, a binary flag denoting whether the atom is in a ring-system, atomic mass as a floating point value, a one hot encoding of the chiral tag, the type of the atom, and a random walk embedding step.

Table 5.3: **Spectra prediction accuracy on NPLIB1**. Results are shown for CFM-ID [73], 3DMolMS [157], FixedVocab [152], NEIMS (FFN) [84], NEIMS (GNN) [150], MassFormer [86], SCARF [130], and ICEBERG. Cosine similarity is calculated at a bin resolution of 0.1 m/z; Coverage indicates the fraction of peaks in the ground truth spectrum explained by the predicted spectrum on average; Valid indicates the fraction of predicted peaks that can be explained as a subset of the precursor molecule’s molecular formula (filtered to formulae with ring-double bond equivalents of greater or equal to zero). Values are shown \pm 95% confidence intervals as measured by three random seeds on a single test split. The best value in each column is typeset in bold.

| Dataset | NPLIB1 | | |
|-------------|-------------------------------------|-------------------------------------|------------------------------------|
| | Cosine sim. | Coverage | Valid |
| CFM-ID | 0.377 \pm 0.000 | 0.235 \pm 0.000 | 1.00 \pm 0.000 |
| 3DMolMS | 0.394 \pm 0.004 | 0.507 \pm 0.001 | 0.92 \pm 0.001 |
| FixedVocab | 0.568 \pm 0.003 | 0.563 \pm 0.002 | 1.00 \pm 0.000 |
| NEIMS (FFN) | 0.491 \pm 0.003 | 0.524 \pm 0.002 | 0.95 \pm 0.001 |
| NEIMS (GNN) | 0.521 \pm 0.003 | 0.547 \pm 0.005 | 0.94 \pm 0.001 |
| MassFormer | 0.568 \pm 0.004 | 0.573 \pm 0.009 | 0.95 \pm 0.003 |
| SCARF | 0.536 \pm 0.013 | 0.552 \pm 0.016 | 1.00 \pm 0.000 |
| ICEBERG | 0.627 \pm 0.002 | 0.549 \pm 0.003 | 1.00 \pm 0.000 |

All models tested are provided the same set of experimental covariates: a onehot vector corresponding to the adduct type. For graph models, this is concatenated as an atom level feature, whereas in other models, it is concatenated to the intermediate hidden representation.

Binned prediction

1. NEIMS (FFN): Wei, Belanger, Adams, *et al.* [84] encode a fingerprint representation of a molecule with a feed forward neural network and predict a fixed dimensional discrete binned output vector. In addition to predicting a fragment vector, Wei, Belanger, Adams, *et al.* also predict the mass differences from the precursor ion as an independent vector and combine the two outputs with a learned, weighted sum. Originally trained on EI mass spectrometry, we replicate this approach on our datasets of ESI mass spectra.

Key hyperparameters are whether or not to use mass differences in the output prediction, the number of feed forward layers, hidden size, and dropout between hidden layers.

2. NEIMS (GNN): Zhu, Liu, and Hassoun [150] extend the NEIMS model to utilize graph neural networks instead of feed forward networks applied to fingerprints. We utilize our default graph neural network layer scheme described above, emphasizing that the novelty of this contribution is to switch to a 2D graph neural network encoder from

Table 5.4: NIST20 spectra retrieval top k accuracy for different values of k , \pm a 95% confidence interval across three random seeds on the same test set.

| Top k | 1 | 2 | 3 | 4 | 5 | 8 | 10 |
|-------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|
| Random | 0.026 \pm 0.001 | 0.052 \pm 0.001 | 0.076 \pm 0.002 | 0.098 \pm 0.001 | 0.120 \pm 0.001 | 0.189 \pm 0.003 | 0.233 \pm 0.004 |
| 3DMolMS | 0.055 \pm 0.003 | 0.105 \pm 0.000 | 0.146 \pm 0.005 | 0.185 \pm 0.007 | 0.225 \pm 0.009 | 0.332 \pm 0.005 | 0.394 \pm 0.008 |
| FixedVocab | 0.172 \pm 0.004 | 0.304 \pm 0.004 | 0.399 \pm 0.002 | 0.466 \pm 0.007 | 0.522 \pm 0.012 | 0.638 \pm 0.009 | 0.688 \pm 0.006 |
| NEIMS (FFN) | 0.105 \pm 0.003 | 0.243 \pm 0.012 | 0.324 \pm 0.013 | 0.387 \pm 0.011 | 0.440 \pm 0.014 | 0.549 \pm 0.010 | 0.607 \pm 0.005 |
| NEIMS (GNN) | 0.175 \pm 0.005 | 0.305 \pm 0.003 | 0.398 \pm 0.002 | 0.462 \pm 0.004 | 0.515 \pm 0.005 | 0.632 \pm 0.007 | 0.687 \pm 0.005 |
| MassFormer | 0.191 \pm 0.008 | 0.328 \pm 0.006 | 0.422 \pm 0.003 | 0.491 \pm 0.002 | 0.550 \pm 0.005 | 0.662 \pm 0.005 | 0.716 \pm 0.003 |
| SCARF | 0.187 \pm 0.008 | 0.321 \pm 0.011 | 0.417 \pm 0.007 | 0.486 \pm 0.008 | 0.541 \pm 0.009 | 0.652 \pm 0.008 | 0.708 \pm 0.009 |
| ICEBERG | 0.189 \pm 0.012 | 0.375 \pm 0.005 | 0.489 \pm 0.007 | 0.567 \pm 0.005 | 0.623 \pm 0.004 | 0.725 \pm 0.003 | 0.770 \pm 0.002 |

Table 5.5: NPLIB1 spectra retrieval top k accuracy for different values of k , \pm a 95% confidence interval across three random seeds on the same test set.

| Top k | 1 | 2 | 3 | 4 | 5 | 8 | 10 |
|-------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|
| Random | 0.033 \pm 0.002 | 0.061 \pm 0.010 | 0.092 \pm 0.007 | 0.118 \pm 0.005 | 0.141 \pm 0.012 | 0.216 \pm 0.012 | 0.258 \pm 0.012 |
| 3DMolMS | 0.087 \pm 0.003 | 0.159 \pm 0.020 | 0.218 \pm 0.008 | 0.268 \pm 0.012 | 0.317 \pm 0.011 | 0.427 \pm 0.016 | 0.488 \pm 0.010 |
| FixedVocab | 0.193 \pm 0.007 | 0.314 \pm 0.008 | 0.390 \pm 0.005 | 0.448 \pm 0.010 | 0.492 \pm 0.003 | 0.587 \pm 0.010 | 0.635 \pm 0.011 |
| NEIMS (FFN) | 0.195 \pm 0.005 | 0.313 \pm 0.005 | 0.388 \pm 0.006 | 0.447 \pm 0.012 | 0.488 \pm 0.003 | 0.585 \pm 0.014 | 0.624 \pm 0.020 |
| NEIMS (GNN) | 0.174 \pm 0.014 | 0.285 \pm 0.008 | 0.362 \pm 0.004 | 0.422 \pm 0.001 | 0.471 \pm 0.003 | 0.586 \pm 0.013 | 0.640 \pm 0.010 |
| MassFormer | 0.198 \pm 0.003 | 0.308 \pm 0.004 | 0.389 \pm 0.001 | 0.454 \pm 0.004 | 0.496 \pm 0.012 | 0.599 \pm 0.012 | 0.653 \pm 0.005 |
| SCARF | 0.135 \pm 0.014 | 0.242 \pm 0.001 | 0.320 \pm 0.003 | 0.389 \pm 0.008 | 0.444 \pm 0.004 | 0.569 \pm 0.003 | 0.630 \pm 0.015 |
| ICEBERG | 0.290 \pm 0.008 | 0.439 \pm 0.013 | 0.528 \pm 0.010 | 0.587 \pm 0.009 | 0.636 \pm 0.004 | 0.723 \pm 0.001 | 0.764 \pm 0.005 |

fingerprints. In such a case, we also concatenate model covariates (i.e., adduct one hot) as additional atom level features to the GNN).

Key hyperparameters are whether or not to use mass differences in the output prediction, the number of GNN layers, the number of random walk embedding steps (i.e., positional encodings of the nodes computed using DGL), and dropout between hidden layers.

3. MassFormer: Young, Wang, and Röst [86] further extend the NEIMS (GNN) approach using a *pretrained* graph neural network model, the Graphormer [137]. Rather than train the model entirely from scratch to make a mass and mass difference prediction vector, MassFormer encodes molecules using a large, 48M parameter model previously pretrained. Our MassFormer implementation replaces the GNN with the Graphormer model. Covariates (i.e., adduct type one hot) are concatenated after the Graphormer encoding. We utilize the same scheme of Young, Wang, and Röst to use the pretrained Graphormer Version 2, reinitialize all 12 intermediate Graphormer layers, reinitialize the layer norm, and keep only the learned token embeddings. We allow this model to predict mass differences as in NEIMS by default.

Key hyperparameters are the number of feed forward layers after the Graphormer encoding, the dropout between these layers, and the hidden dimension size.

4. 3DMolMS: Hong, Li, Welch, *et al.* [157] replace the 2D graph neural network architecture with a 3D point cloud graph neural network to predict binned spectra. However, they do not predict mass differences as in MassFormer or NEIMS, substantially reducing model performance in our reimplementing setting.

Table 5.6: **Dataset details.**

| Name | # Spectra | # Molecules | Description |
|--------|-----------|-------------|---------------------------------|
| NIST20 | 35,129 | 24,403 | Standards library |
| NPLIB1 | 10,709 | 8,533 | Natural products from GNPS [28] |

Table 5.7: **Graph neural network (GNN) atom features.**

| Name | Description |
|-------------------------|---|
| Element type | one-hot encoding of the element type |
| Hydrogen number | one-hot encoding the number of hydrogens on each atom |
| Adduct type | one-hot encoding of the ionization adduct |
| Random walk embed steps | positional encodings of the nodes computed using DGL |

Key hyperparameters are the hidden size of their encoder, number of layers in the 3D neural network, number of top layers after the 3D neural network, and the number of neighbors to utilize for message passing in 3D space.

Fragmentation prediction

1. CFM-ID: Wang, Liigand, Tian, *et al.* [51] parametrize a Markov model of fragmentation by first combinatorially fragmenting a molecule of interest then subsequently estimating transition probabilities across the fragmentation graph to predict peak intensities. We utilize the pretrained version 4 available for public release, rather than retraining the model.

Formula prediction

1. FixedVocab: Murphy, Jegelka, Fraenkel, *et al.* [152] recently introduced GRAFF-MS, a model that predicts spectra at the level of molecular formula by pre-defining a list of 10,000 common molecular formula and molecular formula differences from the precursor. They encode the input molecule with a graph neural network and predict intensities at each predefined molecular formula. Due to unavailability of source code, we define a variation of their model we term FixedVocab that predicts intensities at a library of the most popular formula fragments and formula losses in the training dataset. We utilize a cosine similarity loss function to train the model. We revise the name of this model to delineate the differences in loss function and consideration of isotopes and adducts between our implementations.

Key hyperparameters are the hidden size of the model, number of graph neural network layers, random walk embedding steps, graph pooling, and the size of the fixed-length formula vocabulary.

An earlier version of this work under-reported the accuracy of this model, as the implementation did not allow for the prediction of precursor compound masses.

2. SCARF: In Chapter 4 [130], we introduced a similar strategy to Murphy, Jegelka, Fraenkel, *et al.* [152], but rather than use a fixed vocabulary, they autoregressively generate formula candidates as a prefix tree first with the SCARF-Thread model. This smaller vocabulary size conditioned on the input molecule enables them to deploy a second neural network model SCARF-Weave to explicitly encode the candidate formulae with a neural network alongside the molecule in a second modeling step to predict intensities.

We refer the reader to [130] for a complete description of hyperparameters, as this model is used without changes from the original work.

5.5.3 Hyperparameters

To fairly compare the various methods, we apply a rigorous hyperparameter optimization scheme. We use RayTune [124] with Optuna [123] and an ASHAScheduler to identify hyperparameters from a grid set of options. All models were allotted 50 different hyperoptimization trials on a 10,000 spectra subset of NIST20. Hyperparameter descriptions for ICEBERG are provided in Table 5.8. Selected baseline and ICEBERG hyperparameters are defined in Tables 5.10 and 5.9 respectively. When possible, we recycle parameters as selected in Goldman, Bradshaw, Xin, *et al.* [130], as the hyperparameter optimization scheme is equivalent.

Table 5.8: **Hyperparameter descriptions for ICEBERG.**

| Name | Model (Generate or Score) | Description |
|-----------------------------------|---------------------------|---|
| learning rate | both | optimizer learning rate |
| learning rate decay (5,000 steps) | both | step-wise learning rate decay every 5,000 model weight update steps |
| dropout | both | model dropout applied in-between linear hidden layers |
| hidden size | both | number of hidden layers |
| gnn layers | both | number of graph neural network layers to encode molecules and fragments |
| mlp layers | both | number of feed forward layers to encode concatenated representations |
| transformer layers | Score | number of set transformer attention layers after mlp encoding |
| batch size | both | number of spectra to include in each training batch |
| weight decay | both | optimizer weight decay |
| random walk embed steps | Generate | number of random walk embedding steps to for graph neural network atom features |
| graph pooling | both | how to combine atom features into a single representation |
| bin size | Score | binned spectrum resolution spacing from 0 Da to 1,500 Da |

Table 5.9: **ICEBERG Generate and Score hyperparameter grid and selected values.**

| Model | Parameter | Grid | Value |
|------------------|-----------------------------------|------------------------------|---------|
| ICEBERG Generate | learning rate | $[1e-4, 1e-3]$ | 0.00099 |
| | learning rate decay (5,000 steps) | $[0.7, 1.0]$ | 0.7214 |
| | dropout | $\{0.1, 0.2, 0.3\}$ | 0.2 |
| | hidden size | $\{128, 256, 512\}$ | 512 |
| | mlp layers | - | 1 |
| | gnn layers | $[1, 6]$ | 6 |
| | batch size | $\{8, 16, 32, 64\}$ | 32 |
| | weight decay | $\{0, 1e-6, 1e-7\}$ | 0 |
| | random walk embed steps | $[0, 20]$ | 14 |
| | graph pooling | $\{\text{mean, attention}\}$ | mean |
| ICEBERG Score | learning rate | $[1e-4, 1e-3]$ | 0.00074 |
| | learning rate decay (5,000 steps) | $[0.7, 1.0]$ | 0.825 |
| | dropout | $\{0.1, 0.2, 0.3\}$ | 0.1 |
| | hidden size | $\{128, 256, 512\}$ | 256 |
| | mlp layers | $[0, 3]$ | 1 |
| | gnn layers | $[1, 6]$ | 4 |
| | transformer layers | $[0, 3]$ | 3 |
| | batch size | $\{8, 16, 32, \}$ | 32 |
| | weight decay | $\{0, 1e-6, 1e-7\}$ | $1e-7$ |
| | bin size | - | 0.1 |
| | graph pooling | $\{\text{mean, attention}\}$ | mean |

Table 5.10: **Baseline hyperparameters.**

| Model | Parameter | Grid | Value |
|-----------------------------------|-----------------------------------|--------------------------|----------------|
| NEIMS (FFN) | learning rate | $[1e-4, 1e-3]$ | 0.00087 |
| | learning rate decay (5,000 steps) | $[0.7, 1.0]$ | 0.722 |
| | dropout | $\{0.0, 0.1, 0.2, 0.3\}$ | 0.0 |
| | hidden size, d | $\{64, 128, 256, 512\}$ | 512 |
| | layers, l | $\{1, 2, 3\}$ | 2 |
| | batch size | $\{16, 32, 64, 128\}$ | 128 |
| | weight decay | $\{0, 1e-6, 1e-7\}$ | 0 |
| | use differences | $\{\text{True, False}\}$ | True |
| | bin size | - | 0.1 |
| | NEIMS (GNN) | learning rate | $[1e-4, 1e-3]$ |
| learning rate decay (5,000 steps) | | $[0.7, 1.0]$ | 0.767 |
| dropout | | $\{0.0, 0.1, 0.2, 0.3\}$ | 0.0 |
| hidden size, d | | $\{64, 128, 256, 512\}$ | 512 |
| layers, l | | $[1, 6]$ | 4 |

Table 5.10 – continued from previous page

| Model | Parameter | Grid | Value |
|--------------|-----------------------------------|-----------------------------------|----------|
| | batch size | {16, 32, 64} | 64 |
| | weight decay | {0, 1e-6, 1e-7} | 1e-7 |
| | use differences | {True, False} | True |
| | bin size | - | 0.1 |
| | random walk embed steps | [0,20] | 19 |
| | graph pooling | {mean, attention} | mean |
| 3DMolMS | learning rate | [1e-4, 1e-3] | 0.00074 |
| | learning rate decay (5,000 steps) | [0.7, 1.0] | 0.86 |
| | dropout | {0.0, 0.1, 0.2, 0.3} | 0.3 |
| | hidden size, d | {64, 128, 256, 512} | 256 |
| | layers, l | [1, 6] | 2 |
| | top layers | [1, 3] | 2 |
| | neighbors, k | [3, 6] | 5 |
| | batch size | {16, 32, 64} | 16 |
| | weight decay | {0, 1e-6, 1e-7} | 1e-6 |
| | bin size | - | 0.1 |
| MassFormer | learning rate | [1e-4, 1e-3] | 0.00014 |
| | learning rate decay (5,000 steps) | [0.7, 1.0] | 0.85 |
| | feed forward dropout | {0.0, 0.1, 0.2, 0.3, 0.4, 0.5} | 0.1 |
| | feed forward hidden size, d | {64, 128, 256, 512, 1024} | 1024 |
| | feed forward layers, l | [1, 6] | 1 |
| | batch size | {32, 64, 128} | 128 |
| | weight decay | {0, 1e-6, 1e-7} | 1e-7 |
| | bin size | - | 0.1 |
| FixedVocab | learning rate | [1e-4, 1e-3] | 0.00018 |
| | learning rate decay (5,000 steps) | [0.7, 1.0] | 0.92 |
| | dropout | {0.0, 0.1, 0.2, 0.3} | 0.3 |
| | hidden size, d | {64, 128, 256, 512} | 512 |
| | layers, l | [1, 6] | 6 |
| | batch size | {16, 32, 64} | 64 |
| | weight decay | {0, 1e-6, 1e-7} | 1e-6 |
| | bin size | - | 0.1 |
| | random walk embed steps | [0,20] | 11 |
| | graph pooling | {mean, attention} | mean |
| | formula library size | {1000, 5000, 10000, 25000, 50000} | 5000 |
| SCARF-Thread | learning rate | [1e-4, 1e-3] | 0.000577 |
| | learning rate decay (5,000 steps) | [0.7, 1.0] | 0.894 |
| | dropout | {0.0, 0.1, 0.2, 0.3} | 0.3 |
| | hidden size, d | {128, 256, 512} | 512 |
| | mlp layers, l_1 | [1, 3] | 2 |
| | gnn layers, l_2 | [1, 6] | 4 |
| | batch size | {8, 16, 32, 64} | 16 |
| | weight decay | {0, 1e-6, 1e-7} | 1e-6 |

Table 5.10 – continued from previous page

| Model | Parameter | Grid | Value |
|-------------|-----------------------------------|---------------------------|-----------|
| | use differences | {True, False} | True |
| | random walk embed steps | [0,20] | 20 |
| | graph pooling | {mean, attention} | mean |
| SCARF-Weave | learning rate | $[1e - 4, 1e - 3]$ | 0.00031 |
| | learning rate decay (5,000 steps) | [0.7, 1.0] | 0.962 |
| | dropout | {0.0, 0.1, 0.2, 0.3} | 0.2 |
| | hidden size, d | {128, 256, 512} | 512 |
| | mlp layers, l_1 | [1, 3] | 2 |
| | gnn layers, l_2 | [1, 6] | 3 |
| | transformer layers, l_3 | [0, 3] | 2 |
| | batch size | {4, 8, 16, 32, 64} | 32 |
| | weight decay | {0, $1e - 6$, $1e - 7$ } | 0 |
| | bin size | - | 0.1 |
| | random walk embed steps | [0,20] | 7 |
| | graph pooling | {mean, attention} | attention |

Chapter 6

Predicting Enzyme-Substrate Compatibility

While metabolomics and the previously outlined chapters provide strategies for *identifying* metabolites, placing such discovered molecules in context is an orthogonal challenge. Specifically, an understanding of how metabolites progress through metabolic pathways would be broadly useful in considering various biochemical interventions (i.e., where in the pathway to inhibit) and also to repurpose the transforming enzymes for other applications. This chapter explores the ability to predict and model enzyme-substrate interactions in the context of families of enzymes with large amounts of screening data available. The vast majority of this text has previously appeared as S. Goldman, R. Das, K. K. Yang, *et al.*, “Machine learning modeling of family wide enzyme-substrate specificity screens,” *PLoS Computational Biology*, vol. 18, no. 2, e1009853, 2022. I conceptualized this work with my advisor, Connor W. Coley, and conducted all experiments and execution; I wrote the initial manuscript, with support from all listed co-authors.

6.1 Introduction

Biology has evolved enzymes that are capable of impressively stereo-selective, regio-selective, and sustainable chemistry to produce compounds and perform reactions that are “the envy of chemists” [187]–[189]. Industrial integration of these enzymes in catalytic processes is transforming our bioeconomy, with engineered enzymes now producing various materials and medicines on the market today [188], [190]. As an exemplar of collective progress in biocatalysis and enzyme engineering, Huffman et al. impressively re-purposed the entire nucleoside salvage pathway for a high yield, 9-enzyme *in vitro* synthesis of the HIV nucleoside analogue drug, islatravir [191]. In an effort to make these types of pathways commonplace, there has been an explosion in new tools for automated computer-aided synthesis planning (CASP) that can include not only traditional organic chemistry reactions [192], but also enzymatic reactions, facilitating further growth of industrial biocatalysis [193]–[195].

Despite this progress in synthesis planning, suggesting an enzyme for each catalytic step in a proposed synthesis pathway remains difficult and limits the practical utility of synthesis planning software. Current enzyme selection methods often use simple similarity searches,

comparing the desired reaction to precedent reactions in a database [196], [197]. Due to the often high selectivity of enzymes, proposed enzymes for a hypothetical reaction step often suffer from low catalytic efficiency. In the extreme case, the proposed enzyme may have zero catalytic effect on the substrate of interest, despite showing moderate activity on a similar natural substrate. Thus, the specificity of enzymatic catalysis can be a double edged sword [198]. As an example, the phosphorylation step in the islatravir synthesis of Huffman et al. required screening a multitude of natural kinase classes to find an enzyme capable of phosphorylating the desired substrate with sufficient activity for subsequent directed evolution [191]. In the less extreme case, the enzyme of interest may have moderate activity but suffer from low initial substrate loadings, proceed slowly, require higher catalyst loadings, and produce low yields [188]. Nevertheless, an enzyme with moderate activity can serve as a “hook” for further experimental optimization and directed evolution efforts.

Machine learning and predictive modeling provide an avenue to accelerate long development cycles and identify enzymes with both initial activity and high efficiency. Sequence-based machine learning methods have already been utilized in “machine learning guided directed evolution” (MLDE) campaigns to help guide the exploration of sequence space toward desirable protein sequences [199]–[201]. MLDE demonstrations often follow a similar paradigm: given a screen of a single enzyme with mutations at select positions, predict the function or activity of enzymes with new mutations. Further developments in pretrained machine learning models can now provide meaningful embeddings at single amino acid positions that capture contextual and structural information about the protein from sequence alone [202]–[207]. Pre-trained machine learning models of protein sequences, specifically masked language models adopted from natural language processing [32] are trained to predict the identity of “masked” input tokens (i.e. amino acids). In doing so, the model learns a meaningful intermediate representation of the protein and distill important structural context around each amino acid position. This intermediate layer can then be extracted and treated as a fixed embedding of the protein [208]. These pre-trained embeddings have proved especially useful for the application of MLDE in low-N settings where the number of protein measurements is small (Figure 6.1A) [209], [210]. Altogether, these approaches provide a way to improve the efficiency of an enzyme given examples of other enzymes with activity on the substrate of interest. However, this paradigm does not extend to meet the the challenge of identifying a “hook” enzyme with sufficient initial activity on a non-native substrate, nor can current approaches incorporate information from enzymes measured on *other* substrates.

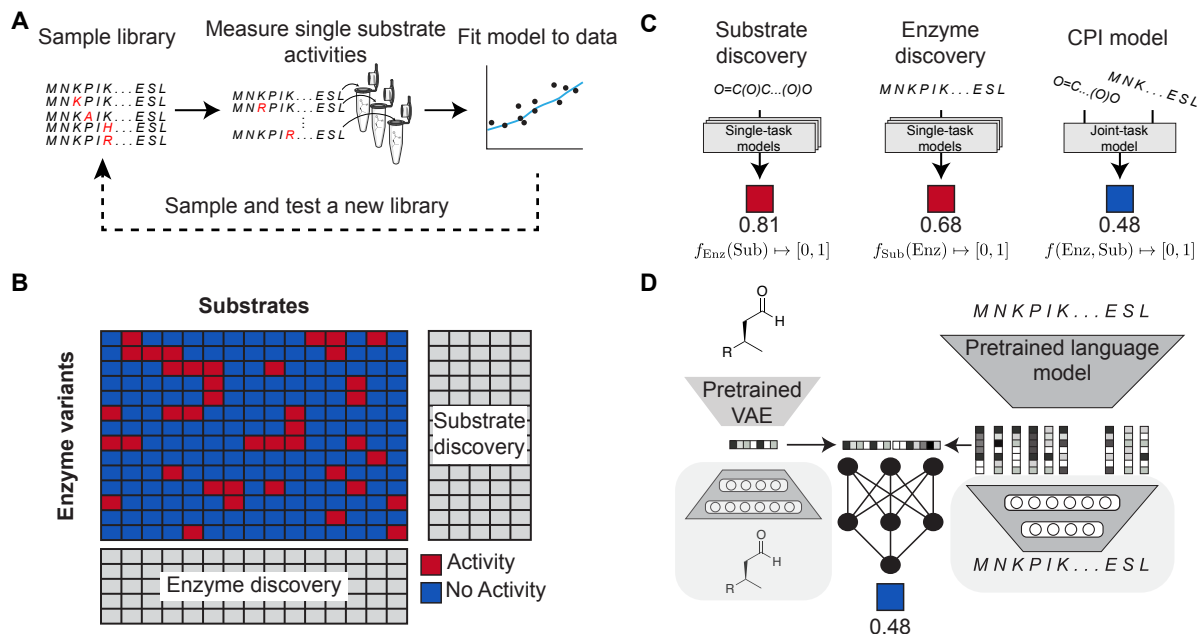


Figure 6.1: **Enzyme-substrate interaction modeling strategies.** (A) Current machine learning-directed evolution strategies, which involve design-build-test-model-learn cycles measuring protein variant activity on a single substrate of interest. (B) The “dense screen” setting where homologous enzyme variants from one protein family are profiled against multiple substrates. In this setting, we can aim to generalize to either new enzymes (“enzyme discovery”) or new substrates (“substrate discovery”). (C) Three different styles of models evaluated in this study, where single task models independently build predictive models for rows and columns from panel (B), whereas a CPI model takes both substrates and enzymes as input. (D) An example CPI model architecture where pretrained neural networks extract features from the substrate and enzyme to be fed into a top-level feed forward model for activity prediction.

Instead, work to date to expand the substrate scope of an enzyme class of interest often relies upon time consuming and *ad hoc* rational engineering based upon structure [211], [212], simple similarity searches between the native substrate and desired substrate of interest [196], or trial-and-error experimental sampling [213]. Once an enzyme with some activity for a substrate of interest is found practitioners can resume directed evolution strategies similar to those described above to increase efficiency [213]–[216].

The last strategy of experimental sampling often involves broad metagenomic sampling [217]–[220], where homologous sequences are chosen for testing [219], [221], [222]. Researchers will test a diverse set of “mined” enzymes for activity against a panel of substrates containing the relevant reactive group. This experimental screening of enzymes against substrates closely mirrors the data setting involved in discovering selective inhibitors in drug discovery, where a panel of similar proteins such as kinases [223] or deubiquinating proteins [224], [225] are screened against a family of compounds. While some work in this field of compound protein interactions (CPI) has attempted to model the drug discovery framing of this problem [226], [227], the CPI modeling framework has not yet been extended to

enzyme promiscuity and there exist few curated datasets to probe our ability to learn from enzyme screens.

In this work, we model enzyme-substrate compatibility as a compound-protein interaction task using a carefully curated set of recent metagenomic enzyme family screens from the literature. We compare state of the art predictive modeling using pretrained embedding strategies (for both small molecules and proteins) and CPI prediction models. Surprisingly, we find that predictive models trained jointly on enzymes and substrates fail to outperform independent, single-task enzyme-only or substrate-only models, indicating that the joint models are incapable of learning interactions. To determine whether this is a quirk specific to our datasets, we reanalyze a recent CPI demonstration and find that this trend generalizes beyond enzyme-substrate data to CPI more broadly: learning interactions from protein family data to go beyond single-task models remains an open problem. Finally, we introduce a new pooling strategy specific to metagenomically-sampled enzymes using a multiple sequence alignment (MSA) and reference crystal structure to enhance enzyme embeddings and improve model performance on the task of enzyme activity prediction. Collectively, this work lays the foundation and establishes dataset standards for the construction of robust enzyme-substrate compatibility models that are needed for various downstream applications such as biosynthesis planning tools.

6.2 Results

6.2.1 Data summary

In order to systematically evaluate our ability to build models over enzyme-substrate interactions, we first need high quality data. Databases of metabolic reactions such as BRENDA [228] describe large numbers of known enzymatic reactions, but are collected from many sources at different concentrations, temperatures, and pH values. Instead, we turn to the literature to find high-throughput enzymatic activity screens with standardized procedures, i.e., exhibiting no variation in the experiments besides the identities of the small molecule and enzyme. We extract amino acid sequences and substrate SMILES strings from six separate studies measuring the activity of halogenase [219], glycosyltransferase [229], thiolase [230], β -keto acid cleavage enzymes (BKACE) [231], esterase [232], and phosphatase enzymes [233] which cover between 1,000 and 36,000 enzyme-substrate pairs (Table 6.1). Enzymatic catalysis (e.g., yield, conversion, activity) in each study was measured using some combination of coupled assay reporters, mass spectrometry, or fluorescent substrate readouts. Data was binarized such that every measured pair is either labeled as active (1) or inactive (0) at thresholds according to standards described in the original papers (Methods). We conceptualize these datasets as “dense screens” insofar as each dataset represents a number of enzyme and substrate pairs measured against each other resembling a grid (Figure 6.1B). While several of the experimental papers presenting these datasets include their own predictive modeling [229]–[231], these demonstrations are not systematically compared with standard splits and are not easily evaluated against new methods due to varied data formats. In compiling these, we expose new datasets to the protein machine learning community. Additional details can be found in Section 6.5.

In order to evaluate the ability of data-driven models to generalize beyond the set of screened enzymes and substrates, we examine extrapolation in two directions: enzyme discovery or substrate discovery (Figure 6.1C). In the former, we consider the setting of a practitioner who is interested in finding enzymes with activity on some set of substrates they have already measured. This parallels the setting where machine learning directed evolution may also be applied, such as increasing the efficiency of an enzyme (Figure 6.1A). On the other hand, for substrate discovery, we are interested in predicting which enzymes from an already-sampled set will act on a new substrate that has not already been measured, a formulation specifically relevant to synthesis planning. We omit the more difficult problem of generalizing to new substrates and new enzymes simultaneously; we posit that we must first be able to generalize in each direction separately in order to generalize jointly and note that empirical performance on joint generalization in compound protein interaction is lower in previous CPI studies [226]. We do not consider the task of interpolation within a dense screen (Figure 6.1B), as this does not reflect any realistic experimental application.

6.2.2 Models

We aim to build a modeling pipeline that accepts both an enzyme and substrate and predicts sequence, with enzymes and substrates specified by sequence and SMILES strings respectively. To featurize enzymes, we turn to pre-trained protein language models, specifically the currently state-of-the-art ESM-1b model [206]. Pre-trained representations are well-suited to low data tasks and have been applied to protein property prediction [206], [207], [209] as well as compound protein interaction [234]. While there has been an explosion in available pre-trained protein representations including UniRep [207], SeqVec [235], MT-LSTM [202], [203], ESM-1b stands out in its performance on contact prediction tasks, ease of use, and also ability to effectively predict the functional effect of sequence variations, likely enabled by its comparatively large scale (i.e., number of parameters and training sequences) [206]. To featurize substrates, we test two primary featurizations: a pretrained Junction-Tree Variational Auto-Encoder (JT-VAE) [236] and the widely used Morgan circular fingerprints (1024 bits) [237]. Despite pre-training compound representations having only marginal benefits on property prediction tasks [238], a recent CPI study [226] extracted compound representations from a pre-trained JT-VAE model and utilized these to identify new kinase inhibitors. Due to the closeness in our proposed task, we follow their methodology and extract the same embeddings for substrates to compare against Morgan fingerprints.

If a single model is able to successfully learn interactions and leverage the full “dense” dataset, it should be able to take as input both enzyme and substrate representations and outperform smaller, single-task models that either use only enzyme inputs or use only substrate inputs (Figure 6.1C). To attempt to model interactions, we consider two simple top models inspired by the CPI literature [226], [227], [234] that either (1) concatenate the representations of the enzyme and substrate before applying a shallow feed forward neural network or (2) project the representations of the enzyme and substrate to smaller and equal length vectors using a shallow multi-layer perceptron (MLP) before taking their dot product (Figure 6.1D). To evaluate these CPI based architectures, we consider three other model classes for comparison: baselines utilizing simple similarity across enzymes and substrates for prediction; multi-task models that learn to predict activity for enzymes against substrates

Table 6.1: **Summary of curated datasets with the number of unique enzymes, unique substrates, and unique pairs in each dataset in addition to an exemplar structure for the protein family.**

| Dataset | # Enz. | # Sub. | Pairs | PDB Ref. |
|---------------------------|--------|--------|--------------------|----------|
| Halogenase [219] | 42 | 62 | 2,604 | 2AR8 |
| Glycosyltransferase [229] | 54 | 90 | 4,298 ^a | 3HBF |
| Thiolase [230] | 73 | 15 | 1,095 | 4KU5 |
| BKACE [231] | 161 | 17 | 2,737 | 2Y7F |
| Phosphatase [233] | 218 | 165 | 35,970 | 3L8E |
| Esterase [232] | 146 | 96 | 14,016 | 5A6V |
| Kinase (inhibitors) [223] | 318 | 72 | 22,896 | 2CN5 |

^a While most datasets we use test all combinations, Yang et al. do not report experiments for some enzyme by substrate interactions

simultaneously but without any feature information about the substrates themselves (and *vice versa* for substrate discovery); and single-task models with no information sharing across substrate (enzyme) tasks. We refer the reader to Section 6.6 for a more complete description of all model classes evaluated (Table 6.4).

6.2.3 Enzyme discovery

We test these various featurizations and model architectures first on the task of enzyme discovery. To do so, we hold out a fraction of the enzymes as a test set and use the training set to make predictions about the interactions between the held out enzymes and the known substrates in the data set. For each dataset, we train the CPI model architectures described above jointly on the entire training set. To test whether CPI models are able to learn interactions, we also train several smaller “single-task” models. These single-task models are specific to each substrate and accept only an enzyme sequence as input. If models are in fact able to learn interactions, the CPI models should outperform the single-task models given their access to more substrate measurements for each enzyme (Figures 6.2A and 6.1C).

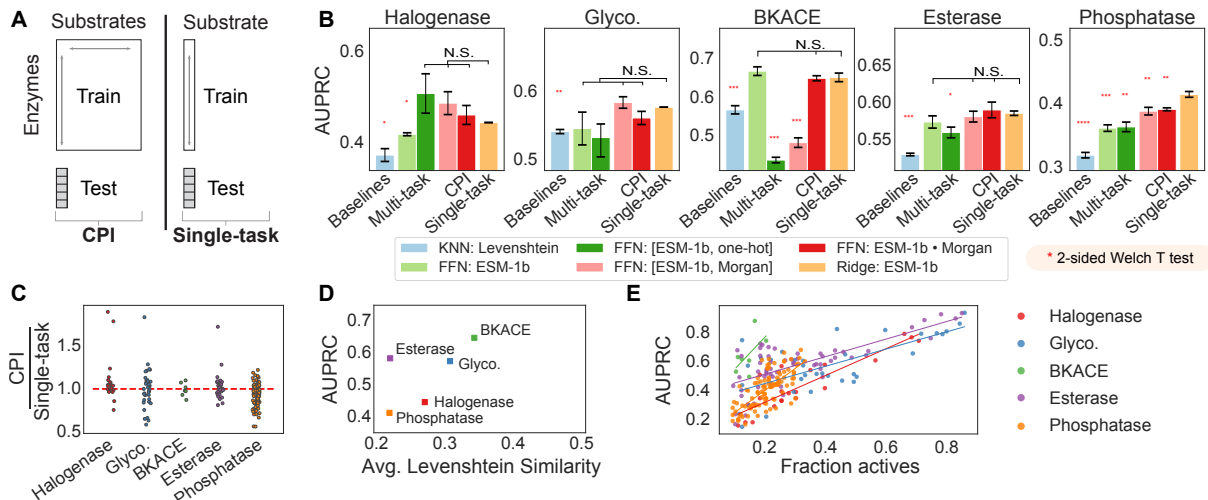


Figure 6.2: **Assessing enzyme discovery in family wide screens.** **(A)** CPI models are compared against the single task setting by holding out enzymes for a given substrate and allowing models to train on either the full expanded data (CPI) or only data specific to that substrate (single-task). **(B)** AUPRC is compared on five different datasets, arranged from left to right in order of increasing number of enzymes in the dataset. Baseline models are compared against multi-task models, CPI models, and single-task models. K-nearest neighbor (KNN) baselines are calculated using Levenshtein edit distances to compare sequences; multi-task models use a shared feed forward network (FFN) to compute predictions against all substrate targets, CPI models utilize FFN with either concatenation (“{prot repr.}, {sub repr.}”]) or dot product interactions (“{prot repr.}•{sub repr.}”), and ridge regression is used for single-task models. ESM-1b features indicate protein features extracted from a masked language model trained on UniRef50 [206]. Halogenase and glycosyltransferase datasets are evaluated using leave-one-out splits, whereas BKACE, phosphatase, and esterase datasets are evaluated with 5 repeats of 10 different cross validation splits. Standard error bars indicate the standard error of the mean of results computed with 3 random seeds. Each method is compared to the single-task L2-regularized logistic regression model (“Ridge: ESM-1b”) using a 2-sided Welch T test, with each additional asterisk representing significance at [0.05, 0.01, 0.001, 0.0001] thresholds respectively after application of a Benjamini-Hochberg correction. **(C)** Average AUPRC on each individual “substrate task” is compared between compound protein interaction models and single-task models. Points below 1 indicate substrates on which single-task models better predict enzyme activity than CPI models. CPI models used are FFN: [ESM-1b, Morgan] and single-task models are Ridge: ESM-1b. **(D)** AUPRC values from the ridge regression model are plotted against the average enzyme similarity in a dataset, with higher enzyme similarity revealing better predictive performance. **(E)** AUPRC values from the ridge regression model broken out by each task are plotted against the fraction of active enzymes in the dataset. Best fit lines are drawn through each dataset to serve as a visual guide.

To evaluate each dataset, we calculate the area under the precision recall curve (AUPRC), computed separately for each substrate column and subsequently averaged. AUPRC is able to better differentiate model performance on highly imbalanced data than the area under

the receiver operating curve (AUROC), which overvalues the prediction of true negatives. Further, AUPRC does not require choosing a threshold to call hits like other metrics like the Matthews Correlation Coefficient. We optimize model hyperparameters on the thiolase dataset [230] prior to training and evaluating on the remaining five datasets. We additionally report benchmarking performance for the thiolase dataset (Tables 6.6 and 6.7 in Section 6.6).

We observe that our supervised models using pretrained protein representations are in fact able to outperform a nearest neighbor sequence-similarity baseline (“KNN: Levenshtein”) that uses the Levenshtein distance, a simple unweighted global alignment distance used in recent protein engineering studies [44], [209], to predict held out enzyme activity (Figure 6.2B and Tables 6.6 and 6.7 in Section 6.6). This affirms the potential of representation learning to improve prediction and protein engineering tasks.

Surprisingly, however, CPI models do not outperform single-task models trained with simple logistic regression (“Ridge: ESM-1b”) (Figure 6.2B). Multi-task models offer a slight benefit on the halogenase dataset, but fail to outperform single-task models across the other four enzyme families tested. Further, upon closer inspection, the comparative performance of CPI based models on the halogenase dataset seem to be driven by relative performance increases only on a small number of substrate tasks as demonstrated by the upper outliers in Figure 6.2C. This is despite the CPI (and multi-task) models having access to a larger number of enzyme-substrate interactions for training. In fact, models trained with CPI can at times perform worse than models that predict the activity of enzymes on each substrate task independently (Figure 6.2C), indicating an inability to learn interactions from the dense screens collected.

The enzymes within each dataset were sampled with different levels of diversity by the studies’ original authors. The phosphatase dataset represents a diverse super-family of enzymes [233], whereas the BKACE dataset represents a more narrowly sampled domain of unknown function (unknown prior to the experimental screen) [231]. We find that the average pairwise Levenshtein similarity between sequences in the dataset is in fact correlated with performance differences across datasets, such that more similar datasets seem to be easier to predict (Figure 6.2D).

In addition to intra-dataset diversity, we also hypothesized that the balance, or fraction of active enzymes, observed for each enzyme could partially explain the observed performance. Plotting the AUPRC metric as a function of number of active enzyme-substrate pairs reveals a strong positive correlation, validating that the number of hits observed in the training set will largely determine the success of the model in generalizing beyond the training set (Figure 6.2E). This is equally a function of both the models and the AUPRC metric, which follows a similar trend for random guesses that favor the majority binary class.

6.2.4 Substrate discovery

We next evaluate generalization in the direction of held out substrates, repeating the same procedures as above. In this case, we restrict our analysis only to the glycosyltransferase and phosphatase datasets where the number of substrates is > 50 , using the halogenase dataset to tune hyperparameters for each model. We report the results comparison on the halogenase dataset (Tables 6.9 and 6.8 in Section 6.6).

Similar to our conclusions in the case of enzyme discovery, we find that the CPI architectures are not able to outperform simpler, single-task logistic regression models with Morgan fingerprints (Figure 6.3 and Figure 6.10 in Section 6.6 and Tables 6.8 and 6.9) in Section 6.6.

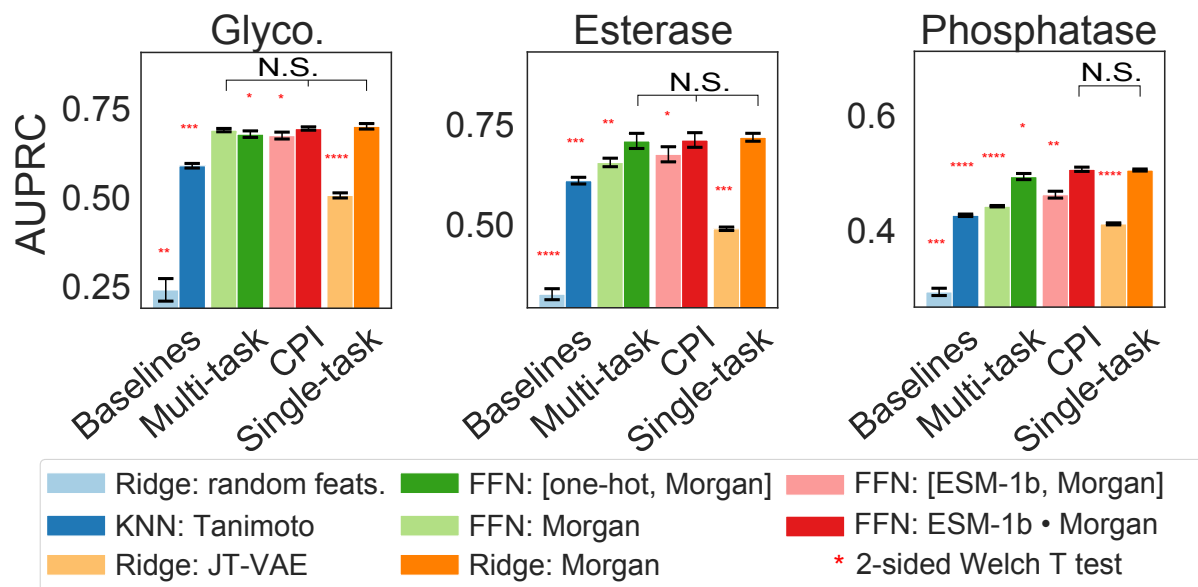


Figure 6.3: **Assessing substrate discovery in family wide screens.** CPI models and single task models are compared on the glycosyltransferase, esterase, and phosphatase datasets, all with 5 trials of 10-fold cross validation. Error bars represent the standard error of the mean across 3 random seeds. Each model and featurization is compared to “Ridge: Morgan” using a 2-sided Welch T test, with each additional asterisk representing significance at [0.05, 0.01, 0.001, 0.0001] thresholds respectively, after applying a Benjamini-Hochberg correction. Pretrained substrate featurizations used in “Ridge: JT-VAE” are features extracted from a junction-tree variational auto-encoder (JT-VAE) [236]. Two compound protein interaction architectures are tested, both concatenation and dot-product, indicated with “{prot repr.}, {sub repr.}” and “{prot repr.}•{sub repr.}” respectively. In the interaction based architectures, ESM-1b indicates the use of a masked language model trained on UniRef50 as a protein representation [206]. Models are hyperparameter optimized on a held out halogenase dataset. AUCROC results can be found in Figure 6.9 in Section 6.6.

Curiously, in both the enzyme discovery (Figures 6.14, 6.15, 6.13, 6.12, and 6.16) in Section 6.6 and substrate discovery (Figures 6.18, 6.17, and 6.19 in Section 6.6) settings, predictions made by CPI models exhibit far more “blocky” characteristics than the respective single task models: when extrapolating to new enzymes, the prediction variance is not sensitive to the paired substrate for CPI models. This indicates that our CPI models struggle to condition their predictions to new enzymes (substrates) based upon the substrate (enzyme) pairing.

6.2.5 Reanalysis of kinase inhibitor discovery

The results for enzyme-substrate activity prediction demonstrate that models designed to learn interactions are seemingly unable to do so in a manner that improves generalization. We therefore wondered to what extent this failing was specific to enzyme-substrate data, as opposed to being symptomatic of a broader problem and shortcoming in the CPI field, including drug discovery. To interrogate this, we re-analyze models from a recent study leveraging an inhibitor screen against the human kinome to discover new inhibitors against tuberculosis [226]. In their study, Hie et al. train CPI models on a dense screen of 442 kinases against 72 inhibitors [223] using concatenated pretrained protein and pretrained substrate features as the input to multi layer perceptrons (MLP), Gaussian processes (GPs), or a combination of the two (GP + MLP), the combination being their most successful (Methods). Unlike the binary classification enzyme activity setting, they predict continuous K_d values.

We compare the MLP and GP+MLP models using pretrained representations against a number of single-task models on two settings matching the original study: drug repurposing and drug discovery. Drug repurposing is analogous to enzyme discovery where certain proteins are held out; drug discovery is analogous to substrate discovery where certain compounds are held out. Single-task models are not presented with training data on other kinase-compound pairs and are therefore unable to learn interactions in a generalizable manner. In addition to the single-task MLP and GP+MLP models, we evaluate a simple single-task, L2-regularized linear regression model (“Ridge”) using Morgan fingerprint features rather than JT-VAE features. In concordance with our results on enzymatic data, we find that single-task models consistently outperform CPI based models in terms of Spearman correlation coefficient between true and predicted K_d on both repurposing (Figure 6.4A) and discovery tasks (Figure 6.4B). This shows that ablating interactions by training single-task models can increase performance over GP+MLP models.

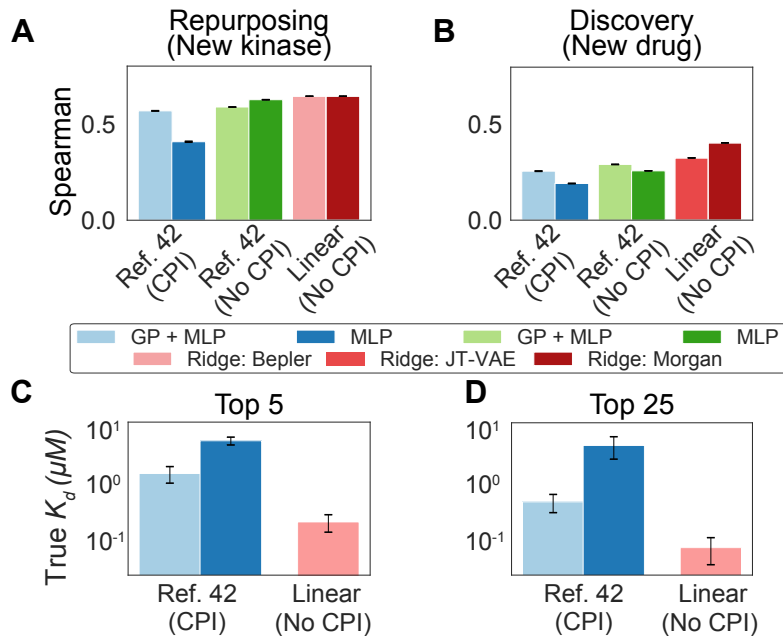


Figure 6.4: **Evaluating single-task models on kinase repurposing and discovery tasks** Kinase data from Davis et al. is extracted, featurized, and split as prepared in Hie et al. Multilayer perceptrons (MLP) and Gaussian process + multilayer perceptron (GP+MLP) models are employed. We add variants of these models without CPI training separate single-task models for each enzyme and substrate in the training set, as well as linear models using both pretrained featurizations (“Ridge: JT-VAE”) and fingerprint based featurizations of small molecules (“Ridge: Morgan”). Spearman correlation is shown for **(A)** held out kinases not in the training set and **(B)** held out small molecules not in the training set across 5 random initializations. **(C)** We repeat the retrospective evaluation of lead prioritization. The top 5 average acquired K_d values are shown for the CPI models in Hie et al. compared against a linear, single-task ridge regression model using the same features. **(D)** The top 25 average acquired K_d values are shown.

Still, increased rank correlation between predictions and true K_d values does not necessarily equate to the ability to select new inhibitors or new drugs. To directly test this, we repeat the retrospective kinase-inhibitor lead prioritization experiments conducted in the original analysis. Models are trained on a set of kinase-inhibitor pairings and used to rank new kinase-inhibitor pairings. The acquisition preference is informed by predictions and, if applicable, predicted uncertainty (Methods). When acquiring either 5 or 25 new data points in cross validation, a single-task ridge regression model with equivalent pretrained features is able to outperform both CPI based models (Figure 6.4C and 6.4D). Our findings are retrospective in nature and do not negate the value of prospective experimental validation [226], but rather make clear that the field requires new methods to leverage the rich information contained in protein, small molecule interaction screens and truly learn interactions. This further reinforces the necessity for simple baseline models in protein engineering studies [239], [240].

6.2.6 Improving enzyme discovery models

Given that single-task models appear to match or even outperform models design for CPI, we next asked if we could improve their generalization in the enzyme discovery direction by leveraging the relationship between different protein sequences within the dataset. That is, working within a single family of proteins *should* be more conducive to generalizations. To directly impart this structural bias on our models, we considered how the construction of the pretrained representation for each protein could be modified. Pretrained language models produce a fixed dimensional embedding at each amino acid position in the protein. To collapse this into a fixed-length protein-level embedding, the *de facto* standard is to compute the mean embedding across the length of the sequence [202], [203], [206], [207], [209]. However, for locally-defined properties, such as enzymatic catalysis or ligand binding at an active site, this mean pooling strategy may be sub-optimal [241]. Previous approaches have largely considered deep mutational scans with few mutations at carefully selected positions [200]. In these settings, mean pooling strategies may be a good approximation of local protein structural changes, as embeddings at distal positions from the mutation would be nearly constant across protein variants. In our setting, however, we have metagenomically sampled sequences with large insertions and deletions, which presents an ideal testing ground to evaluate pooling strategies.

We test 3 alternative pooling strategies to mean pooling, where we first compute a multiple sequence alignment (MSA) and pool only a subset of residues in each sequence corresponding to a subset of columns in the MSA (Figure 6.5A). We rank order the columns in the MSA to be pooled based upon the (i) proximity to the active site of a single “reference structure” (Section 6.5.7) (ii) coverage (i.e., pooling columns with the fewest gaps), or (iii) conservation (i.e., pooling columns that have the highest frequency of any single amino acid type) (Table 6.1). To expand this analysis beyond catalysis to drug discovery, we also consider a portion of the kinase inhibitor dataset from Davis et al. [223], subsetted down to a single kinase family (PF00069), rather than the whole human kinome (Table 6.1). We compare these pooling strategies across the enzyme discovery datasets tested, as well as the protein-inhibitor kinase dataset. We use ridge regression models with pretrained ESM-1b [206] embeddings, and split the data as in the enzyme discovery setting, varying only the pooling strategy from our previous analysis (Figure 6.2). In the case of the kinase regression dataset, we use the Spearman rank correlation to evaluate performance.

In all cases tested, the active site pooling performance peaks when pooling only a small number of residues around the active site (< 60 amino acids), showing gains in performance over other pooling strategies as well as the mean pooling baseline (Figure 6.5B). This corresponds to a distance of < 10 angstroms away from the active site (Figure 6.5Ai). This may indicate an optimal range at which residue positioning is relevant to promiscuity. In the case of the kinases, the Levenshtein distance baseline outperforms the mean pooling method with respect to Spearman rank correlation, but using active site aware pooling outperforms both. We find that the performance increase is steeper in the kinase repurposing setting compared to other datasets, potentially due to both the regression nature of the dataset and also the non-dynamic binding of compound inhibitors in comparison to the other tasks, which focus on enzymatic catalysis not protein inhibition. Interestingly, for the halogenase dataset, no alternative pooling strategies outperform mean pooling (Figure 6.11 in Section 6.6), likely

because the halogenase enzymes have high variance in solubility, a global property that could be driving enzyme activity [219]. Similarly, for the esterase dataset, coverage pooling is far more effective, indicating that a combination of targeted pooling residue strategies may be most effective (Figure 6.11 in Section 6.6). While performance gains from active site aware pooling are modest, this strategy provides a simple but principled way to incorporate a structural prior into enzyme prediction models, particularly for metagenomic data with high numbers of sequence indels which may introduce unwanted variance into a mean pooled protein representation.

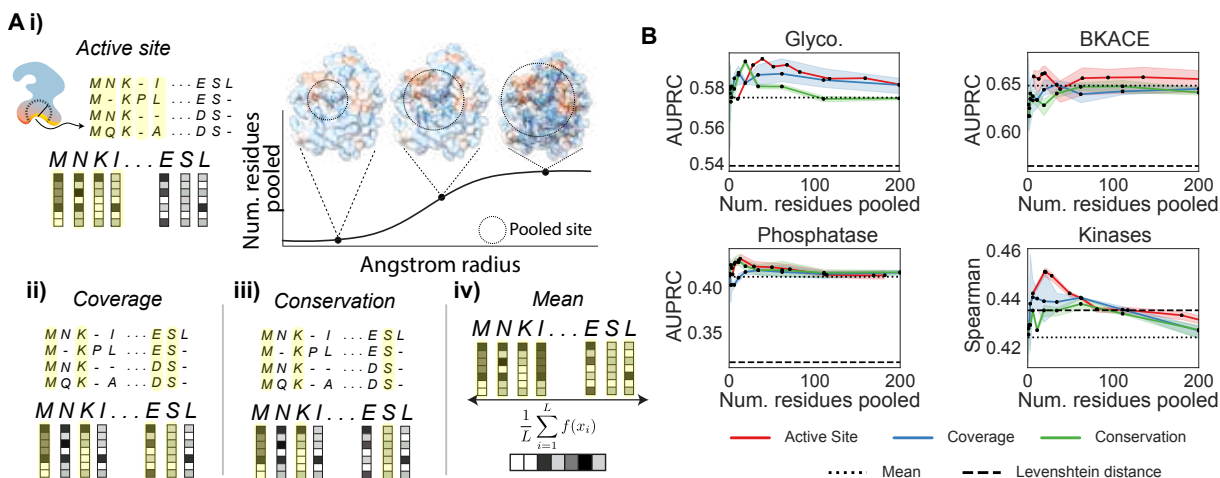


Figure 6.5: **Structure-based pooling improves enzyme activity predictions.** (A) Different pooling strategies can be used to combine amino acid representations from a pre-trained protein language model. Yellow coloring in the schematic indicates residues that will be averaged to derive a representation of the protein of interest. (i) We introduce active site pooling, where only embeddings corresponding to residues within a set radius of the protein active site are averaged. By increasing the angstrom radius from the active site, we increase the number of residues pooled. Crystal structures shown are taken from the BKACE reference structure, PDB: 2Y7F rendered with Chimera [242]. (ii, iii) We also introduce two other alignment based pooling strategies: coverage and conservation pooling average only the top- k alignment columns with the fewest gaps and highest number of conserved residues respectively. (iv) Current protein embeddings often take a mean pooling strategy to indiscriminately average over all sequence positions. (B) Enzyme discovery AUPRC values are computed for various different pooling strategies. Each strategy is tested for different thresholds of residues to pool, comparing against both KNN Levenshtein distance baselines and a mean pooling baseline. The same hyperparameters are used as set in Figure 6.2 for ridge regression models. The kinase repurposing regression task from Hie et al. is shown with Spearman’s ρ instead of AUPRC as interactions are continuous, not binarized. All experiments and are repeated for 3 random seeds.

6.3 Discussion

Data-driven models of enzyme-substrate compatibility have the potential to drive new insights in basic biology research and also to accelerate engineering efforts focused on the design of new enzymatic synthesis routes. In addition, the same classes of models can be used for compound-protein interaction prediction for both drug discovery and drug repurposing efforts.

In this work, we take a critical step toward opening up this suite of problems to machine learning researchers by providing several high quality, curated datasets and standardized splits to evaluate model performance and generalizability. While the small number of unique enzymes and unique substrates in each dataset makes quantitative performance sensitive to hyperparameters and dataset splitting decisions, this collection of data is an essential starting point to develop new modeling strategies and motivate future, higher throughput enzyme activity screening.

Our experiments show that pretrained representations for proteins, coupled with structure-informed pooling techniques, can go beyond standard sequence similarity based approaches to predict protein function, an exciting demonstration of how representation machine learning can impact protein engineering. Nevertheless, despite this excitement, our analysis makes clear that current CPI modeling strategies cannot consistently leverage information from multiple substrate measurements effectively, a problem broadly applicable to CPI models. That is, models designed to learn interactions do not outperform single-task models.

6.4 Conclusion

To predict enzyme-substrate compatibility or design selective inhibitors against a protein family, we need new strategies to jointly embed proteins and compounds to enable more robust extrapolation to new combinations thereof. Such a scheme would allow learned interactions to be more explicitly transferred from larger databases onto smaller, but higher quality screen. This will be an exciting frontier in protein and compound representation learning, as the field seeks to go beyond protein structural prediction to protein function prediction. Further, with the exception of our structure informed (MSA-informed) pooling, our analysis remains sequence based. The relative performance of structure-based tools such as molecular dynamics for the prediction of enzyme-substrate scope remains an exciting question that this data, coupled with recent advances in protein structure prediction [35], [37], can help to address.

6.5 Methods

6.5.1 Dataset preparation

Each dataset is collected from their respective papers [219], [223], [229]–[233]. Activity binarizations are chosen to closely mirror the original dataset preparation with exact cutoff

thresholds described in Section 6.6. Additionally, certain enzymes were filtered based upon low solubility or activity that may result from screening decisions (Section 6.6).

Davis kinase filtering

Kinases used in reanalysis of Hie et al. are tested exactly as prepared [226]. To evaluate structure based pooling using this dataset, we further subset the original dataset such that each entry only contains one domain from the PFAM family, PF00069, described in the SI with dataset statistics in Table 6.1.

6.5.2 Hyperparameter optimization

All hyperparameters are set on a held out enzyme-substrate dataset using the hyperparameter optimization framework Optuna [123]. Hyperparameter optimization is set using up to 10 trials of leave-one-out cross validation on the thiolase dataset and halogenase dataset for enzyme discovery and substrate discovery tasks, respectively. Hyperparameters are chosen to maximize the average area under the precision recall curve. For nearest neighbor models, the number of neighbors is treated as a hyperparameter between 1 and 10.

For logistic ridge regression models, the regularization coefficient, α is set from $\{1e-3, 1e-2, 1e-1, 1e0, 1e1, 1e2, 1e3, 1e4\}$. For both feed-forward dot product and concatenation models, hyperparameters for dropout ($[0, 0.2]$), weight decay ($[0, 0.01]$), hidden dimension ($[10, 90]$), layers ($[1, 2]$), and learning rate ($[1e-5, 1e-3]$) are chosen. All neural network models are trained for 100 epochs using the Adam optimizer and Pytorch [174].

For linear ridge regression used in reanalysis of kinase data, a default hyperparameter regularizer value of $\alpha = 1e1$ is set.

6.5.3 Evaluation metrics

To evaluate models in the enzyme discovery direction, activity on each substrate is considered to be its own “task”. The data is divided up into a set number of folds, and models are re-trained to make predictions on each held out fold. A single, separate AUPRC value is computed for the activity on each substrate task and then averaged across substrate tasks. AUPRC values are computed using the average precision function from `sklearn` [243]. For the halogenase, thiolase, and glycosyltransferase datasets, this is done with leave-one-out cross validation. For the phosphatase, BKACE, and kinase datasets, to limit the number of trials, we use 10 fold cross validation repeated 5 times. This procedure is repeated for 3 random seeds.

An identical procedure is conducted on the task of substrate discovery, where each enzyme is separately evaluated as its own “task”. In this case, the glycosyltransferase, esterase, and phosphatase datasets are evaluated with 5 repetitions of 10 fold cross validation.

6.5.4 Filtering imbalanced tasks

Certain enzymes have activity on only a few substrates and certain substrates have activity on only a few enzymes. To avoid computing AUPRC values on these tasks, we filter to

only incorporate tasks with a maximum fraction of either 0.9 positives or negatives. The remaining tasks can be found in Table 6.3 in Section 6.6.

6.5.5 Kinase inhibitor reanalysis

We modify the code from Hie et al. directly to reproduce their GP, GP + MLP models, and add our no-interaction models. GP + MLP models involve first fitting an MLP model followed by a GP to predict residual loss. First, all kinases are converted into features using a pretrained language model [202], [203] and all inhibitors are converted into features using a pretrained JT-VAE [236] or Morgan fingerprints. We create a training set of kinase, inhibitor pairs with labeled K_d values, and establish 3 separate segments of the test data: new kinases (repurposing), new inhibitors (discovery), and new kinase+inhibitor pairs. The data is split into four quadrants and one quadrant is used for training models. GP and linear models are implemented with `scikit-learn` [243] and MLP models are implemented with `Keras` [244], following parameter choices from the original study [226]. Linear regression models are parametrized with $\alpha = 10$ and normalization set to `True`. Prior to training single-task models, we standardize the regression target values based upon the training set to have a mean of 0 and variance of 1, as we find it helps with stability with less training data.

To test the ability of models to prioritize candidates, we repeat the train/test split and rank the entire test set by predicted K_d , using an additional upper confidence bound ($\beta = 1$) metric to adjust rank for the GP + MLP model that uses uncertainty. The top $k = 5$ and $k = 25$ compound-kinase pairs are evaluated by their average true K_d . All kinase-inhibitor reanalysis experiments were repeated for five random seeds.

6.5.6 Pooling strategies

We use the program `Muscle` with default parameters to compute a multiple sequence alignment (MSA) on each dataset for pooling. Because many positions in certain datasets have high coverage, ties are randomly broken when choosing a priority for pooling residue embeddings. This explains the large variance in the glycosyltransferase results shown (Figure 6.5) taken over several seeds. Coverage and conservation based pooling strategies are sampled for $i \in \{1, 2, 3, 6, 11, 19, 34, 62, 111, 200\}$ pooling residues, each repeated for 3 random seeds.

6.5.7 Active site pooling

To pool over active sites of proteins, we identify a reference crystal structure within each protein family or super family (Table 6.1). For each of these crystal structures, we select either an active site bound ligand or active site residue(s) from the literature. We attempt to pool residues within a Cartesian distance from these sites ranging from 3 to 30 angstroms to roughly mirror the number of residues pooled for coverage and conservation methods. Angstrom shells are calculated using `Biopython` [245] and an in depth description of active sites used can be found in Table 6.2 in Section 6.6.

6.6 Additional results

6.6.1 Data

Halogenase data

Halogenase data was prepared as described by Fisher et al. They measure the activity of 87 different proteins against 62 substrates using high throughput LC-MS based screening [219]. However, many sampled enzymes are either insoluble or display no halogenation activity. We subset the proteins to a smaller set of 42 proteins that have some halogenation activity on at least one of the substrates tested. We binarize data at the 8% conversion threshold, which Fisher et al. report as removing false positives. Further, rather than test activity on both the chlorination and bromination activity labels, we opt to use only the dataset measuring bromination, which has more a higher percentage of active conversions.

SMILES strings are extracted from the ChemDraw file provided by Fisher et al. All protein sequences with greater than 1000 amino acids were filtered from the enzyme dataset.

Phosphatase data

Phosphatase data is extracted from SI tables provided by Huang et al. Compounds listed in the results using common names are converted using a combination of PubChem’s name converter, cirpy, and manual re-drawing according to compounds in the SI [128], [246].

Sequence IDs are mapped to amino acid sequences using the UniProt database. All entries that are no longer valid are identified using the UniParc database [247]. Enzyme-substrate hits are called at a binary threshold cutoff of 0.2 OD as described in the original paper to correct for background noise.

BKACE data

β-ketoacid cleavage enzyme (BKACE) substrates are manually re-drawn and SMILES strings are extracted from ChemDraw [231]. Enzyme sequences are extracted from the SI, and all hits are binarized according to original procedure from Bastard et al. using a mixture of Gaussians.

Thiolase data

Binary thiolase data is used and extracted as prepared by Robinson et al. and binarized at a threshold of 1×10^{-8} [230].

Esterase data

Binary esterase data is used and extracted as prepared by Martínez-Martínez et al. All enzymes that display > 0 activity are considered to be hits after binarization [232].

Table 6.2: **Active site structure references used in pooling.** All structure informed pooling strategies require a catalytic center in order to define various angstrom shells of residues to pool over. This table provides the PDB reference crystal structure as well as the reference residues or structural elements used to define the pooling center, from which spherical radii originate.

| Dataset | PDB Ref. | Ref. type | Ref. |
|---------------------------|----------|---------------------|--|
| Halogenase [219] | 2AR8 | ligand | 7-chlorotryptophan |
| Glycosyltransferase [229] | 3HBF | ligands | UDP and 3,5,7-TRIHYDROXY-2-(3,4,5-TRIHYDROXYPHENYL)-4H-CHROMEN-4-ONE |
| Thiolase [230] | 4KU5 | catalytic residue | C143 |
| BKACE [231] | 2Y7F | ligand | (5S)-5-amino-3-oxo-hexanoic-acid |
| Phosphatase [233] | 3L8E | ligand | acetic acid |
| Esterase [232] | 5A6V | catatlytic residues | S105, D187D, H224 |
| Kinase (inhibitors) [223] | 2CN5 | ligand | ADP |

Glycosyltransferase data

Glycosyltransferase acceptors and donors were originally measured and classified as having no, intermediate, or strong activity using a “green”, “amber”, or “red” classification system [229]. We make the simplifying assumption to treat all intermediate activity enzymes as a positive example in our binary classification formulation. Further, many more glycosyltransferase acceptor substrates are tested than donors, and so we choose to predict activity of glycosyltransferase-glycosyl acceptor substrate pairs. We use ChemDraw to extract acceptor substrate SMILES strings.

Kinase data

The kinase data in this study is a panel of inhibitors screened against kinases, originally collected by Davis et al. [223]. To compare models directly to Hie et al., we use identical data preprocessing and featurization [226].

For the analysis of structure based pooling, we further processed this dataset for consistency with family-wide protein screens. We subset the data to represent a single PFAM family, PF00069 [248], and we use the `hmmsearch` tool to identify all proteins that satisfy this domain [249].

Due to the large size of the kinase proteins, individual kinase domains were experimentally cloned separately. Some protein entries have multiple measurements corresponding to the first and second kinase domains within the protein. To account for this, we use the envelope returned from `hmmsearch` to subset each protein down to its relevant domain(s). In the interest of maintaining a dataset without point mutations, we further remove all proteins that have specific deletions or insertions. Finally, all kinase-inhibitors without a measured K_d are given a default value of 10,000 as set by Hie et al.

Table 6.3: **Summary of valid substrate and sequence “tasks”**. In each dataset, only certain substrates and sequences are defined as valid “tasks” based upon the balance between active and inactive examples. Each substrate or sequence used for an enzyme or substrate discovery task respectively requires at least 2 positive examples and at a minimum, 10% of examples in that task must be part of the minority class. This table defines the number of valid substrate and sequence tasks.

| Dataset | Num entries | # Seqs. | # Subs. | Valid subs. | Valid seqs. |
|-------------|-------------|---------|---------|-------------|-------------|
| Thiolase | 1095 | 73 | 15 | 11 | 70 |
| Halogenase | 2604 | 42 | 62 | 20 | 17 |
| BKACE | 2737 | 161 | 17 | 7 | 54 |
| Glyco. | 4347 | 54 | 91 | 35 | 48 |
| Esterase | 14016 | 146 | 96 | 59 | 106 |
| Phosphatase | 35970 | 218 | 165 | 103 | 102 |

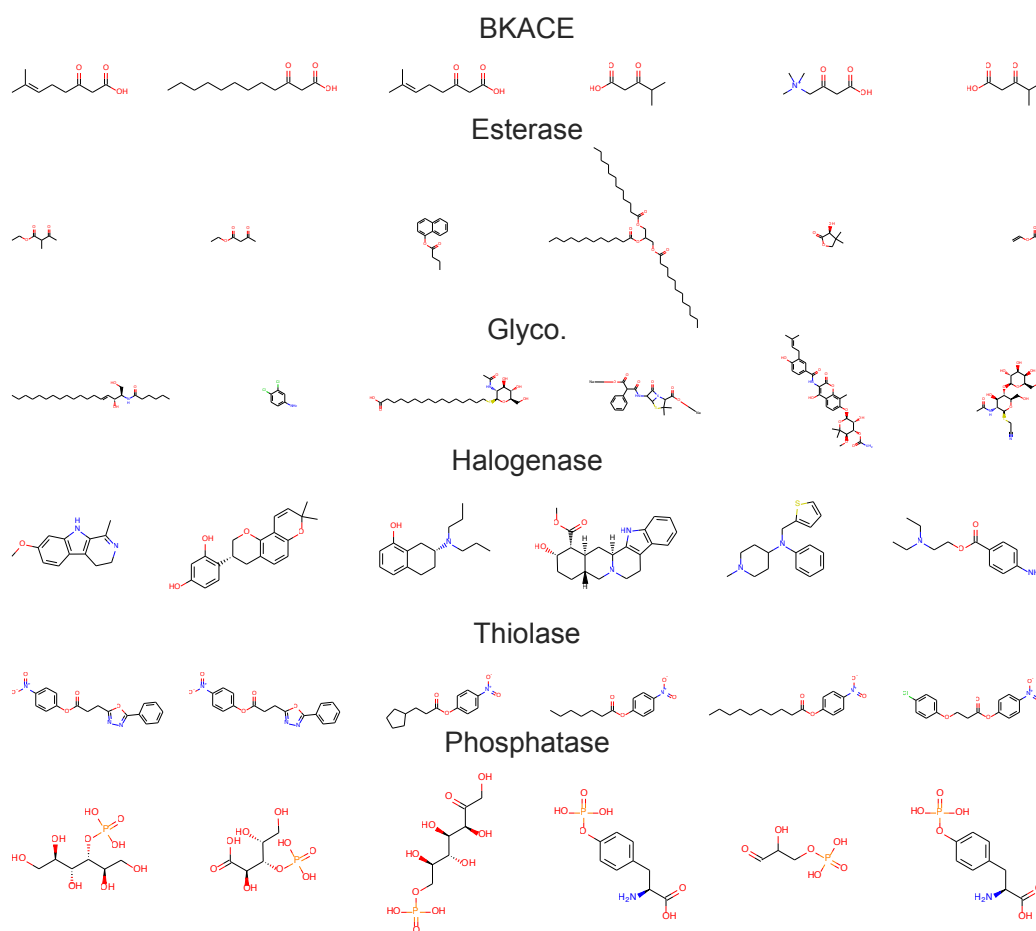


Figure 6.6: **Dataset substrates** 6 exemplar molecule substrates are randomly chosen from each dataset and displayed.

Enzyme and Substrate Dataset Diversity

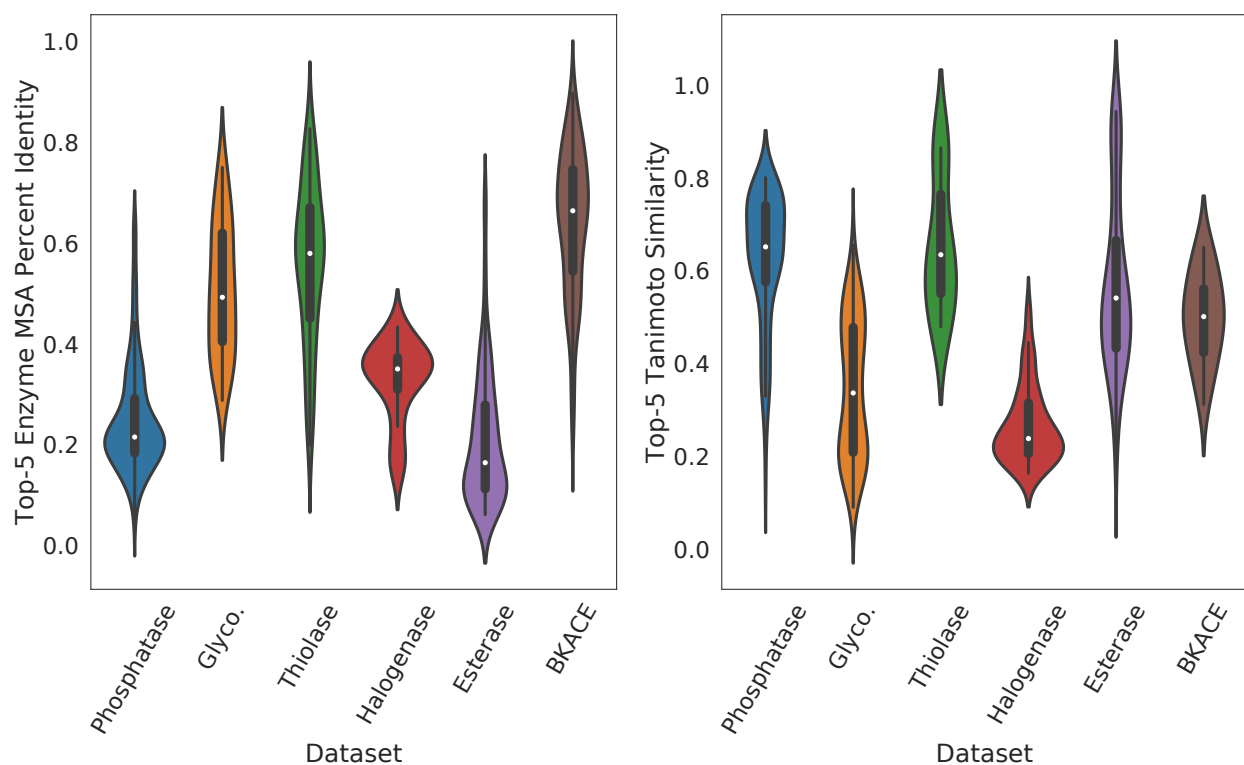


Figure 6.7: **Dataset diversity** Distributions of top-5 enzyme similarity (left) and substrate similarity (right) are shown across enzyme datasets collected. Enzyme similarity is calculated as the percent overlap between two sequences in their respective multiple sequence alignment, excluding positions where both sequences contain gaps. Substrate similarity is computed using Tanimoto similarity between 2048-bit chiral Morgan fingerprints.

6.6.2 Models

Enzyme-substrate models

We consider 4 classes of models to learn from family-wide enzyme-substrate models:

1. **Baseline:** For simple baseline models, we consider non-parametric approaches such as K-nearest neighbors using sequence edit distance and Tanimoto similarity. In addition, we consider the case where we have random features, rather than meaningful ones. With such inputs, we expect models will be able to learn *only* statistical biases in the data.
2. **Multi-task:** In a multi-task approach, we consider models that share some intermediate representation of the input. Considering the case of enzyme discovery, each task represents measurements for enzymes $\{x_1, x_2, \dots, x_n\}$ against a separate substrate. For a single enzyme, x_i , we denote $\{y_{i,1}, y_{i,2}, \dots, y_{i,m}\}$ as the activities of x_i against all m substrates. Thus, the goal is to learn some set of $j \in \{1, 2, \dots, m\}$ functions $f_j(x_i) = \hat{y}_{i,j}$. In multi-task learning, we force each function to have a shared intermediate representation, decomposing f_j as:

$$f_j(x_i) = H_j(G(x_i)) = \hat{y}_{i,j}$$

Therefore, each model must learn a common transformation $G(x_i)$ and is able to share information across tasks. In practice, we can accomplish this multi-task learning setup by using a multi-layer perceptron (MLP) model that has a single shared intermediate layer and outputs m different values in the final layer corresponding to the input enzyme’s activity against all m substrates.

For enzyme discovery, while this model class can share information across substrates, there is no meaningful representation of the substrate itself included in this model. The same properties hold for substrate discovery, with the tasks featuring measurements against different enzymes instead.

3. **CPI:** On the other hand, the CPI based models are able to consider *both* meaningful representations of the enzyme and the substrate. We use primarily feed forward neural networks with concatenation and dot product layers that fuse representations of the substrate and enzyme.
4. **Single-task:** We consider single task models that act on enzymes (substrates) measured against a single substrate (enzyme) target, independently from other substrate (enzyme) targets. Unlike CPI and multi-task models, single-task models have no ability to share information across tasks.

In addition to model classes, we test several different representations for both enzymes and substrates. For enzyme representations, we consider three different featurization schemes:

1. **Sequence:** When considering sequence similarity in K-nearest neighbor approaches, we featurize the enzyme as a sequence

2. **ESM-1b:** ESM-1b [206] is a pretrained deep learning model that can extract meaningful featurizations of full enzyme sequences at each respective position.
3. **One-hot:** In the multi-task learning setting, rather than force models to pass through a shared intermediate, we can also encode each enzyme as a unique "one-hot" vector. That is, the featurization for the i^{th} enzyme in the dataset will be a zero vector with a value of 1 at the i^{th} position only, such that the sum of the vector is itself 1. Because such an encoding does not give the model any structural information about the enzyme or substrate, we utilize such a featurization in a second multi-task learning scheme.

To featurize substrates, we consider:

1. **Morgan:** Circular Morgan fingerprints have been a staple of cheminformatics based regression.. We use 1024-bit Morgan fingerprints [237] to represent structural features of each molecule.
2. **Random:** For baseline methods, we also consider randomly sampled compound features and learn simple ridge regression methods.
3. **One-hot:** As with enzyme sequences, we consider one-hot based encodings of substrates for multi-task enzyme discovery.
4. **JT-VAE:** In analogy to how we extract pre-trained representations from enzymes, we also consider using the encoding from Jin et al.’s JT-VAE model [236].

Table 6.4: **Summary and classifications of different models utilized.**

| Model class | Model | Enz. features | Sub. features | Enz. discovery? | Sub. discovery? | Model name |
|-------------|---------------|---------------|---------------|-----------------|-----------------|------------------------|
| Baseline | KNN | Sequence | - | Yes | No | KNN: Levenshtein |
| | KNN | - | Morgan | No | Yes | KNN: Tanimoto |
| | Ridge | - | Random | No | Yes | Ridge: random feats. |
| Multi-task | FFN | ESM-1b | - | Yes | No | FFN: ESM-1b |
| | FFN (concat.) | ESM-1b | One-hot | Yes | No | FFN: [ESM-1b, one-hot] |
| | FFN | - | Morgan | No | Yes | FFN: Morgan |
| | FFN (concat) | One-hot | Morgan | No | Yes | FFN: [one-hot, Morgan] |
| CPI | FFN (concat) | ESM-1b | Morgan | Yes | Yes | FFN: [ESM-1b, Morgan] |
| | FFN (dot) | ESM-1b | Morgan | Yes | yes | FFN: ESM-1b • Morgan |
| Single-task | Ridge | ESM-1b | - | Yes | No | Ridge: ESM-1b |
| | Ridge | - | Morgan | No | Yes | Ridge: Morgan |
| | Ridge | - | JT-VAE | No | Yes | Ridge: JT-VAE |

Compound-protein models

When conducting our re-analysis of the CPI study from Hie et al., we use similar combinations of models and features. For consistency with their work, we utilize their model architectures, protein featurizations, and substrate featurizations. We showcase the differences among these in Table 6.5.

Table 6.6: **Full enzyme discovery area under the precision recall curve (AUPRC) results.** On the 6 different datasets tested (thiolase datasets used for hyperparameter optimization), K-nearest neighbor baselines with Levenshtein edit distance are compared against feed-forward networks using various featurizations and ridge regression models. Pre-trained features (“ESM-1b”) indicate protein features extracted from a masked language model trained on UniRef50 [206]. Two compound protein interaction architectures are tested, both concatenation and dot products, indicated with “[{prot repr.}, {sub repr.}]” and “[{prot repr.}•{sub repr.}]” respectively. Halogenase and glycosyltransferase datasets are evaluated using leave-one-out splits, whereas BKACE, phosphatase, and esterase datasets are evaluated with 5 repeats of 10 different cross validation splits. Average precision is calculated using scikit-learn for each substrate task separately before being averaged. Average values are presented across 3 random seeds \pm standard error.

¹Used for hyperparameter optimization

| Method Type | Dataset Method | BKACE | Esterase | Glyco. | Halogenase | Phosphatase | Thiolase ¹ |
|-------------|------------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|
| Baselines | KNN: Levenshtein | 0.564 \pm 0.011 | 0.528 \pm 0.002 | 0.539 \pm 0.003 | 0.375 \pm 0.014 | 0.316 \pm 0.004 | 0.499 \pm 0.001 |
| CPI | FFN: [ESM-1b, Morgan] | 0.478 \pm 0.012 | 0.579 \pm 0.007 | 0.581 \pm 0.008 | 0.489 \pm 0.025 | 0.386 \pm 0.006 | 0.536 \pm 0.027 |
| | FFN: ESM-1b • Morgan | 0.645 \pm 0.007 | 0.588 \pm 0.011 | 0.559 \pm 0.010 | 0.463 \pm 0.021 | 0.389 \pm 0.002 | 0.541 \pm 0.003 |
| Multi-task | FFN: [ESM-1b, one-hot] | 0.433 \pm 0.007 | 0.557 \pm 0.007 | 0.526 \pm 0.024 | 0.510 \pm 0.043 | 0.361 \pm 0.008 | 0.552 \pm 0.013 |
| | FFN: ESM-1b | 0.664 \pm 0.011 | 0.572 \pm 0.008 | 0.543 \pm 0.024 | 0.420 \pm 0.003 | 0.359 \pm 0.005 | 0.487 \pm 0.003 |
| Single-task | Ridge: ESM-1b | 0.648 \pm 0.011 | 0.583 \pm 0.003 | 0.575 \pm 0.000 | 0.446 \pm 0.000 | 0.413 \pm 0.005 | 0.519 \pm 0.000 |

Table 6.5: **Summary and classifications of different models utilized in our reanalysis of Hie et al. [226]**

| Model class | Model | Prot. features | Sub. features | Repurposing? | Discovery? | Original study? |
|-------------|----------|----------------|---------------|--------------|------------|-----------------|
| CPI | MLP | Bepler | JT-VAE | Yes | Yes | Yes |
| | GP + MLP | Bepler | JT-VAE | Yes | Yes | Yes |
| No CPI | MLP | Bepler | - | Yes | No | No |
| | GP + MLP | Bepler | - | Yes | No | No |
| | MLP | - | JT-VAE | No | Yes | No |
| | GP + MLP | - | JT-VAE | No | Yes | No |
| Linear | Ridge | Bepler | - | Yes | No | No |
| | Ridge | - | JT-VAE | No | Yes | No |
| | Ridge | - | Morgan | No | Yes | No |

6.6.3 Extended results

Table 6.7: **Full enzyme discovery area under the receiver operating curve (AUC-ROC) results.** On the 6 different datasets tested (thiolase datasets used for hyperparameter optimization), K-nearest neighbor baselines with Levenshtein edit distance are compared against feed-forward networks using various featurizations and ridge regression models. ESM-1b features indicate protein features extracted from a masked language model trained on UniRef50 [206]. Two compound protein interaction architectures are tested, both concatenation and dot products, indicated with “[{prot repr.}, {sub repr.}]” and “[{prot repr.}•{sub repr.}]” respectively. Halogenase and glycosyltransferase datasets are evaluated using leave-one-out splits, whereas BKACE, phosphatase, and esterase datasets are evaluated with 5 repeats of 10 different cross validation splits. AUC ROC is calculated using scikit-learn for each substrate task separately before being averaged. Average values are presented across 3 random seeds \pm standard error.

¹Used for hyperparameter optimization

| Method Type | Dataset Method | BKACE | Esterase | Glyco. | Halogenase | Phosphatase | Thiolase ¹ |
|-------------|------------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|
| Baselines | KNN: Levenshtein | 0.896 \pm 0.003 | 0.686 \pm 0.002 | 0.623 \pm 0.002 | 0.506 \pm 0.013 | 0.635 \pm 0.003 | 0.560 \pm 0.003 |
| CPI | FFN: [ESM-1b, Morgan] | 0.793 \pm 0.001 | 0.723 \pm 0.003 | 0.636 \pm 0.019 | 0.568 \pm 0.036 | 0.676 \pm 0.008 | 0.637 \pm 0.025 |
| | FFN: ESM-1b • Morgan | 0.884 \pm 0.002 | 0.730 \pm 0.004 | 0.647 \pm 0.021 | 0.587 \pm 0.025 | 0.678 \pm 0.001 | 0.637 \pm 0.011 |
| Multi-task | FFN: [ESM-1b, one-hot] | 0.768 \pm 0.006 | 0.714 \pm 0.006 | 0.586 \pm 0.018 | 0.616 \pm 0.043 | 0.657 \pm 0.002 | 0.651 \pm 0.032 |
| | FFN: ESM-1b | 0.892 \pm 0.002 | 0.713 \pm 0.002 | 0.624 \pm 0.025 | 0.564 \pm 0.003 | 0.669 \pm 0.004 | 0.569 \pm 0.024 |
| Single-task | Ridge: ESM-1b | 0.892 \pm 0.004 | 0.725 \pm 0.001 | 0.653 \pm 0.000 | 0.567 \pm 0.000 | 0.703 \pm 0.002 | 0.589 \pm 0.000 |

Table 6.8: **Full substrate discovery area under the precision recall curve (AUPRC) results.** CPI models and single task models are compared on the glycosyltransferase, esterase, and phosphatase datasets, all with 5 trials of 10-fold cross validation. Each model and featurization is compared to “Ridge: Morgan” using a 2-sided Welch T test, with each additional asterisk representing significance at [0.05, 0.01, 0.001, 0.0001] thresholds respectively after applying a Benjamini-Hochberg correction. Pretrained substrate featurizations used in “Ridge: JT-VAE” are features extracted from a junction-tree variational auto-encoder (JT-VAE) [236]. Two compound protein interaction architectures are tested, both concatenation and dot-product, indicated with “[{prot repr.}, {sub repr.}]” and “[{prot repr.}•{sub repr.}]” respectively. In the interaction based architectures, ESM-1b indicates the use of a masked language model trained on UniRef50 as a protein representation [206]. Average precision is calculated using scikit-learn for each substrate task separately before being averaged. Models are hyperparameter optimized on a held out halogenase dataset. Values represent mean values across 3 random seeds \pm standard error.

¹Used for hyperparameter optimization

| Method Type | Dataset Method | Esterase | Glyco. | Halogenase ¹ | Phosphatase |
|-------------|------------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|
| Baselines | KNN: Tanimoto | 0.609 \pm 0.008 | 0.588 \pm 0.006 | 0.464 \pm 0.000 | 0.426 \pm 0.002 |
| | Ridge: random feats. | 0.327 \pm 0.014 | 0.239 \pm 0.032 | 0.258 \pm 0.037 | 0.294 \pm 0.006 |
| CPI | FFN: [ESM-1b, Morgan] | 0.674 \pm 0.019 | 0.673 \pm 0.010 | 0.471 \pm 0.005 | 0.462 \pm 0.006 |
| | FFN: ESM-1b • Morgan | 0.709 \pm 0.018 | 0.693 \pm 0.005 | 0.478 \pm 0.033 | 0.506 \pm 0.004 |
| Multi-task | FFN: Morgan | 0.654 \pm 0.011 | 0.689 \pm 0.005 | 0.362 \pm 0.019 | 0.442 \pm 0.001 |
| | FFN: [one-hot, Morgan] | 0.707 \pm 0.019 | 0.677 \pm 0.009 | 0.469 \pm 0.009 | 0.493 \pm 0.005 |
| Single-task | Ridge: Morgan | 0.716 \pm 0.010 | 0.699 \pm 0.008 | 0.525 \pm 0.000 | 0.504 \pm 0.002 |
| | Ridge: JT-VAE | 0.489 \pm 0.004 | 0.505 \pm 0.007 | 0.468 \pm 0.000 | 0.411 \pm 0.002 |

Table 6.9: **Full substrate discovery area under the receiver operating curve (AUC-ROC) results.** CPI models and single task models are compared on the glycosyltransferase, esterase, and phosphatase datasets, all with 5 trials of 10-fold cross validation. Each model and featurization is compared to “Ridge: Morgan” using a 2-sided Welch T test, with each additional asterisk representing significance at [0.05, 0.01, 0.001, 0.0001] thresholds respectively after applying a Benjamini-Hochberg correction. Pretrained substrate featurizations used in “Ridge: JT-VAE” are features extracted from a junction-tree variational auto-encoder (JT-VAE) [236]. Two compound protein interaction architectures are tested, both concatenation and dot-product, indicated with “[{prot repr.}, {sub repr.}]” and “[{prot repr.}•{sub repr.}]” respectively. In the interaction based architectures, “ESM-1b” indicates the use of a masked language model trained on UniRef50 as a protein representation [206]. Models are hyperparameter optimized on a held out halogenase dataset. Values represent mean values across 3 random seeds \pm standard error.

¹Used for hyperparameter optimization

| Method Type | Dataset Method | Esterase | Glyco. | Halogenase ¹ | Phosphatase |
|-------------|------------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|
| Baselines | KNN: Tanimoto | 0.807 \pm 0.001 | 0.855 \pm 0.004 | 0.739 \pm 0.000 | 0.680 \pm 0.002 |
| | Ridge: random feats. | 0.513 \pm 0.019 | 0.483 \pm 0.042 | 0.440 \pm 0.051 | 0.481 \pm 0.011 |
| CPI | FFN: [ESM-1b, Morgan] | 0.808 \pm 0.005 | 0.883 \pm 0.004 | 0.726 \pm 0.006 | 0.689 \pm 0.004 |
| | FFN: ESM-1b • Morgan | 0.831 \pm 0.003 | 0.892 \pm 0.001 | 0.728 \pm 0.021 | 0.715 \pm 0.003 |
| Multi-task | FFN: Morgan | 0.784 \pm 0.004 | 0.880 \pm 0.007 | 0.583 \pm 0.022 | 0.675 \pm 0.002 |
| | FFN: [one-hot, Morgan] | 0.833 \pm 0.006 | 0.880 \pm 0.002 | 0.750 \pm 0.007 | 0.711 \pm 0.005 |
| Single-task | Ridge: Morgan | 0.841 \pm 0.003 | 0.878 \pm 0.005 | 0.745 \pm 0.000 | 0.710 \pm 0.001 |
| | Ridge: JT-VAE | 0.724 \pm 0.003 | 0.751 \pm 0.004 | 0.682 \pm 0.000 | 0.641 \pm 0.001 |

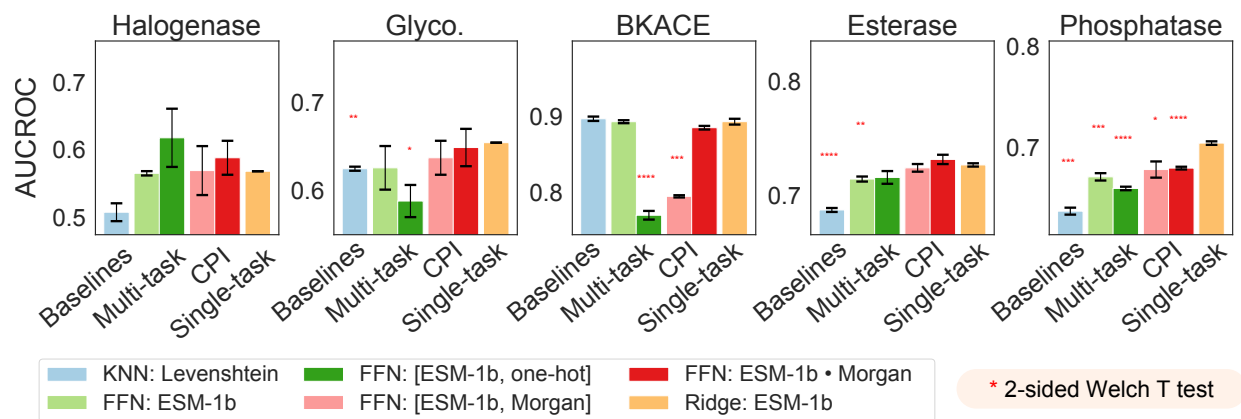


Figure 6.8: **Enzyme discovery benchmarking with AUCROC** On the 5 different datasets tested, K-nearest neighbor baselines with Levenshtein edit distance are compared against feed-forward networks using various featurizations and ridge regression models in terms of AUC ROC performance. ESM-1b features indicate protein features extracted from a masked language model trained on UniRef50 [206]. Concatenation and dot product architectures are indicated with “[{prot repr.}, {sub repr.}]” and “[{prot repr.}•{sub repr.}]” respectively. Halogenase and glycosyltransferase datasets are evaluated using leave-one-out splits. BKACE, phosphatase, and esterase datasets are evaluated with 5 repeats of 10 different cross validation splits. AUC ROC is calculated using scikit-learn for each substrate task separately before being averaged. Error bars represent the standard error of the mean across 3 random seeds. Each model and featurization is compared to “Ridge: ESM-1b” using a 2-sided Welch T test, with each additional asterisk representing significance at [0.05, 0.01, 0.001, 0.0001] thresholds respectively after applying a Benjamini-Hochberg correction.

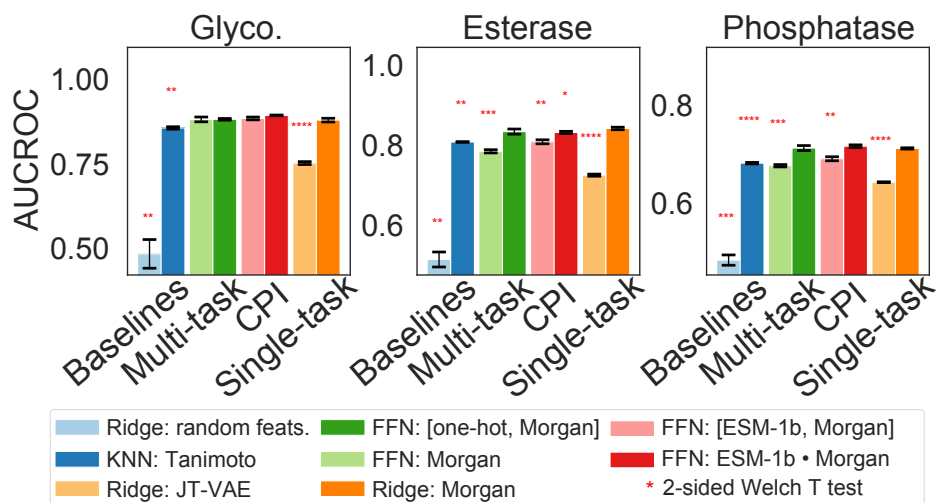


Figure 6.9: **Full substrate discovery AUC ROC results.** CPI models and single task models are compared on the glycosyltransferase, esterase, and phosphatase datasets, all with 5 trials of 10-fold cross validation. Error bars represent the standard error of the mean across 3 random seeds. Each model and featurization is compared to “Ridge: Morgan” using a 2-sided Welch T test, with each additional asterisk representing significance at [0.05, 0.01, 0.001, 0.0001] thresholds respectively after applying a Benjamini-Hochberg correction. Pretrained substrate featurizations used in “Ridge: JT-VAE” are features extracted from a junction-tree variational auto-encoder (JT-VAE) [236]. Concatenation and dot-product architectures are indicated with “[{prot repr.}, {sub repr.}]” and “[{prot repr.}•{sub repr.}]” respectively. In the interaction based architectures, “ESM-1b” indicates the use of a masked language model trained on UniRef50 as a protein representation [206]. Models are hyper-parameter optimized on a held out halogenase dataset.

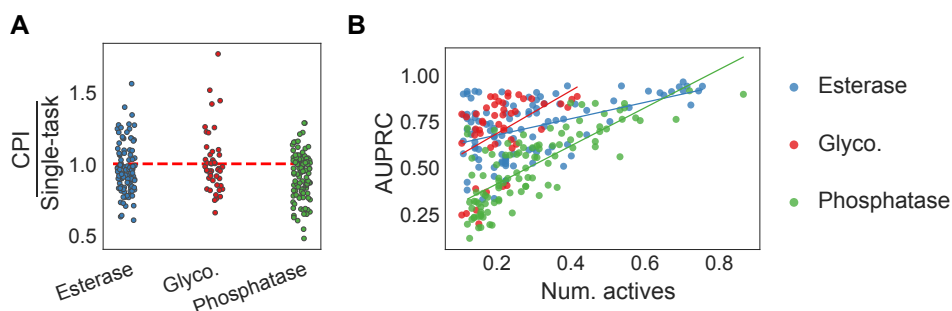


Figure 6.10: **Substrate Discovery Extended Analysis (A)** Average AUPRC on each individual “enzyme task” is compared between compound protein interaction models and single-task models. Points below 1 indicate substrates on which single-task models better predict enzyme activity than CPI models. CPI models used are “FFN: [ESM-1b, Morgan]” and single-task models are “Ridge: Morgan”. **(B)** AUPRC values from the ridge regression model broken out by each task are plotted against the fraction of active enzymes in the dataset. Best fit lines are drawn through each dataset to serve as a visual guide.

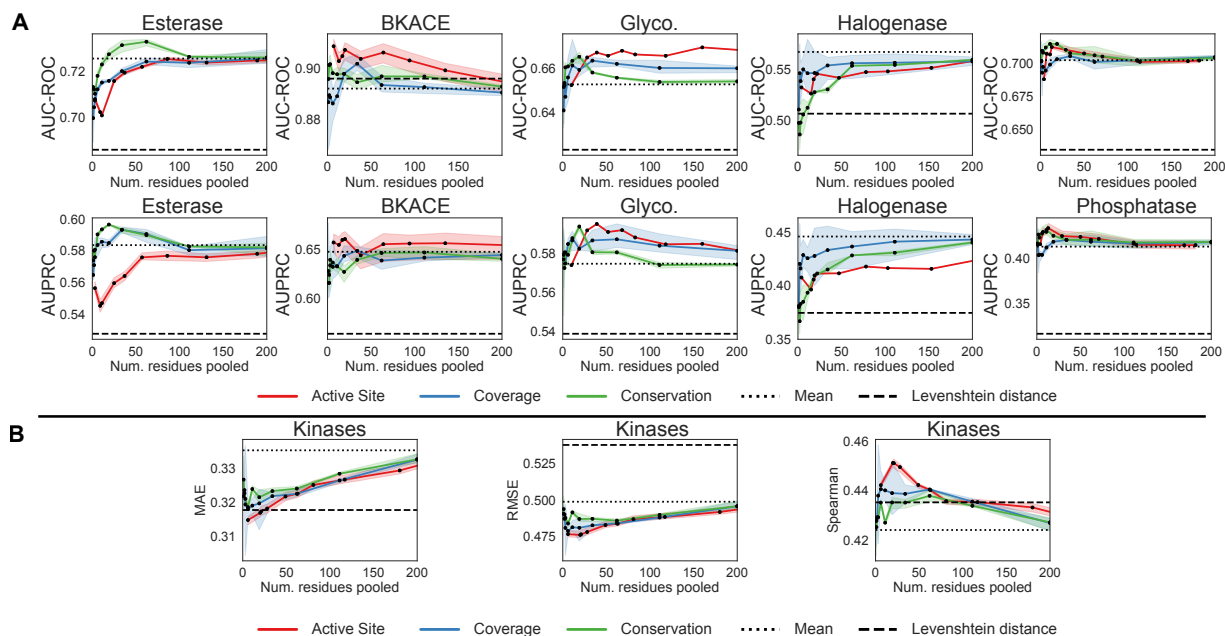


Figure 6.11: **MSA and structure based pooling across all datasets tested (A)** Active site, coverage, conservation, and mean pooling are plotted for all 5 enzyme discovery datasets tested. Both AUCROC and AUPRC values are shown. These are compared against the Levenshtein distance baseline (dotted). **(B)** Equivalent analysis is conducted on the filtered kinase dataset extracted from Davis et al. with MAE, RMSE, and Spearman rank correlation shown [223]. The same hyperparameters are used as set in Figure 6.2 for ridge regression models. All experiments are repeated for 3 random seeds following the same split evaluation as in other enzyme discovery model benchmarking.

Prediction outputs

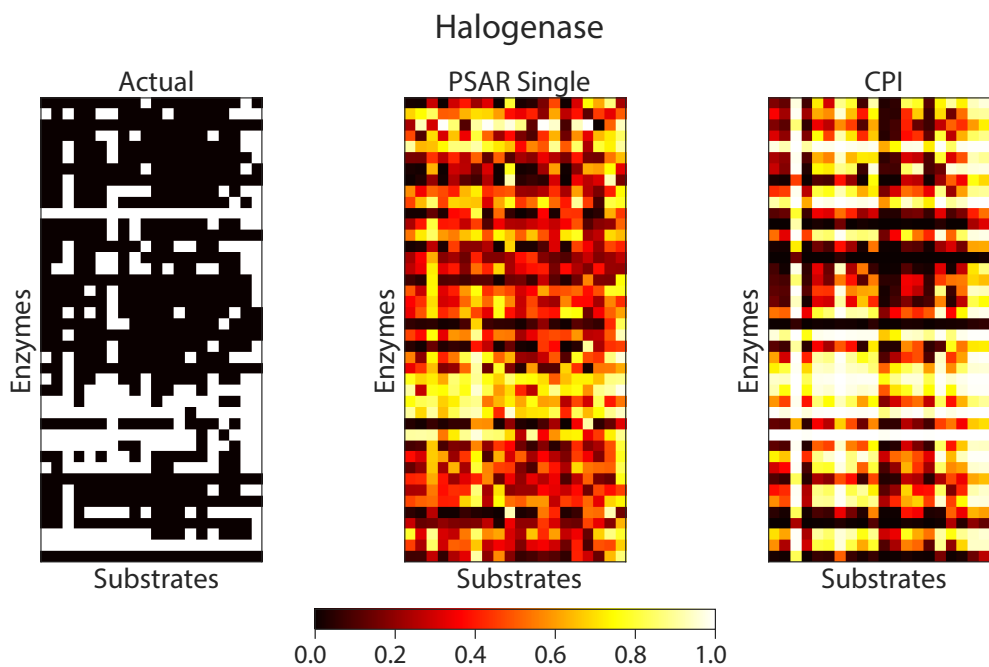


Figure 6.12: **Enzyme discovery halogenase prediction results** Ground truth binary enzyme-substrate activities (left) are compared against a single seed of predictions made through cross validation using a single-task ridge regression model (middle) and a CPI based model, FFN: [ESM-1b, Morgan] (right).

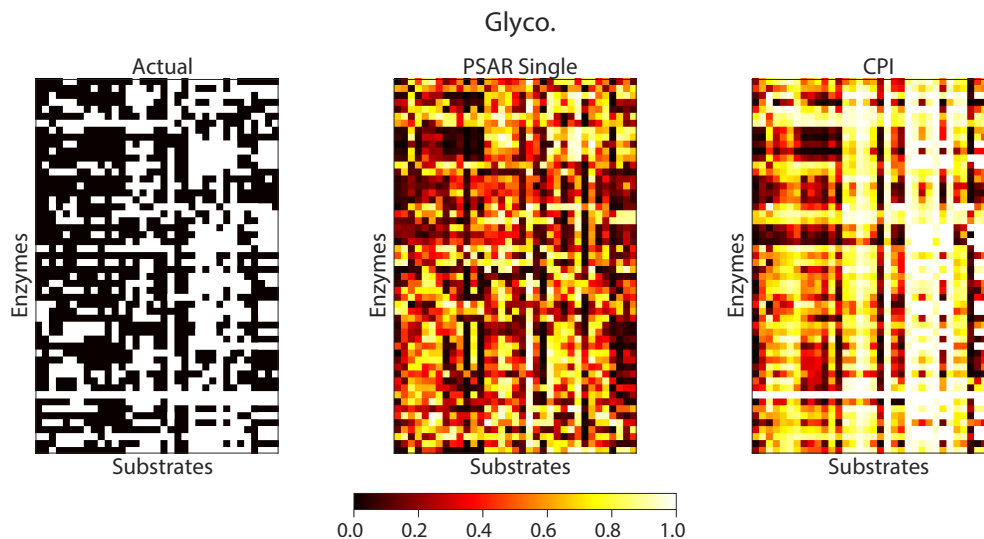


Figure 6.13: **Enzyme discovery glycosyltransferase prediction results** Ground truth binary enzyme-substrate activities (left) are compared against a single seed of predictions made through cross validation using a single-task ridge regression model (middle) and a CPI based model, FFN: [ESM-1b, Morgan] (right).

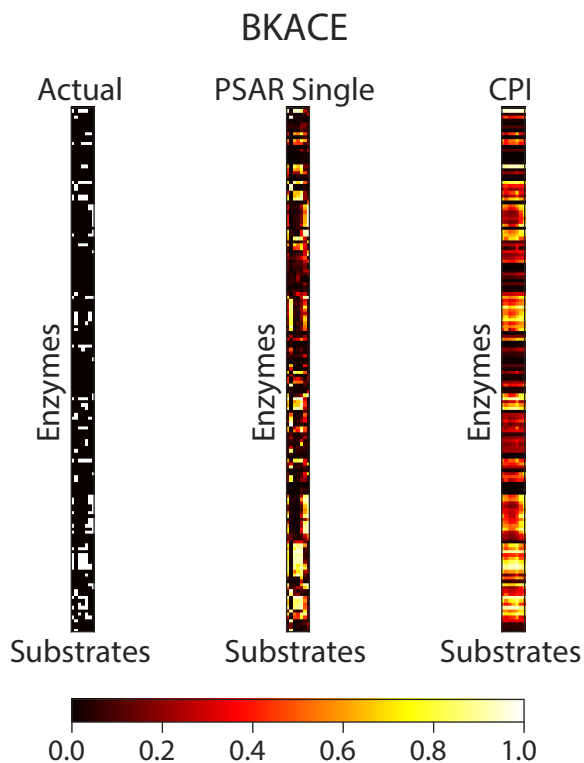


Figure 6.14: **Enzyme discovery BKACE prediction results** Ground truth binary enzyme-substrate activities (left) are compared against a single seed of predictions made through cross validation using a single-task ridge regression model (middle) and a CPI based model, FFN: [ESM-1b, Morgan] (right).

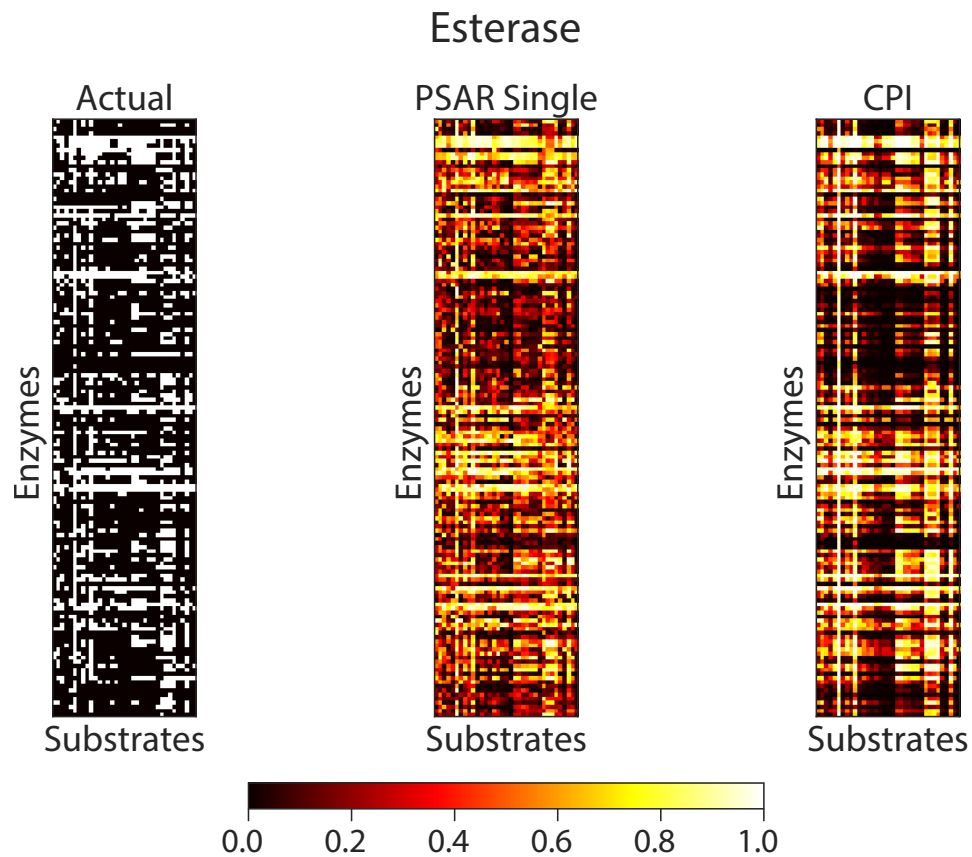


Figure 6.15: **Enzyme discovery esterase prediction results** Ground truth binary enzyme-substrate activities (left) are compared against a single seed of predictions made through cross validation using a single-task ridge regression model (middle) and a CPI based model, FFN: [ESM-1b, Morgan] (right).

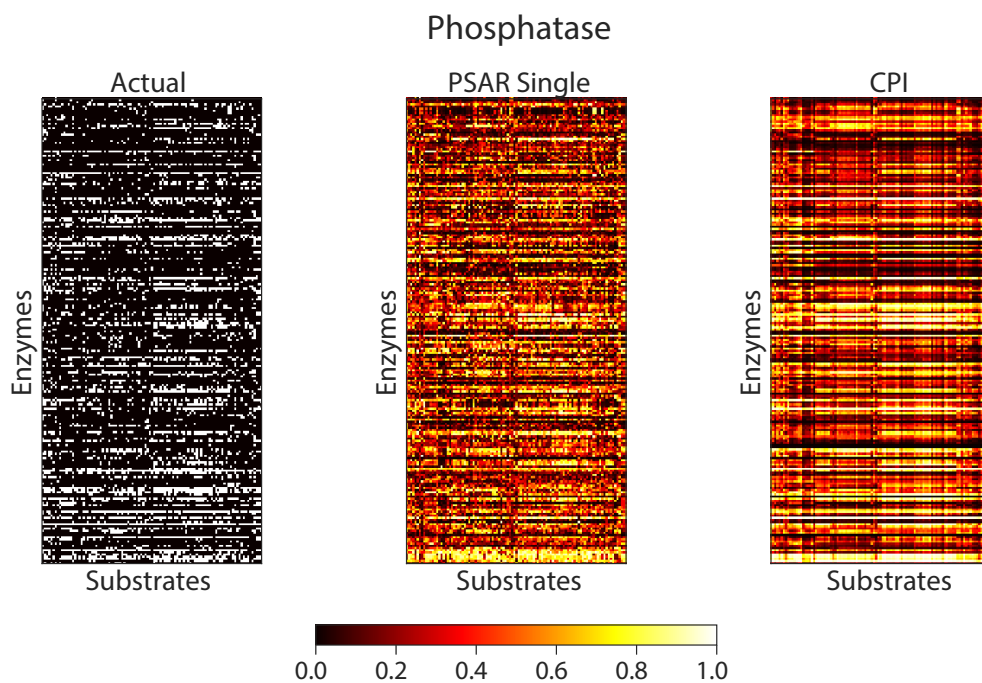


Figure 6.16: **Enzyme discovery phosphatase prediction results** Ground truth binary enzyme-substrate activities (left) are compared against a single seed of predictions made through cross validation using a single-task ridge regression model (middle) and a CPI based model, FFN: [ESM-1b, Morgan] (right).

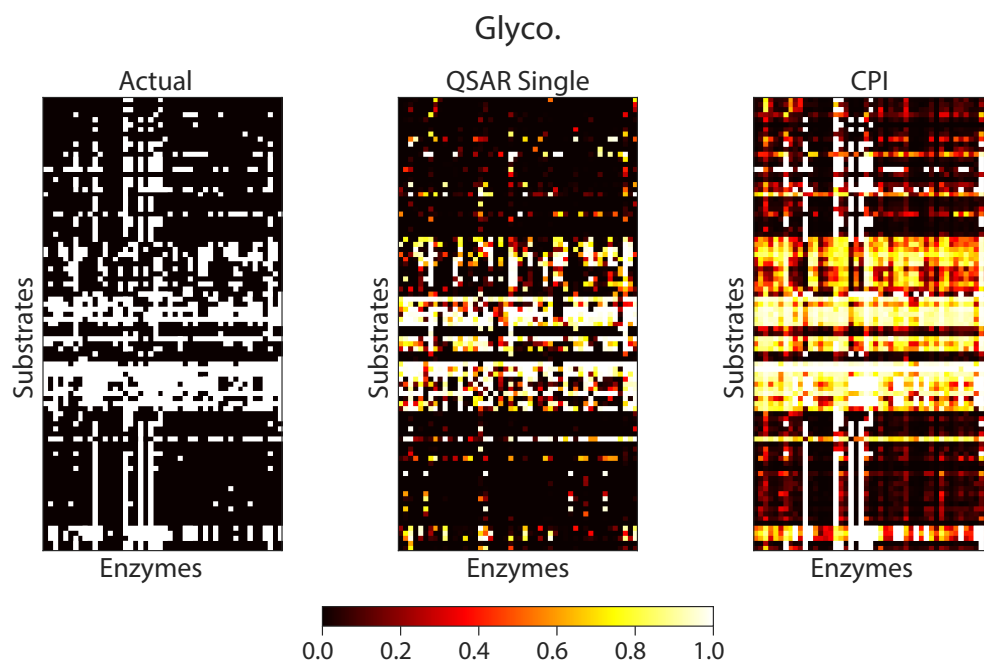


Figure 6.17: **Substrate discovery glycosyltransferase prediction results** Ground truth binary enzyme-substrate activities (left) are compared against a single seed of predictions made through cross validation using a single-task ridge regression model (middle) and a CPI based model, FFN: [ESM-1b, Morgan] (right).

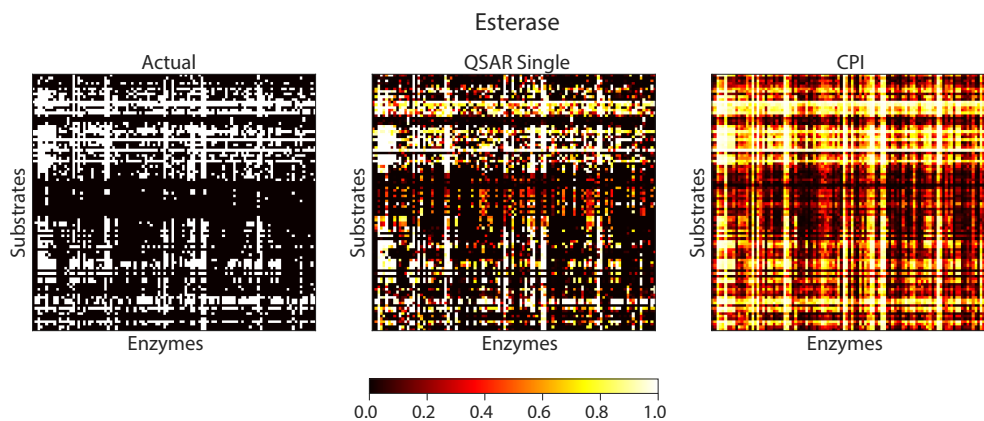


Figure 6.18: **Substrate discovery esterase prediction results** Ground truth binary enzyme-substrate activities (left) are compared against a single seed of predictions made through cross validation using a single-task ridge regression model (middle) and a CPI based model, FFN: [ESM-1b, Morgan] (right).

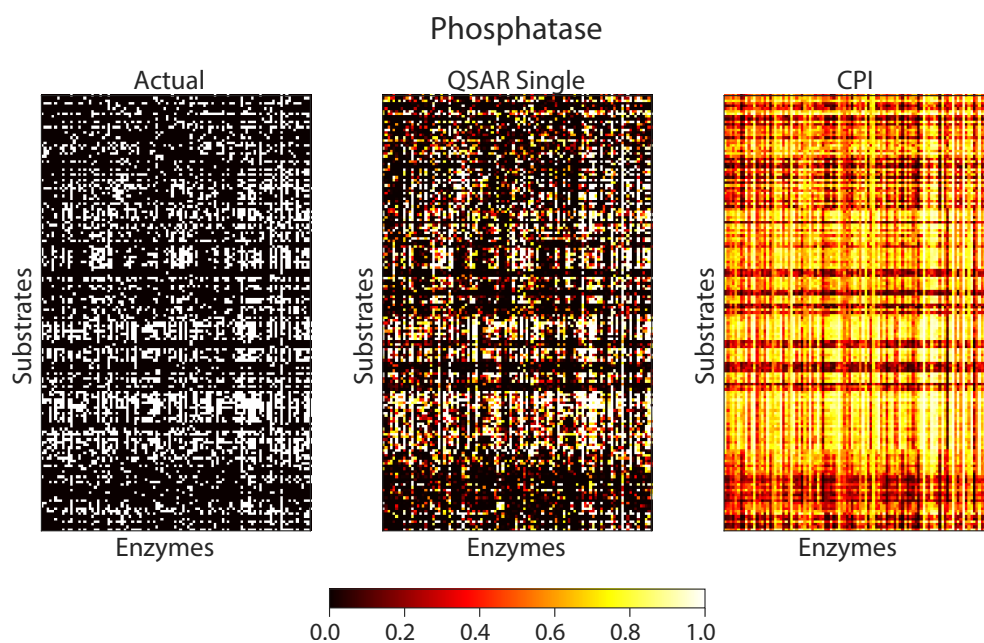


Figure 6.19: **Substrate discovery phosphatase prediction results** Ground truth binary enzyme-substrate activities (left) are compared against a single seed of predictions made through cross validation using a single-task ridge regression model (middle) and a CPI based model, FFN: [ESM-1b, Morgan] (right).

Chapter 7

Evidential Molecular Prediction and Discovery

This chapter continues the theme of predicting functional attributes of small molecules, metabolites, and their interactions. Specifically, we consider how new computational methods for quantifying uncertainty can be used for molecular property prediction. This text has previously appeared as A. P. Soleimany, A. Amini, S. Goldman, *et al.*, “Evidential Deep Learning for Guided Molecular Property Prediction and Discovery,” *ACS Central Science*, vol. 7, no. 8, pp. 1356–1367, 2021. I contributed equally alongside the first two authors. All co-first authors contributed to project conceptualization, the implementation of the methods, and presentation of the work, with additional input from supervising advisors.

7.1 Introduction

As quantitative structure-activity relationship (QSAR) models are increasingly applied across the chemical and physical sciences to guide time- and resource-intensive experimentation, an understanding of when to trust model predictions is of critical importance [251]–[253]. Though neural networks have shown tremendous promise in QSAR modeling [254], [255], they remain difficult to interpret, are susceptible to pathological failures in out-of-domain regimes, and lack guarantees on their robustness. Therefore, a better understanding of predictive confidence of neural models is essential, particularly for drug discovery and virtual screening applications where model predictions can inform safety-critical experimental pipelines.

Uncertainty quantification (UQ) methods can help meet this critical need to facilitate the robust application of neural models in the chemical sciences. Indeed, significant work has been done in establishing general methods to estimate epistemic uncertainties (i.e., model uncertainty due to uncertainty in parameters and predictions) and aleatoric uncertainties (i.e., data uncertainty due to noise inherent in the observations) of neural network predictions [256]. Recent studies have demonstrated the importance of focusing explicitly on epistemic uncertainty in the contexts of property and reaction prediction in the chemical sciences [257], discovery in the biological sciences [226], as well as in healthcare more broadly [258]. While a plethora of distance-based and non-parametric methods for UQ have

been developed [259], [260], Bayesian neural networks [261] and sampling based approaches, such as model ensembling [262] and dropout sampling [263], are still accepted as state of the art for epistemic UQ in neural networks, due in part to their model-agnostic nature and ease of implementation [257], [264], [265].

However, these approaches only generate approximations to the underlying uncertainty functions via stochastic sampling, incurring computational costs and runtimes that are routinely an order of magnitude higher than those of single models. This poses a significant challenge to using these epistemic uncertainty models in iterative active learning procedures, scans of very large chemical libraries, and molecular dynamics simulations [265], [266]. Additionally, the most recent adaptations of atomistic neural networks for prediction of potential energy surfaces and quantum mechanical properties have achieved state-of-the-art results by using more expressive, larger network architectures that sacrifice speed for predictive accuracy [267], [268]. Large model sizes compound the computational expense of deploying sampling-based UQ methods and necessitate the development of more efficient approaches. Though neural networks can be trained to obtain closed-form solutions of only aleatoric uncertainty without sampling, these methods fail to provide estimates of epistemic uncertainties, limiting their broad utility [269]–[273]. More generally, recent analyses have revealed an overwhelming lack of consensus as to the top performing UQ methods across molecular property prediction datasets [264], [274]. Thus, there remains a need for fast, calibrated, and scalable UQ methods for QSAR models that provide estimates of model uncertainty and can be deployed across a range of molecular property prediction and discovery tasks.

Emerging evidential deep learning algorithms have the potential to address these limitations in their ability to directly learn grounded representations of epistemic uncertainty without the need for sampling [275], [276]. Specifically, these methods formulate learning as an evidence acquisition process, wherein new training examples add support to a learned evidential distribution that parameterizes a probability distribution over the network’s likelihood function. Evidential learning therefore offers the promise of efficient and calibrated uncertainty learning without the need for sampling. Furthermore, evidential neural networks can be implemented without significant architecture changes, but rather via modifications to the training loss function, and could thus enable tight integration with domain-specific architectures. However, while evidential learning formulations for both regression [276] and classification [275] have recently been presented, the utility of these methods on complex, non-uniform inputs, such as molecular graphs pervasive throughout the chemical sciences, has yet to be shown.

In this work, we establish evidential deep learning as a new approach to UQ for molecular structure-property prediction (Figure 7.1). Specifically, this work makes the following contributions:

1. Development of evidential message passing and atomistic networks that learn 2D or 3D molecular representations, respectively, and return well-calibrated epistemic uncertainties without any sampling;
2. Evaluation of evidential uncertainties on benchmark QSAR regression tasks against gold-standard, sampling-based UQ methods;
3. Validation of the relevance of evidential deep learning to key molecular discovery ap-

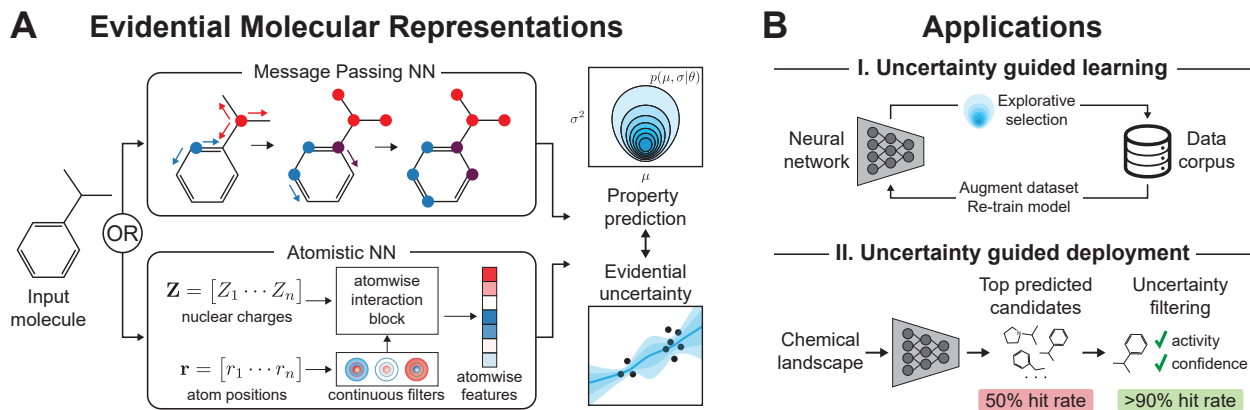


Figure 7.1: **Evidential uncertainty for molecular prediction and discovery.** (A) Evidential direct message passing or atomistic neural networks learn molecular representations, predict target properties, and infer the parameters of an underlying evidential distribution that captures the evidence in support of each prediction and enables uncertainty estimation. (B) Uncertainties are applied during learning (I) to guide sample acquisition and during deployment (II) to discover high confidence candidates with high empirical success rates.

applications that require sample prioritization from a larger screening library, namely:

- (i) Uncertainty-guided learning for sample-efficient model training and accelerated property optimization;
- (ii) Uncertainty-guided deployment for prioritization of high confidence candidates in virtual screening.

Taken together, our experiments validate a framework to use evidential deep learning as a powerful and flexible replacement for UQ in molecular property prediction and discovery tasks across the chemical sciences.

7.2 Approach

7.2.1 Formulating evidential learning for molecules

Evidential deep learning models [275], [276] are a recent approach to training single networks to estimate predictive uncertainties. While neural networks have been trained to output probabilities, for example with Softmax [277] for classification or Gaussian distributions (MVE) [269] for regression, these approaches estimate the probability of an output but neglect the model’s uncertainty associated with that output. Evidential deep learning extends the idea of learning the parameters of a probability distribution further to predict higher-order distributions over the original likelihood parameters themselves. These higher-order parameters define the evidential distribution and capture both the model’s prediction as well as the degree of evidence associated with that prediction. These models estimate

uncertainty by directly learning the parameters defining this evidential distribution, and are closely related to Bayesian neural networks [263] and ensembling approaches [262] which estimate the model’s uncertainty by sampling from the likelihood distribution, instead of directly learning to output it.

In the regression setting (e.g., prediction of a continuous target), we are given a dataset of paired training examples $\mathcal{D} = \{x_i, y_i\}_{i=1}^N$, for which the targets, $y_i \in \mathbb{R}$, can be assumed to be drawn i.i.d. from a Gaussian distribution with mean and variance $\theta = \{\mu, \sigma^2\}$. In the case of MVE networks, the likelihood parameters θ are deterministic and fixed, such that the model is optimized during training to predict these values directly, preventing estimation of model uncertainty. As an extension to this approach, evidential models assume these parameters are unknown and must instead be probabilistically estimated. This is done by placing priors over the likelihood parameters, such that the mean μ is drawn from a Gaussian distribution and the variance σ^2 is drawn from an Inverse-Gamma distribution. The resulting higher-order distribution (also referred to as the evidential distribution) thus can be represented by a Normal-Inverse-Gamma distribution, $p(\theta|\mathbf{m})$. This evidential distribution is specified by four parameters $\mathbf{m} = \{\gamma, \lambda, \alpha, \beta\}$. For continuous targets, evidential models directly learn these parameters \mathbf{m} which in turn define full distributions on top of the likelihood parameters $\{\mu, \sigma^2\}$, thus capturing the uncertainty in the model’s prediction (Figure 7.2A, B). Accordingly, the model outputs four values per target, corresponding to the four parameters of \mathbf{m} , and is trained using a multi-objective loss that aims to jointly maximize model fit while minimizing evidence on errors (Figure 7.2C).

In this work, we demonstrate that this evidential learning framework can be coupled with molecular feature extraction networks to predict target properties and also estimate uncertainty (Figure 7.2). We accomplish this by taking the molecular representation learned by a feature extractor (e.g., a neural network operating on 2D molecular graphs) and feeding it into a dense evidential layer which maps these higher dimensional learned features to the four evidential parameters \mathbf{m} which yield both property prediction and uncertainty estimates. All models are trained end-to-end, from molecule input to evidential property output, via backpropagation by optimizing the evidential loss function. Full model details and code for implementation can be found in Section 7.6.

7.3 Results

7.3.1 Uncertainty benchmarking

We first sought to demonstrate that our evidential learning algorithm could produce desirable uncertainties across both molecular and atomistic property prediction tasks. Given our emphasis on downstream tasks that require choosing the correct molecule from a larger screening library (Figure 7.1B), we evaluated whether predicted uncertainties are correctly ranked with respect to error; that is, predictions with the lowest uncertainty should also be expected to have the lowest error.

We integrated the evidential regression method into a directed message passing neural network (D-MPNN) architecture and assessed its performance in the “lower-N” (dataset of $\leq 10,000$ molecules) regression setting on commonly used benchmarking datasets [278] of

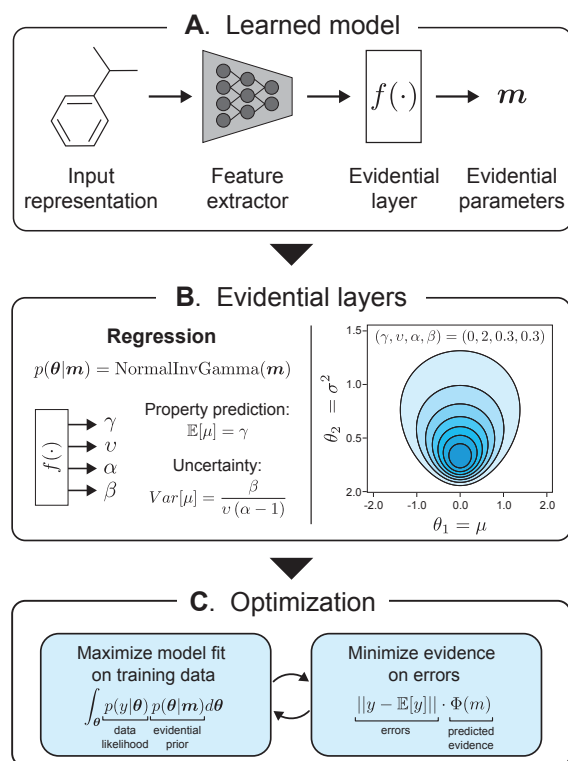


Figure 7.2: **Building and training evidential models.** (A) Evidential layers can be added to the end of existing molecular feature extractor neural networks. The output of the evidential layer is the parameters (\mathbf{m}) defining the molecule’s evidential distribution. (B) For continuous regression learning problems, the evidential distribution $p(\boldsymbol{\theta}|\mathbf{m})$ can take the form of a Normal-Inverse-Gamma distribution, where $\mathbf{m} = \{\gamma, v, \alpha, \beta\}$, to yield both property prediction and uncertainty estimates. Color represents likelihood density (darker = greater density). (C) The model (feature extractor and evidential layer) is trained end-to-end using backpropagation with a dual-objective loss that jointly maximizes model fit and inflates uncertainty (i.e., minimizes evidence) on errors.

aqueous solubility (Delaney), solvation energy (Freesolv), lipophilicity (Lipo), and atomization energy (QM7) (Figure 7.3A). With smaller datasets, sampling approaches such as model ensembling are not prohibitively expensive in practice, so evidential regression must demonstrate more calibrated uncertainty predictions relative to standard sampling-based UQ methods to justify its adoption.

Evidential regression performed well in its ability to rank uncertainties with respect to error (Table 7.1, Figure 7.7). Specifically, the evidential method achieves the lowest error across all methods tested when considering only the top 5% most certain predictions for three of the four datasets tested (Delaney, Freesolv, QM7; Table 7.1). On both the Delaney and QM7 datasets, error returned by the evidential model is well below the second best performing method by the 50% confidence cutoff (Figure 7.3B, C). The drastic improvement over ensembles in QM7 is consistent with previous observations that single neural network

models are more accurate than ensembles in the top confidence percentiles on QM7 [264]. Still, in the lower-N setting, there is some variance in performance across datasets. On the lipophilicity dataset, the RMSE computed at uncertainty cutoff percentiles of 0.25 and below for evidential regression is higher (worse) than the dropout-based sampling method, showing no advantage in selecting the most accurately predicted test set molecules over dropout (Table 7.1, Figure 7.7). Furthermore, the evidential method yields higher rank correlation between uncertainty and error than both ensembles and dropout on two of the four lower-N datasets tested and is at least within one standard deviation of the ensemble method for three of the four datasets tested, supporting its ability to better rank predictions (Figure 7.9).

To further evaluate performance in the lower-N setting, we conducted a similar analysis on three additional lower-N datasets acquired from the Therapeutics Data Commons [238] with tasks to predict hepatocyte clearance (“Clearance”), median lethal dose (“LD50”), and plasma protein binding rates (“PPBR”). We find that evidence is competitive with sampling-based methods here as well, with RMSE that is at least as low as the top performing sampling method on all three datasets tested (Table 7.2). Furthermore, error for the evidential method decreases steeply as a function of predicted certainty (Figure 7.8), and rank correlation between error and uncertainty is highest for the evidential method on two out of the three datasets tested (Figure 7.9).

Given these promising results in the lower-N setting, we next evaluated the generalizability of the evidential method to larger scale datasets of $\geq 50,000$ data points (Figure 7.3D). Datasets of this size often represent large-scale chemical libraries [279] or are generated via expensive physics-based simulations [280]. In these settings, the ability to quickly rank a larger library based upon confidence is advantageous for guiding downstream analyses and experimentation [281]. To this end, we compared UQ methods for 2D message passing networks evaluated on each of two “higher-N” ($\geq 50,000$ molecules) datasets: the QM9 dataset containing computer-generated quantum mechanical properties for small organic molecules [278], [280], and a ligand docking dataset containing scores of 50,240 molecules docked against thymidylate kinase using AutoDock Vina [282], [283].

For both these datasets, evidential regression predictions have lower error than both ensemble and dropout-based methods for all confidence percentile cutoffs greater than 50%, demonstrating steeper declines in error as a function of confidence (Figure 7.3E, Table 7.1). Compared to the ensemble-based method, evidential regression also displays higher rank correlation between uncertainty and error on both the docking dataset ($\rho_{\text{evidence}} = 0.163 \pm 0.009$, $\rho_{\text{ensemble}} = 0.040 \pm 0.018$) and the QM9 dataset ($\rho_{\text{evidence}} = 0.469 \pm 0.084$, $\rho_{\text{ensemble}} = 0.244 \pm 0.097$) (Figure 7.9). On the QM9 dataset with 2D molecular representations, the evidential method exhibits steeper declines in error as a function of confidence cutoffs relative to the ensemble baseline, both for individual tasks (Figure 7.10) and on the aggregated uncertainty averaged across tasks (Figure 7.11A, B).

To demonstrate the utility of the evidential learning approach for a variety of chemistry-specific neural model architectures, we integrated the evidential regression loss function into an atomistic neural network, implemented via the SchNet pack software*, that operates on 3D molecular conformers (Figure 7.3D). While ensembles are more accurate with no cutoff calculations ($\text{MAE}_{\text{ensemble}} = 2.04 \times 10^{-2}$, $\text{MAE}_{\text{evidence}} = 2.98 \times 10^{-2}$; Table 7.1), the

*<https://github.com/atomistic-machine-learning/schnetpack>

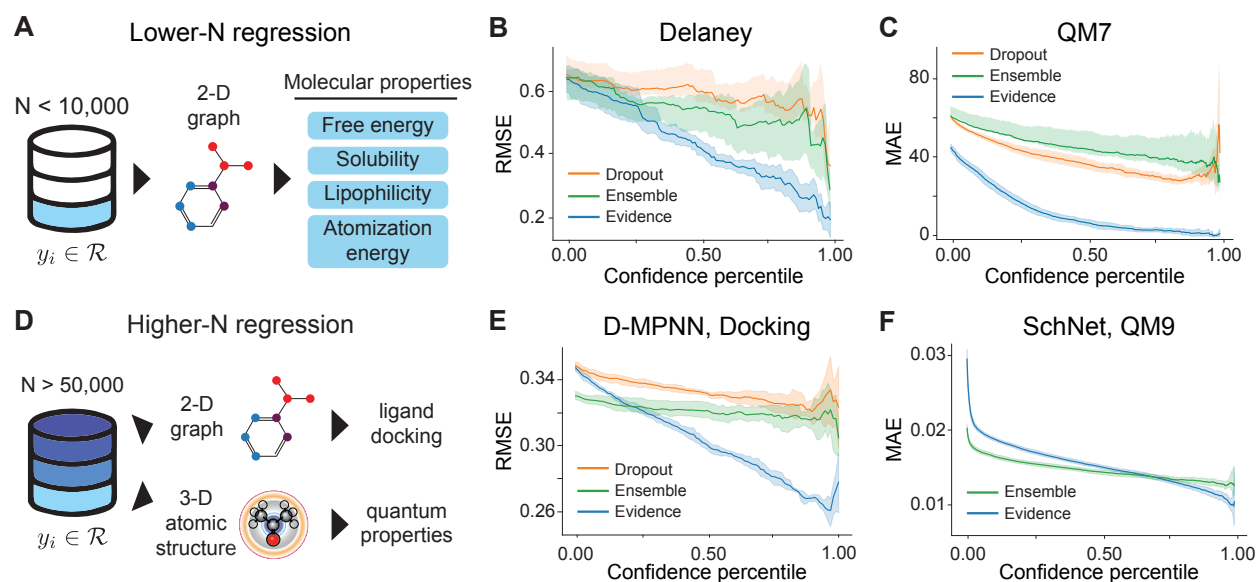


Figure 7.3: **Benchmarking evidential uncertainty for molecular property prediction.** (A) Lower-N regression tasks using 2D molecular representations for uncertainty benchmarking. (B, C) Prediction error, measured as root mean squared error (RMSE) or mean average error (MAE), at different confidence percentile cutoffs for the Delaney (B) and QM7 (C) datasets. Mean \pm 95% confidence interval (c.i.), $n = 10$ independent trials. (D) Higher-N regression tasks using 2D or 3D molecular representations. (E, F) Prediction error at different confidence percentile cutoffs for the Docking (E) and QM9 (F) datasets for 2D direct message passing (D-MPNN; E) and 3D atomistic (SchNet; F) neural networks, respectively. Mean \pm 95% c.i., $n = 5$ independent trials.

evidential method still produces uncertainties that correlate well with error (Figure 7.3F). When considering only predictions in the 95% confidence percentile, the evidential method displays a quantitative improvement over the ensemble method ($\text{MAE}_{\text{ensemble}} = 1.29 \times 10^{-2}$, $\text{MAE}_{\text{evidence}} = 1.12 \times 10^{-2}$; Figure 7.3F, Table 7.1). Rank correlation between error and uncertainty also reflects the steeper decrease of the evidential method relative to the ensemble based method ($\rho_{\text{evidence}} = 0.361 \pm 0.007$ vs. $\rho_{\text{ensemble}} = 0.220 \pm 0.012$; Figure 7.9C). Taken together, these results demonstrate the promise of evidential regression in achieving well-ranked uncertainty estimates across different dataset sizes and molecular representations, highlighting the modularity of this method.

7.3.2 Calibration and tunability

After observing steep reductions in error with increasing confidence, we next investigated the calibration of the predicted uncertainties – a critical property for the translation of any UQ method. With a perfectly calibrated classifier, we expect to find the true target value in the 90% credible interval 90% of the time [284]. However, if a regression model is overconfident, we would find the true value in the 90% credible interval less than 90% of the time and vice

| Cutoff | Delaney | | | Freesolv | | | Lipo | | | QM7 ($\times 10^2$) | | |
|--------|-----------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|-----------------------|-----------------------------------|-----------------------------------|
| | Dropout | Ensemble | Evidence | Dropout | Ensemble | Evidence | Dropout | Ensemble | Evidence | Dropout | Ensemble | Evidence |
| 0.0 | 0.68 \pm 0.02 | 0.65 \pm 0.03 | 0.66 \pm 0.02 | 1.00 \pm 0.06 | 0.94 \pm 0.06 | 0.96 \pm 0.07 | 0.55 \pm 0.01 | 0.53 \pm 0.02 | 0.55 \pm 0.02 | 1.18 \pm 0.02 | 1.12 \pm 0.02 | 1.15 \pm 0.03 |
| 0.5 | 0.62 \pm 0.03 | 0.55 \pm 0.03 | 0.44 \pm 0.01 | 0.79 \pm 0.07 | 0.45 \pm 0.04 | 0.42 \pm 0.04 | 0.52 \pm 0.01 | 0.40 \pm 0.01 | 0.50 \pm 0.01 | 0.88 \pm 0.06 | 0.88 \pm 0.06 | 0.39 \pm 0.03 |
| 0.75 | 0.59 \pm 0.03 | 0.50 \pm 0.05 | 0.35 \pm 0.02 | 0.85 \pm 0.12 | 0.41 \pm 0.05 | 0.36 \pm 0.04 | 0.50 \pm 0.02 | 0.38 \pm 0.02 | 0.51 \pm 0.02 | 0.65 \pm 0.03 | 0.81 \pm 0.06 | 0.23 \pm 0.04 |
| 0.90 | 0.55 \pm 0.03 | 0.51 \pm 0.09 | 0.28 \pm 0.02 | 0.66 \pm 0.20 | 0.40 \pm 0.06 | 0.35 \pm 0.08 | 0.46 \pm 0.03 | 0.38 \pm 0.02 | 0.53 \pm 0.03 | 0.69 \pm 0.05 | 0.71 \pm 0.11 | 0.10 \pm 0.04 |
| 0.95 | 0.53 \pm 0.06 | 0.45 \pm 0.06 | 0.22 \pm 0.02 | 0.75 \pm 0.30 | 0.27 \pm 0.04 | 0.38 \pm 0.12 | 0.49 \pm 0.04 | 0.36 \pm 0.03 | 0.50 \pm 0.04 | 0.73 \pm 0.08 | 0.69 \pm 0.11 | 0.10 \pm 0.04 |

| Cutoff | Enamine D-MPNN | | | QM9 D-MPNN | | | QM9 Atomistic ($\times 10^{-2}$) | |
|--------|-----------------|-----------------------------------|-----------------------------------|-----------------|-----------------------------------|-----------------------------------|------------------------------------|-----------------------------------|
| | Dropout | Ensemble | Evidence | Dropout | Ensemble | Evidence | Ensemble | Evidence |
| 0.0 | 3.40 \pm 0.12 | 4.47 \pm 0.18 | 5.60 \pm 0.20 | 0.35 \pm 0.00 | 0.33 \pm 0.00 | 0.35 \pm 0.00 | 2.04 \pm 0.03 | 2.98 \pm 0.08 |
| 0.5 | 3.64 \pm 0.05 | 2.12 \pm 0.02 | 1.55 \pm 0.12 | 0.33 \pm 0.00 | 0.32 \pm 0.00 | 0.30 \pm 0.00 | 1.45 \pm 0.02 | 1.52 \pm 0.02 |
| 0.75 | 3.42 \pm 0.04 | 1.94 \pm 0.04 | 1.04 \pm 0.13 | 0.33 \pm 0.00 | 0.32 \pm 0.00 | 0.28 \pm 0.00 | 1.36 \pm 0.02 | 1.33 \pm 0.02 |
| 0.90 | 3.30 \pm 0.06 | 1.80 \pm 0.03 | 0.63 \pm 0.12 | 0.32 \pm 0.00 | 0.32 \pm 0.01 | 0.27 \pm 0.00 | 1.31 \pm 0.03 | 1.18 \pm 0.03 |
| 0.95 | 3.26 \pm 0.05 | 1.79 \pm 0.05 | 0.42 \pm 0.01 | 0.33 \pm 0.01 | 0.32 \pm 0.01 | 0.26 \pm 0.01 | 1.29 \pm 0.03 | 1.12 \pm 0.03 |

Table 7.1: **Model error at various confidence percentile cutoffs.** For a given confidence percentile cutoff, top performing methods based on prediction standard error of the mean (\pm s.e.m.) are bolded. A cutoff of 0.95 indicates that only the top 5% most confident predictions are considered. Full confidence plots for all datasets are shown in Figure 7.3 and Figures 7.7, 7.10, 7.11. Mean \pm s.e.m. (RMSE for all D-MPNN models, MAE for atomistic); $n = 10$ independent trials for lower-N datasets, $n = 5$ independent trials for higher-N datasets.

versa for an underconfident model. Here, we explore the calibration properties of evidential learning for molecular property prediction.

Evidential learning methods introduce regularization terms that minimize model evidence in instances of high predictive error [275], [276]. Specifically, in the evidential regression method, the training loss takes the form:

$$L(x) = L_{NLL}(x) + \lambda L_{REG}(x)$$

where λ controls the strength to which overconfident predictions are penalized by the regularization term L_{REG} . Thus, λ provides a tunable hyperparameter capable of modulating the calibration of any model trained with evidential loss (Figure 7.4A). To investigate the effect of λ on uncertainty calibration, we trained separate models with different regularization strengths and computed empirical calibration curves that compare the fraction of test set points that fall within a credible interval against the fraction of test set points expected to fall within the predicted credible interval [260]. As expected, trained evidential D-MPNNs move from overconfident to underconfident regions as λ is increased, as shown on the Delaney dataset (Figure 7.4B) and across other lower-N datasets (Figure 7.12).

To quantify the calibration accuracy, we calculated the area between the observed calibration curve and the parity line (perfect calibration) for each value of λ evaluated across all lower-N datasets (Figure 7.4C). For all lower-N datasets except QM7, there exists a value of λ at which evidential regression is more calibrated than the ensemble baseline (Figure 7.4C). Based on these results, we choose a default of $\lambda = 0.2$, as cutoff RMSE is robust to small changes in λ (Figure 7.13), and use this value for all evaluations unless stated otherwise. This regularization strength yields predictions calibrated at least as well as those of ensemble-based methods across all datasets and tasks tested (Figures 7.7, 7.11, 7.14). While methods have been developed to re-calibrate and augment uncertainty predictions [284], the ability to tune λ via a hyperparameter search with the evidential regression formula-

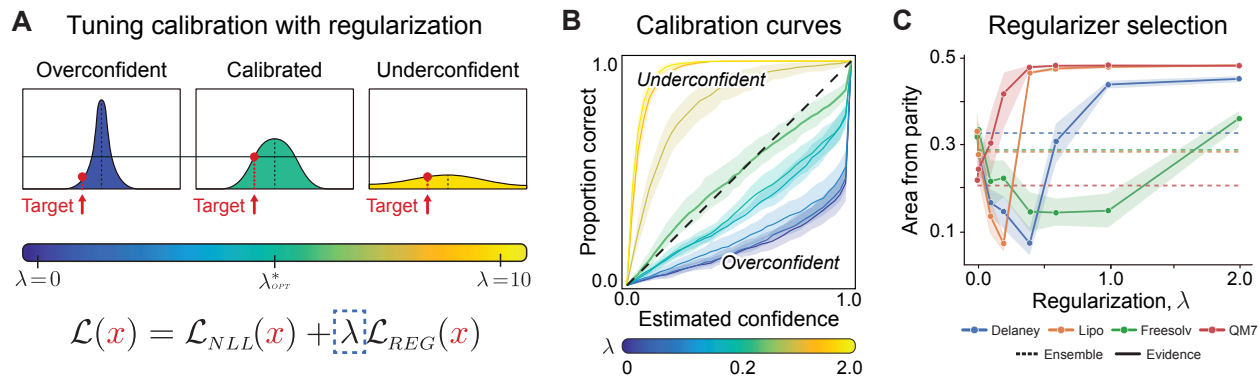


Figure 7.4: **Tunability of the evidential uncertainty.** (A) The evidential regression method can be fine tuned with a single hyperparameter, λ , in order to achieve more calibrated predictions for a given dataset. (B) Estimated confidence (cumulative probability) against the observed proportion correct for an evidential D-MPNN evaluated on the Delaney dataset. Dotted line represents perfect calibration. Mean \pm 95% c.i., $n = 5$ independent trials. (C) Area between the observed calibration curve and the perfect calibration line across several lower-N datasets for evidential D-MPNNs trained with varying λ . Dotted lines represent calibration of an ensemble of models. Mean \pm 95% c.i., $n = 10$ independent trials.

tion presents an additional, attractive option for chemical science practitioners to quickly calibrate uncertainty before applying general purpose re-calibration techniques.

7.3.3 Application I: Uncertainty-guided learning

Having verified that evidential uncertainties were well-calibrated to errors on property prediction tasks, we next sought to use these uncertainties to guide learning towards improved sample efficiency or accelerated molecular optimization. Concretely, in this section we investigate two applications, active learning and Bayesian optimization, that utilize UQ to intelligently prioritize sample acquisition (Figure 7.5A).

Active learning for sample efficient training

As a first validation of the utility of evidential uncertainties for guided learning, we turned to the QM9 dataset [280], a standard dataset for molecular property prediction that captures geometric, energetic, electronic, and thermodynamic properties, and asked whether uncertainty-guided sample acquisition could yield a more sample-efficient learning process. For QM9 active learning experiments, data acquisition was simulated as iterative selection from the library repeated six times after initialization with a random 15% subset of the training data. At each step, the uncertainty was evaluated across the remainder of the training data (i.e., samples that had not yet been selected). For explorative selection, the k samples with the greatest estimated uncertainties at each iteration were added to the training set, and the model was subsequently re-trained using this expanded dataset and then evaluated

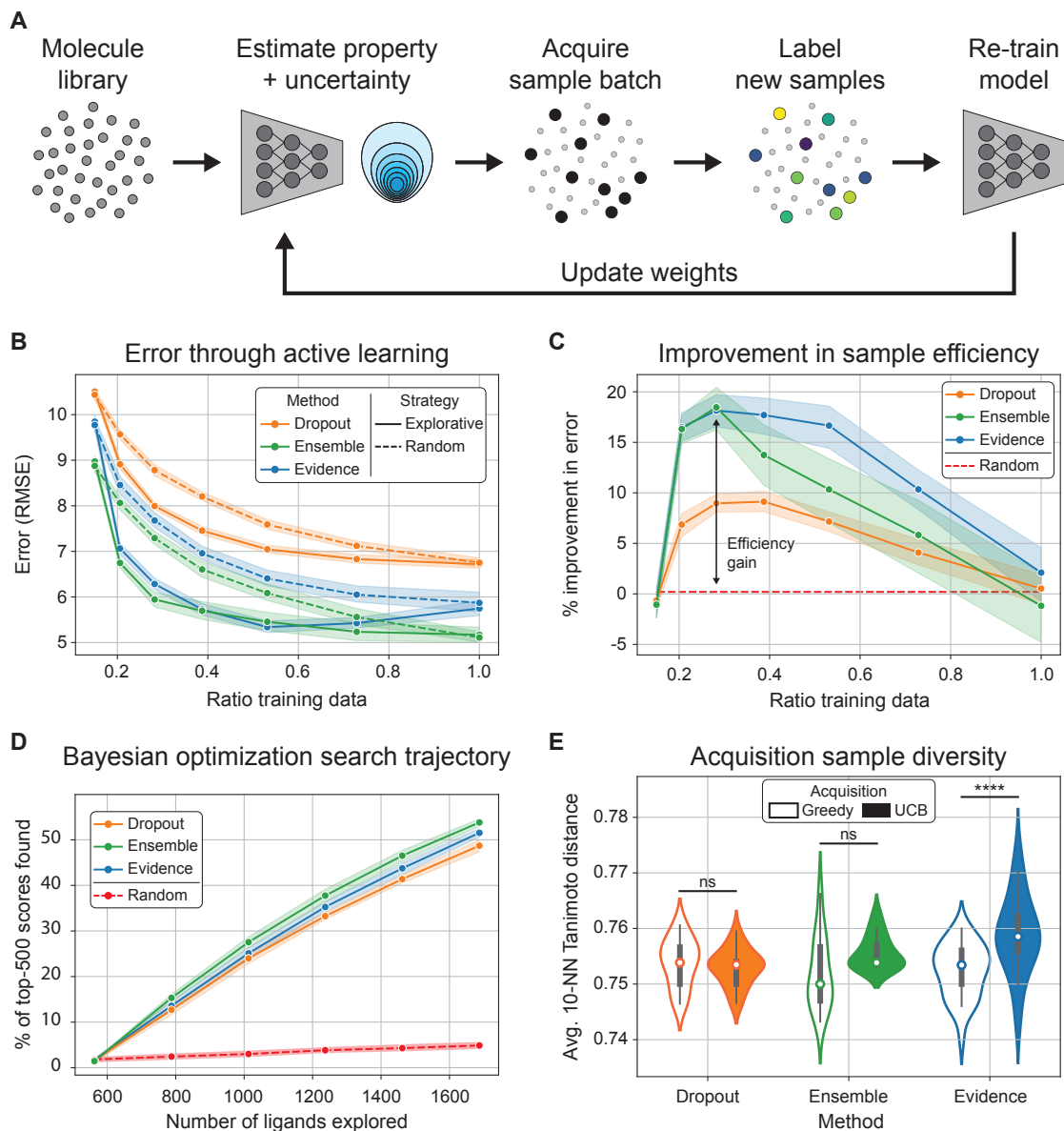


Figure 7.5: **Evidential active learning and Bayesian optimization.** (A) Experimental scheme. (B) Active learning with explorative (solid) versus random (dashed) sampling for D-MPNN evaluated on the QM9 dataset. Mean \pm 95% c.i., $n = 10$ independent trials. (C) Change in sample efficiency for explorative acquisition in (B), evaluated as the percent decrease in predictive error relative to a randomly-selected training set. (D) Bayesian optimization performance on Enamine 50k data, measured by the percentage of top-500 scores found as a function of the number of ligands explored. Solid traces represent an upper confidence bound (UCB) acquisition strategy. Mean \pm 95% c.i., $n = 10$ independent trials. (E) Average 10-nearest training set neighbors (10-NN) Tanimoto distance for batch samples after the first round of acquisition in Bayesian optimization experiments. Dots represent median; bars represent interquartile range; lines represent upper and lower adjacent values. $n = 10$ independent trials, two-tailed unpaired t-test, **** $P < 0.0001$.

on a held out test set (Figure 7.5A). For all evaluations, random sample selection served as a baseline for each uncertainty quantification method considered.

We find that active selection based on evidential uncertainties yields significantly improved sample efficiency, reaching the same level of performance of the full training dataset with over 60% less data (Figure 7.5B). Further, acquisition using evidential uncertainties results in increased data efficiency relative to dropout-based selection. For example, to achieve an RMSE of 7.0, evidence-guided models required an average of 21% of the entire training data compared to 55% for dropout-guided models (Figure 7.5B). We observe, consistent with prior literature, that stochastically trained networks (e.g., trained with dropout) also suffer from a baseline performance drop relative to their deterministically trained counterparts, even when all data is considered. As expected, ensemble-based selection shows the greatest improvement over random selection, which is consistent with the advantages that training multiple independent models affords. However, this comes at a significant computational cost, given that multiple independent models must be completely re-trained at each active learning step. With growing interest in using uncertainties to inform full molecular dynamics simulations [266] and to decide when to perform density functional theory simulations [265], retraining 5 to 10 times as many models can drastically increase the expense and overhead of a simulation. In contrast, evidential learning enables resource-efficient uncertainty estimation at the cost of just a single deterministic model, while still achieving increased training efficiency relative to random acquisition at a level nearly on par with ensemble-based selection (Figure 7.5C). Specifically, while model ensembling achieves a better overall error, the evidential method exhibits superior improvements in sample efficiency across different stages of learning, attaining up to a 18% improvement in error relative to the random baseline (Figure 7.5C). At any stage in learning, evidence-guided explorative sampling yields an equal to or greater drop in error over random acquisition when compared to either dropout or ensembling.

Bayesian optimization for accelerated molecular discovery

Instead of acquiring samples solely based on uncertainty, as with purely explorative active learning, Bayesian optimization provides a framework to discover high performing compounds (e.g., those with desired molecular properties) from a large search space by incorporating both predicted property scores and uncertainties to guide sample acquisition [282], [285]. In this scheme, uncertainties can be used to explore the search space more conservatively and to broadly enhance the overall diversity of acquired sample batches [286].

To this end, we investigated the utility of evidential uncertainty for Bayesian optimization settings, where the aim is to rapidly discover compounds with target molecular properties. We turned to the ligand docking dataset, previously benchmarked in Figure 7.3D, of 50,240 molecules docked *in silico* against thymidylate kinase [282]. Given this library of ca. 50k molecules, we aim to identify those with the best ligand docking scores by only observing ground truth docking scores for a small subset of the library. Data acquisition is initiated by training on a random 1% subset (ca. 500 molecules) and then simulated as the iterative selection of new samples based on an upper confidence bound (UCB) acquisition function according to a given UQ method. In these experiments, a D-MPNN is used as the surrogate model by which docking scores and uncertainties are estimated.

For all three UQ methods evaluated (dropout, ensemble, and evidence), UCB acquisition yields clear improvements over the random baseline, representative of a brute-force search, as measured by the percentage of top-500 (ca. top 1%) of scores found as a function of the number of compounds explored (Figure 7.5D). Specifically, the evidential method discovers over 50% of the top-500 docking molecules from the pool of 50k molecules after exploring fewer than 2k molecules (less than 4% of the search space). Similar to the active learning experiments, we also observe that the evidential method outperforms dropout sampling, but does not exceed the performance of ensembling (Figure 7.5D, Table 7.3). While previous studies on this dataset have shown that using greedy sampling based upon predicted docking score outperforms UCB [282], we additionally evaluate the structural diversity of the newly acquired pool, relative to the training set, in both greedy and UCB sampled molecules after one round of acquisition (Figure 7.5E). Relative to its greedy baseline, the evidential UCB method results in a statistically significant increase in the average Tanimoto distance between sampled molecules and their respective 10-nearest training set neighbors, while the dropout- and ensemble-based UCB methods do not (Figure 7.5E). Together, these results support the use of evidential uncertainties within Bayesian optimization frameworks for accelerated virtual screening and molecular discovery.

7.3.4 Application II: Uncertainty-guided inference for virtual screening

Though virtual screening is a common tool in computer-aided molecular discovery, all *in silico* predictions of QSAR models must be experimentally validated, and often only a small fraction of predictions or candidates nominated by QSAR models hold true in the real world [254], [255]. Therefore, there remains a need for integrated methods that can help ensure the robustness of QSAR predictions. Fast and scalable UQ methods have the potential to meet this need by guiding *in silico* discovery towards molecular candidates associated with greater predictive confidences, based on the hypothesis that high confidence candidates are better suited for downstream experimental validation. To this end, we next investigate the potential for evidential deep learning to discover high confidence drug candidates in retrospective virtual screening campaigns. Concretely, we consider a virtual screening pipeline for antibiotic discovery [43] and demonstrate how evidential uncertainties can be integrated to more accurately prioritize drug repurposing candidates for use as antibiotics by additionally filtering large screening libraries based on confidence in addition to predicted activity.

To achieve this, we develop a framework for uncertainty-guided prioritization in virtual screening (Figure 7.6A). A large, labeled dataset of small molecules is used to train an evidential model which is in turn applied to a smaller, unlabeled discovery dataset to predict both molecular properties as well as uncertainties. From these predictions, candidate molecules are subsequently ranked by their associated property values, and then filtered further based on confidence thresholds ranging from the 50th to 100th percentiles of greatest predictive confidence (i.e., lowest uncertainties). Finally, among this filtered subset, experimental hit rates (i.e., correlation of true versus predicted activities) are determined either retrospectively, as in this work, or prospectively to assess the relative benefit of uncertainty-guided prioritization versus naive nomination (i.e., without confidence filtering).

To concretely demonstrate the utility of this approach, we considered the question of antibiotic discovery and leveraged a recent dataset [43] of small molecules and their *in vitro* growth inhibition against *Escherichia coli* (*E. coli*, measured as OD₆₀₀) (Figure 7.15A). We frame the prediction as a regression problem and train a D-MPNN with evidential loss and outputs on this dataset of 2,335 small molecules to predict their OD₆₀₀ values. The accuracy of the resulting model is shown on a held-out validation set (Figure 7.6B).

Next, extending off the virtual screening pipeline presented by Stokes et al. [43], we applied the trained evidential D-MPNN model to an independent, unlabeled discovery dataset with the aim of identifying high confidence candidate antibiotics. We leverage the Broad Drug Repurposing Hub [287] as the discovery dataset and generate predictions for both estimated antibiotic activities as well as learned evidential uncertainties (Figure 7.15B). To begin to understand how uncertainties scale and extend to this discovery dataset, we visualize the structural overlap in chemical space between molecules in the discovery (Broad) dataset compared to those in the training dataset, and annotate molecules in the discovery dataset with their estimated evidential uncertainties (Figure 7.6C). Qualitatively, this analysis revealed select regions of chemical space associated with higher evidential uncertainties that also exhibit less overlap with the training set, consistent with the expected inflation of uncertainties for out-of-distribution or distribution shifted domains. Furthermore, comparison of predicted growth inhibition to evidential uncertainty demonstrates that predicted active molecules (lower predicted OD₆₀₀) trended towards higher uncertainties (Figure 7.15B), an observation consistent with the stark imbalance and skewness of the training set (Figure 7.15A).

We then utilized evidential uncertainties to prioritize high confidence candidate antibiotics from the discovery dataset, with the goal of identifying molecule sets with high experimental hit rates, i.e., high likelihoods of having true growth inhibitory activity in the real world (Figure 7.6A). To this end, following prediction of both antibiotic activity and the associated evidential uncertainty, we rank molecules in the discovery dataset according to their predicted antibiotic activities (i.e., lowest to highest predicted OD₆₀₀, where lower is better), and then select the top k ranking molecules based on predicted activity, as outlined in a previous virtual screen on this dataset [43]. We subsequently filtered the resulting set of k molecules filtered based on confidence estimates for varying confidence thresholds. Specifically, for a given confidence threshold p , molecules with estimated confidences below the associated p^{th} percentile are removed from the list of top k molecules, with p ranging from the 50th to 100th percentiles of greatest predictive confidence.

Experimental hit rates (true OD₆₀₀ < 0.2) for these model-nominated compounds were then determined using the subset of candidates for which growth inhibitory activity against *E. coli* had been determined [43] (Figure 7.15C). This analysis revealed that augmenting network predictions with confidence-based filtering with evidential uncertainties can increase the experimental hit rate relative to that of an unfiltered set of candidates (Figure 7.6D). Increasing confidence percentiles enriched the candidate set for experimental hits, from a hit rate of 78% for naive filtering to over 95% after confidence filtering using our evidential method (Figure 7.6D). While filtering with ensemble-derived uncertainties also increased the experimental hit rate above baseline, this difference was not as great as the relative increase provided by the evidential method. These results suggest that evidential uncertainties can be used to prioritize high confidence drug candidates in virtual screens in order to ultimately

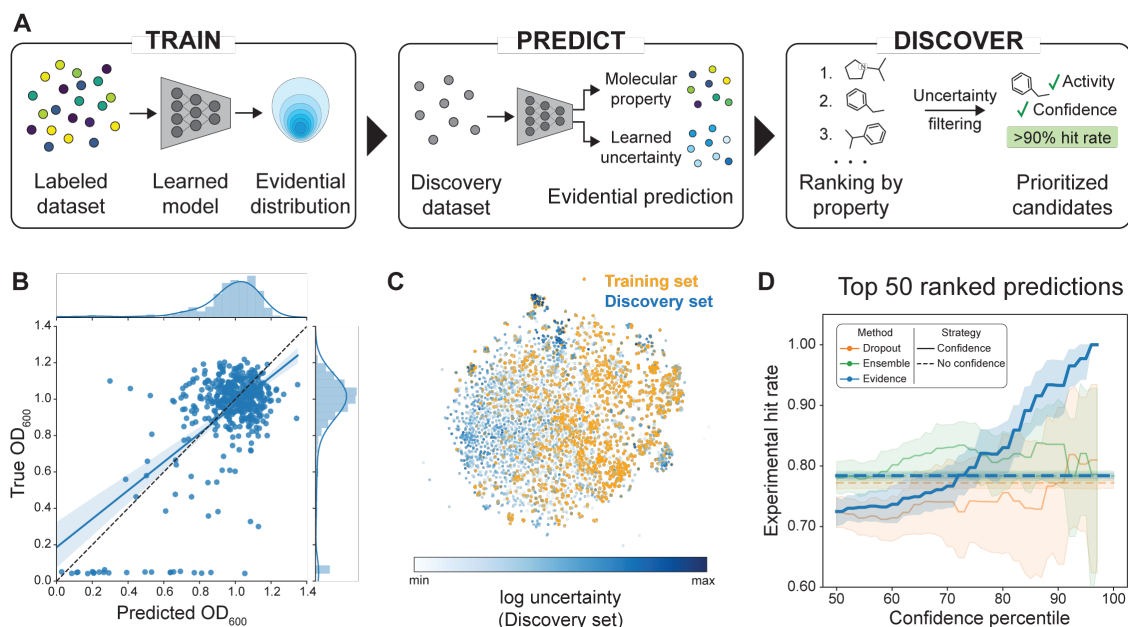


Figure 7.6: **Uncertainty guided nomination in virtual screens.** (A) Experimental framework in which trained uncertainty-aware models are deployed on a discovery dataset, and molecules are ranked based on predicted properties. Uncertainty filtering is used to prioritize candidates among the top ranking molecules. (B) Performance of evidential D-MPNN after training to predict *E. coli* growth inhibition. (C) t-SNE visualization of training set (orange) and discovery dataset (Broad library), colored by predicted evidential uncertainties (blue). (D) Application of confidence filters to prioritize sets of antibiotic candidates with high experimental hit rates. Mean \pm 95% c.i., $n = 10$ independent trials.

guide discovery towards greater likelihoods of experimental success.

7.4 Discussion

7.4.1 Contributions

In this work, we establish evidential deep learning as a scalable, efficient, and easy-to-use uncertainty quantification method for molecular property prediction in the chemical and physical sciences. By integrating our algorithm into both message passing and atomistic networks operating on 2D graphs and 3D conformers, respectively, we demonstrate its modularity across different network architectures and its applicability across a range of tasks in both lower-N and higher-N settings. Through benchmarking experiments against model ensembling and dropout sampling methods, we show that our evidential algorithm exhibits strong calibration and yields uncertainties that scale with prediction errors, supporting the utility of our method for prioritizing candidate molecules from large screening libraries. Furthermore, we validate the utility of evidential deep learning for uncertainty-guided learning

and compound prioritization in virtual screening. We find that evidential uncertainties can effectively guide sample acquisition to improve training efficiency and to accelerate virtual screening in active learning and Bayesian optimization settings. Finally, by leveraging an evidential message passing network to identify high confidence candidate antibiotics, we show that evidential uncertainties can be used to direct retrospective virtual screening campaigns towards compound sets with increased experimental validation rates.

7.4.2 Advantages of evidential learning in chemical science

The evidential deep learning framework offers several advantages relative to existing UQ approaches for neural models in the chemical sciences. In contrast to other methods such as Bayesian neural networks that require modifying architectures to output probability distributions over network weights [288], our algorithm can be incorporated into an existing network architecture by modifying the loss function and the network’s final output layer. Furthermore, our method presents key scalability and efficiency advantages over sampling-based approaches for QSAR UQ, namely traditional model ensembling [262] and dropout sampling [263], which necessitate training and/or evaluation of multiple surrogate models in order to obtain approximations of epistemic uncertainty. While widely used, these methods can incur high computational costs which may be prohibitive in settings that are resource-constrained, that require iterative training, or that use large networks or large datasets. Our method overcomes this limitation by directly modeling a higher-order probability distribution over the likelihood function and requires only a single forward pass through a network to obtain uncertainty estimates [275], [276].

7.4.3 Opportunities and applications in molecular property prediction

Because of its efficiency and ease of use, evidential deep learning may be particularly relevant to uncertainty-guided virtual screening of large scale chemical libraries. Virtual screening workflows often involve exhaustive prediction of the properties and performance of compounds in large virtual libraries prior to prioritization of candidates for experimental validation [279]. In this setting, evidential deep learning may be used to efficiently obtain uncertainty estimates to understand when the predictions of QSAR models may not be trusted, and furthermore to accelerate downstream sample annotation or acquisition, for example via active learning. We envision that evidential learning can be incorporated as an efficient, modular UQ method for virtual screening and compound discovery campaigns.

Opportunities for future work also exist in the context of neural networks as surrogate models for quantum mechanical and molecular dynamics simulations [289], [290], where there is increasing interest in using uncertainties to actively guide simulation experiments to determine when machine learned predictions can no longer be trusted [265], [266]. Furthermore, though in this work we integrated evidential deep learning into the SchNet architecture [291], continued development and integration of evidential methods into new atomistic machine learning architectures [267] could help advance their deployment for prediction of potential energy surfaces and quantum mechanical properties, tasks that demand computational scalability and efficiency.

7.4.4 Scope and future work

While evidential deep learning provides key advantages over existing methods for UQ in neural models, there are several considerations that motivate opportunities for future work. First, the vast majority of our analyses here focus on regression problems in which networks are trained to predict a continuous molecular property. Evidential deep learning models were originally presented in the setting of multiclass classification [275], and as such are also applicable to these domains as well. However, since the natural form of many molecular properties is continuous (not discrete), we focus our analysis in this work on the applicability of evidential methods specifically in the regression domain. We hypothesize that the benefits of evidential UQ for classification will also be apparent in *multiclass* settings, where a single input is being classified as one discrete class from a set of options, such as in protein secondary structure or amino acid prediction, among other applications [292]. Future research to this end will be important to validate the generality of evidential UQ for classification settings.

Furthermore, in this work, we focus on several metrics to evaluate the quality of uncertainty estimates: confidence percentile cutoff errors, Spearman rank correlation coefficients between error and uncertainty, and miscalibration area. We also acknowledge that these UQ methods would be best evaluated in terms of their performance on realistic applications, and accordingly demonstrate the use of evidential uncertainties for efficient Bayesian optimization, active learning, and virtual screening. Thus, future work is needed to identify and formulate other impactful applications where effective UQ methods yield improvements in downstream performance. This could help solidify the utility of uncertainty for machine learning in the chemical, biological, and physical sciences, where guarantees in not only model performance but also confidence are ultimately needed for wide-scale adoption.

Finally, in this work we use our evidential UQ method to guide learning and compound screening in the retrospective setting, through active learning experiments and evaluation on an antibiotic discovery dataset. While these evaluations support the utility of evidential deep learning, and UQ more broadly, for similar analyses, they remain retrospective. Further work to explore the utility of evidential uncertainties in the prospective setting [293], for example to identify new, high confidence drug candidates that may in turn be experimentally tested in the real world, could help facilitate the adoption of UQ approaches in discovery and engineering pipelines.

7.5 Conclusion

In summary, we have developed a flexible, scalable, and efficient approach to uncertainty estimation in neural networks for molecular property prediction in the chemical and physical sciences. We demonstrated that evidential deep learning provides well-calibrated uncertainties in structure-property prediction, and validated its relevance to uncertainty-guided active learning and to prioritization of candidates in virtual screening. We expect that evidential deep learning, which can readily be incorporated into existing network architectures and be applied to a variety of predictive learning tasks, could help facilitate the robust and reliable deployment of uncertainty-aware neural models for molecular property prediction, discovery, and design.

7.6 Methods

All code to reproduce experiments and results can be found at <https://github.com/aamini/chemprop>. All scripts were written in Python; PyTorch [174] was used for building all machine learning architectures; RDKit [158] was used for various cheminformatics calculations

7.6.1 Evidential deep learning formulations

Evidential deep learning approaches seek to directly learn prediction uncertainties by formulating learning as an evidence acquisition process [275], [276]. This is achieved by training models to infer the parameters of a higher-order *evidential* distribution that models the evidence behind individual predictions. That is, individual observations of training examples lend support to this higher-order distribution, such that the predictions of the neural network learner are represented as a distribution over the prediction likelihood function itself. Estimates of uncertainty are then formulated using the parameters of the learned evidential distribution and thus can be obtained directly from a single forward pass through the model.

Evidential learning is achieved through two modifications to a standard forward prediction model. First, the network’s output layer is modified to output the parameters of the evidential distribution, rather than a point estimate of a target label. Second, the resultant model is trained with a specific loss function that jointly maximizes the model’s fit to the data and also minimizes its evidence on errors, i.e., increases uncertainty when predictions should not be trusted.

Evidential learning for regression In regression, we are given a dataset of paired training examples $\mathcal{D} = \{x_i, y_i\}_{i=1}^N$ where the targets are assumed to be drawn i.i.d. from a Gaussian distribution with unknown mean and variance $\theta = \{\mu, \sigma^2\}$. We seek to probabilistically estimate the mean and variance assuming that the mean is drawn from a Gaussian and the variance is drawn from an Inverse-Gamma distribution. The joint higher-order, evidential distribution is thus represented as a Normal-Inverse-Gamma. Specifically, the Normal-Inverse-Gamma distribution, $p(\theta|\mathbf{m})$, which is also the conjugate prior to the Gaussian [294], is parametrized by $\mathbf{m} = \{\gamma, v, \alpha, \beta\}$ and represents a distribution over $\theta = \{\mu, \sigma^2\}$. Therefore, in this work, the final layers of evidential D-MPNN and atomistic networks were modified to output these Normal-Inverse-Gamma hyperparameters. Thus, the network has 4 outputs for every target task. Given a Normal-Inverse-Gamma distribution, the prediction and uncertainty are formulated according to the distribution moments:

$$\underbrace{\mathbb{E}[\mu]}_{\text{prediction}} = \gamma, \quad \underbrace{\text{Var}[\mu]}_{\text{uncertainty}} = \frac{\beta}{v(\alpha - 1)}.$$

Evidential models are trained using a dual-objective loss $\mathcal{L}(x)$ that consists of two loss terms, to both maximize model fit according to the negative log-likelihood and regularize evidence on errors:

$$\mathcal{L}(x) = \mathcal{L}^{\text{NLL}}(x) + \lambda \mathcal{L}^{\text{R}}(x)$$

where $\mathcal{L}^{\text{NLL}}(x)$ is the negative log-likelihood and $\mathcal{L}^R(x)$ is an evidence regularizer [276]. The regularization coefficient λ controls the strength of uncertainty inflation relative to model fit. All evidential models were trained according to this loss function, with λ values specified in the figure captions and corresponding Methods sections. We refer to the work of Amini et al. for more details on the evidential regression formulation [276]. We also note that the evidential method has been demonstrated in the context of discrete, multiclass classification problems [275], despite the focus of this work (along with relevant prior literature in molecular property prediction) being on continuous regression tasks.

7.6.2 Network architectures

To show its broad applicability in molecular property prediction, evidential regression was integrated into networks operating on 2D molecular graphs and 3D conformers – directed message passing neural network (D-MPNN) and atomistic neural network models, respectively.

Directed message passing neural networks To investigate performance on 2D molecular graphs, evidential methods were integrated into a state-of-the-art D-MPNN model [42]. The D-MPNN architecture is a variant of a message passing neural network (MPNN). MPNNs operate on molecular graphs to first learn an encoded molecular representation (i.e., molecule-level feature vector) by passing “messages” between atoms and/or bonds and their direct neighbors. These messages build up a hidden state for each atom and/or bond, and repeated message passing iterations yield a molecule-level feature vector. A feed-forward network operating on this feature is used to produce a task-specific representation of an input molecule.

D-MPNN models were implemented in PyTorch [174] within the Chemprop library [42]. D-MPNNs were implemented using standard settings: messages passed on directed bonds, messages subjected to ReLU activation, a learned hidden dimension of 300, 3 layers, no dropout, and the output of the message-passing phase fully connected to the output layer. For evidential regression models, the final output layer was modified to infer a single evidential distribution for each task, with each task parametrized by four outputs (e.g., prediction of 12 tasks uses 48 outputs). Models were trained using the Adam optimization algorithm. Target values were normalized with a standard scaler for training. For evaluations, model state was reloaded from the epoch with the lowest validation score after training was completed.

Atomistic neural networks To investigate performance on 3D conformer representations, evidential regression was integrated into the end-to-end atomistic neural network SchNet [291]. Rather than operating only on the 2D graph, SchNet instead builds internal representations of a molecule using 3D coordinate positions of each atom as inputs to the model. SchNet employs “continuous filters”, where information is shared between molecules not based upon discrete edges but rather continuous spatial distance between molecules (Figure 7.1). This model architecture has proved beneficial for predicting energies and forces, as it is rotationally invariant and equivariant [291]. Similar to the message passing networks,

SchNet alternates between transforming atom-wise representations individually and integrating interaction information. At the end of these internal hidden layers, information at each atom is aggregated through summation to produce a fixed dimension hidden representation, upon which a single feed forward layer is applied to produce an output value. In the case of energy predictions, SchNet predicts energy at each atom separately and adds all energies to predict the energy for the molecule. We refer the reader to the work of Schutt et al. for more details [291], [295].

The final layer of SchNet was modified to infer a single evidential distribution over the task of predicting U_0 for the QM9 dataset, outputting four values rather than one, corresponding to the parameters of the higher-order evidential distribution.

Unlike in the Chemprop library, where normalization is conducted prior to any training, the SchNet model uses a customized scaler that normalizes both according to the target values *and* a precomputed “atomref” value. Instead of transforming the target values and training the model with modified target values, SchNet adds these normalization numbers to outputs for each atom separately as part of its final layer by default. SchNet builds up final energy values for each atom m , $G(x_i)_m$, and the final prediction for molecule x_i is computed according to:

$$\hat{y}_i = \sum_{m \in \text{Atoms}(x_i)} [\sigma_y G(x_i)_m + \text{atomrefs}(x_{i,m}) + \mu_y]$$

where

$$\mu_y = \frac{1}{N} \sum_{i=1}^N \frac{1}{|x_i|} \left(y_i - \sum_{m \in \text{Atoms}(x_i)} \text{atomrefs}(x_{i,m}) \right)$$

$$\sigma_y = \sqrt{\frac{1}{N} \sum_{i=1}^N \left(\mu_y - \frac{1}{|x_i|} \left(y_i - \sum_{m \in \text{Atoms}(x_i)} \text{atomrefs}(x_{i,m}) \right) \right)^2}$$

Here, SchNet tries to learn residual values at each atom. Because the evidential regression method is most easily adapted for predicting and training on data with standard normalization procedures, a custom standardization method was implemented to incorporate atomref values, which we found to be empirically helpful. Specifically, the value U_0 for each input molecule was adjusted before inference according to:

$$y'_i = \frac{y - \mu_y - \sum_{m \in \text{Atoms}(x_i)} \text{atomrefs}(x_{i,m})}{\sigma_y}$$

At inference time, the inverse of this scaler was computed to predict \hat{y} . Accounting for such discrepancies in the scaling function could provide opportunities for future development of atomistic neural network architectures.

7.6.3 Datasets

Lower-N 2D datasets The lower-N 2D datasets used in this study were extracted from MoleculeNet [278] and used as prepared by Wang et al. [42]. For regression tasks (Table 7.1), datasets of aqueous solubility (Delaney), solvation energy (freesolv), lipophilicity

(lipo), and atomization energy (QM7) were evaluated. The datasets were split randomly using a 80/10/10 split for training/validation/testing. SMILES strings were used as input to D-MPNN models.

Lower-N 2D TDC Datasets In addition to the MoleculeNet lower-N datasets, evidential regression was additionally evaluated on three datasets from the Therapeutics Data Commons [238] (TDC). For regression tasks (Table 7.2), datasets of hepatocyte clearance (“clearance”) [296], plasma protein binding rates (“PPBR”) [296], and the lethal dosage of drugs (“LD50”) [297] were used. The datasets were extracted from the TDC and split randomly using a 80/10/10 split for training/validation/testing. SMILES strings were used as input to D-MPNN models.

Higher-N 2D datasets For the 2D setting, the QM9 [280], [298] dataset with SMILES strings input was extracted from MoleculeNet [278]. For both benchmarking and active learning D-MPNN experiments using this dataset, all 12 output tasks of the QM9 dataset, which reflect computer-generated quantum mechanical properties [278], [280], were predicted. A ligand docking dataset based on Enamine’s Diversity Collection of 50,240 molecules was used as a second higher-N 2D dataset. Target values consisted of docking scores of compounds against thymidylate kinase (PDB ID: 4UNN31) using AutoDock Vina [283]. This dataset was used as prepared by Graff et al. [282].

Higher-N 3D (atomistic) datasets For the 3D setting, a variation of the QM9 dataset was utilized wherein atomic coordinates, not molecular graphs generated from SMILES strings, were used as inputs. Atomic coordinates correspond to the coordinates of molecular conformers in 3D space. In this setting, the single task of total formation energy at 0K, U_0 [295], was predicted for a given molecule input.

Antibiotic discovery datasets For virtual screening experiments for antibiotic discovery, D-MPNN models were trained on a dataset of 2,335 small molecules and their *in vitro* growth inhibitory activity against *Escherichia coli*, as generated and reported by Stokes et al. [43]. In this dataset, growth inhibitory activity is reported as endpoint OD₆₀₀, where lower OD₆₀₀ values correspond to stronger growth inhibitory activity, and models were trained to predict this as a continuous target (i.e., formulated as a regression problem). The Broad Drug Repurposing Hub [287] was used as a discovery dataset, as prepared by Stokes et al. [43]. Model predictions were compared to empirically determined growth inhibitory activity against *E. coli* for a subset of molecules from the Broad Drug Repurposing Hub, as measured by Stokes et al. [43].

7.6.4 Uncertainty quantification baselines

Evaluations are focused on regression tasks, defined by a dataset \mathcal{D} containing data points (x_i, y_i) where $y_i \in \mathbb{R}$ is a scalar-valued target property and x_i is a molecule representation, represented as either a SMILES string or a set of coordinates for the 2D or 3D settings, respectively.

For baselines, we use gold-standard epistemic uncertainty quantification methods that rely on sampling, i.e., creating a set of predictions that together constitute an ensemble from which estimates of predictive variance can be obtained. Specifically, a set of predictions $\mathcal{E} = \{G_1(x), G_2(x), \dots, G_n(x)\}$, where each $G_i(x)$ is an inference sample, is obtained such that the individual samples (e.g., individual model predictions) can yield a final prediction defined by:

$$\hat{G}(x) = \sum_{G \in \mathcal{E}} \frac{G(x)}{n}.$$

As previously proposed [262], [263], from the multiple samples obtained from this set of models, the uncertainty $U(x)$ is defined as the variance across predictions:

$$U(x) = \sum_{G \in \mathcal{E}} \frac{(\hat{G}(x) - G(x))^2}{n}.$$

Traditional model ensembling For the ensemble baseline, distinct models $G_i \in \mathcal{E}$ were trained on different splits of the same training data and initialized with different sets of randomly selected weights, as previously proposed [262]. Greater variance among model outputs reflects greater uncertainties, due to the fact that for out-of-distribution regions not well represented in the training data, each ensemble member will be more significantly affected by its initialization, ultimately resulting in more variable predictions. The computational cost of this approach scales linearly with the size of the ensemble and clearly exceeds the cost of training a single model. All evaluations utilized an ensemble size of 5.

Monte Carlo dropout sampling For the dropout baseline, a single model G is trained with dropout, in which individual network weights are randomly set to zero at every training step with probability p , also known as the “dropout rate”. At inference time, a set of predictions \mathcal{E} is obtained for an input x_i by application of randomly-generated dropout masks to a trained model G . This strategy approximates Bayesian inference and can thus be used to obtain prediction samples from which uncertainty can be estimated [263]. All evaluations utilized a dropout rate of 0.2 and a set size of 5.

7.6.5 Uncertainty benchmarking experiments

Each dataset was randomly partitioned using an 80/10/10 split for training, validation, and testing set splits respectively. For D-MPNN experiments, test set target values were normalized using statistics from the train set such that the train set target values have mean 0 and standard deviation 1 in each case (atomistic normalization previously described separately).

For the D-MPNN models, each model was trained for 100 epochs using the Adam optimizer with default parameters, and the best model was selected based upon validation loss. We use the default Chemprop architecture and training procedure parameters as set by Yang et al. [42], including a Noam learning rate scheduler with final learning rate of 10^{-4} , batch size of 50, hidden size of 300, depth of 3, and ReLU activations.

For the atomistic neural network models, the default SchNet architecture parameters and preparations of the U_0 calculations from the QM9 dataset were used. For training, the same training procedure as in **Chemprop** (optimizer, learning rate, and learning rate scheduler) was used, with the learning rate in the Noam scheduler modified to 2×10^{-4} , rather than 1×10^{-4} , for greater stability in early model epochs.

Test set predictions and corresponding uncertainties were generated in each experiment for downstream analysis.

Error vs. confidence cutoffs Test set predictions from each model run were sorted by uncertainty, such that r_i represents the index of the test set molecule, x_{r_i} that has the i^{th} highest predictive uncertainty (e.g., x_{r_1} is predicted with highest uncertainty and x_{r_n} with most confidence where n is the total number of molecules evaluated). For every value of i , we compute the error for the set of all predicted test set molecules of $\{x_{r_j} : j \geq i\}$. We compute cutoff mean average error (MAE) and root mean squared error (RMSE) for different cutoff values i , corresponding to different confidence cutoffs (e.g., 50% confidence cutoff can be computed by setting $i = \lceil 0.5n \rceil$):

$$\text{MAE}_i = \frac{1}{n-i} \sum_{j \geq i}^n |y_{r_j} - G(x_{r_j})|$$

$$\text{RMSE}_i = \sqrt{\frac{1}{n-i} \sum_{j \geq i}^n (y_{r_j} - G(x_{r_j}))^2}$$

Confidence cutoff errors were computed at even intervals of 30, exclusive, for all datasets except for the lower-N datasets, in which case cutoff errors for *all* values of i were computed.

This procedure was repeated for different random starts of the model and random training/validation/testing splits, corresponding to independent experimental trials ($n = 10$ for lower-N data, $n = 5$ for all else). Cutoff values at each confidence percentile were computed separately for each random initialization in order to estimate standard deviations over multiple trials at each confidence percentile.

In the multitask setting, cutoff errors were computed separately across each task to avoid artifacts due to averaging uncertainties across tasks. Thus, for $n = 5$ trials on the QM9 dataset, cutoff errors were computed separately for each of the $d = 12$ tasks in each of the n trials, leading to $n * d = 60$ computed cutoff error values at each confidence percentile. Performance on separate tasks in QM9 is shown in Figure 7.10.

Spearman’s rank correlation coefficient We desire that, on average, predictions for molecules made with higher certainty should also be more accurate. Therefore, Spearman’s rank correlation coefficient was used as an additional metric to assess the ability of models to rank errors, and was computed following the procedure outlined in Hirschfeld et al. [264]. Spearman’s rank correlation coefficient is defined by creating two vectors L_1 and L_2 and corresponding rank vectors r_{L_1} and r_{L_2} that sort the dataset in ascending order. The correlation measures the agreement between these two ranking lists:

$$\rho(L_1, L_2) = \frac{\text{cov}(r_{L_1}, r_{L_2})}{\sigma(r_{L_1})\sigma(r_{L_2})}$$

If these two lists are perfectly correlated, then $\rho = 1$, whereas if they are perfectly inversely correlated, $\rho = -1$. We desire an uncertainty estimation method where ρ between error and predictive uncertainty is highest. Spearman’s rank correlation coefficients were computed using the `scipy.stats` module [299].

Calibration analysis & miscalibration area In the regression case, the empirical probability of observing the true target values y_i around the predicted values \hat{y}_i should match the posterior predictive probability distribution $p(\hat{y}_i)$ defined by the uncertainty quantification method for a well-calibrated model [284]. That is, if we create 50% credible intervals around each predicted point value \hat{y}_i , the true value y_i should fall within that credible interval 50% of the time. To evaluate this, we assess uncertainty calibration following the procedure outlined in Tran et al. [260].

We assume the posterior predictive distributions to be Gaussian around the predicted point value \hat{y}_i . We find the lower and upper bound values between which we expect to observe a fraction e of the true values. For each data point x_i , given an inverse CDF function for the predictive distribution, F_i^{-1} , we define lower bound, Lb, and upper bound, Ub, values for each predicted point in the test set:

$$\text{Lb}_i = F_i^{-1}(0.5 - e)$$

$$\text{Ub}_i = F_i^{-1}(0.5 + e)$$

We count the fraction of true predictive values where $\text{Lb}_i < y_i < \text{Ub}_i$, which we denote as the “estimated confidence”. We repeat this value at various expected probability values, e , to create the calibration plots which show the expected proportion correct e against the estimated confidence, i.e. the estimated cumulative probability.

To quantify the degree of calibration, we measure each model’s deviation from ideal calibration by computing the area away from the parity line in which the estimated confidence and observed proportion correct are matched. This integration was computed using the `scipy.stats.simps` function [299].

For the lower-N datasets, the effect of the evidential regularizer strength was determined by varying the regularization parameter λ . Individual evidential D-MPNNs were trained with varying regularization coefficients. Lower λ leads to overconfident predictors, whereas higher λ leads to underconfident predictions, as they are penalized more for attributing higher evidence to erroneous predictions.

7.6.6 Active learning on QM9 dataset

Experiments were conducted on the QM9 dataset with the objective of improving the predictive accuracy of the model in terms of its learning efficiency. In other words, the objective of these experiments was to select a training dataset intelligently such that higher predictive accuracy could be achieved with less data. Indeed, because data in the chemical sciences

can often be limited and expensive to acquire, active learning can yield powerful predictive models that require minimal data generation.

The general experimental set up consists of iterative rounds of model training, data acquisition, and model assessment. All active learning experiments were conducted in the 2D setting using D-MPNN models. Briefly, models were initially trained on a randomly selected subset of data from QM9, and performance was assessed on a held-out test set and quantified via RMSE. At each iteration, m new data samples were selected and added to the training set on the basis of an acquisition function α . Models were then re-trained from scratch, and performance was again assessed on the held-out test set. This process was then repeated until the entire training dataset (consisting of 80% of all of QM9) was acquired. All experiments were conducted with $n = 10$ independent trials.

Two general acquisition strategies were tested in this study: first, an explorative strategy in which the algorithm chooses to acquire instances about which it is most uncertain, and second, a baseline strategy in which new data instances are acquired at random. These strategies are labeled as "Explorative" and "Random", respectively, in Figure 7.5. Evidential deep learning, model ensembling, and dropout sampling were used as the UQ methods for explorative selection. In all three instances, the value of the acquisition function for a particular point x is given by the estimated uncertainty: $\alpha(x) = \hat{\sigma}(x)$. For each method, the m samples with the greatest uncertainties were acquired at each iteration. Explorative acquisition was compared to each method’s respective random acquisition baseline, given differences in training and model output for each UQ approach.

7.6.7 Bayesian optimization on docking dataset

Bayesian optimization is an active learning approach that seeks to optimize an objective function by iteratively selecting experiments to perform according to a model’s predictions. In this work, Bayesian optimization was performed on a set of candidate molecules with the objective of selecting molecules that optimize a target property $f(x)$, in this case the ligand docking score against thymidylate kinase from AutoDock Vina. We follow the general procedure outlined by Graff et al. [282]. The objective function $f(x)$ is calculated for n randomly-selected molecules $\{x\}_{i=1}^n$, yielding a dataset $\mathcal{D} = \{(x_i, f(x_i))\}_{i=1}^n$. A D-MPNN model is then trained on these data, and predictions $\hat{f}(x)$ are passed to an acquisition function α which describes the utility of acquiring a new point. A set of m new points are subsequently selected according to the acquisition function; the objective function for each of these points is calculated; and these points are used to grow the dataset \mathcal{D} . This process was repeated iteratively for a fixed number of iterations. We refer the reader to recent works for more details on the Bayesian optimization procedure [282], [285].

The following acquisition functions were tested in this study:

$$\text{Random}(x) \sim \mathcal{U}(0, 1); \quad \text{Greedy}(x) = \hat{\mu}(x); \quad \text{UCB}(x) = \hat{\mu}(x) + \beta \hat{\sigma}(x).$$

Here $\hat{\mu}(x)$ and $\hat{\sigma}(x)$ are the model’s predicted mean and uncertainty at point x , respectively. $\beta = 2$ for all experiments reported in the paper.

Search trajectory performance was evaluated by the fraction of top- k scores identified during Bayesian optimization, calculated as the size of the intersection of the list of true

top- k scores and the list of top- k scores found, then divided by k . Acquisition sample diversity was measured as the average 10-nearest training set neighbors (10-NN) Tanimoto distance for batch samples after the first round of acquisition in Bayesian optimization.

7.6.8 Uncertainty-guided virtual screening for antibiotic discovery

Evidential D-MPNNs were trained to predict a molecule’s growth inhibitory effect on *E. coli*, following the general virtual screening pipeline presented by Stokes et al. [43]. Growth inhibitory activity was measured as *in vitro* OD₆₀₀ of *E. coli* following incubation with a compound, where lower values correspond to more potent inhibition, and to estimate the uncertainty associated with that prediction. The evidential D-MPNN model was trained following the procedure outlined by Stokes et al. [43], with the notable exception that the prediction task was formulated as a regression problem. Briefly, the molecular representation learned by the D-MPNN was augmented with 200 additional molecule-level features computed in RDKit [158]. Models were trained on the primary dataset of the OD₆₀₀ (target values y) of *E. coli* following incubation with each of 2,335 small molecules (input values x in SMILES representation) using a 80/20 training/validation split. Models were trained for 30 epochs with five-fold cross validation and a regularization coefficient $\lambda = 0.1$. Due to the imbalanced distribution of OD₆₀₀ values in this dataset (Figure 7.15A), the training dataset was rebalanced by sampling molecules with OD₆₀₀ > 0.2 with probability 0.1.

The trained evidential D-MPNN was applied to the Broad Drug Repurposing Hub [287] to predict OD₆₀₀ values and uncertainties for molecules in this discovery dataset. t-SNE analysis was conducted using scikit-learn’s implementation of t-Distributed Stochastic Neighbor Embedding. Morgan fingerprints for each molecule using a radius of 2 and 2048-bit fingerprint vectors were first computed in RDKit, and t-SNE with Tanimoto (Jaccard) distance metric and default parameters was then used to reduce the data from 2048 dimensions to two dimensions. The distance between points in the t-SNE plots thus reflects the Tanimoto distance of the corresponding molecules.

For the uncertainty-guided virtual screen, molecules from the Broad discovery dataset were first ranked based on predicted OD₆₀₀ values (lower is better), and the top 50 ranking molecules were downselected. Confidence percentiles across this set were computed based on uncertainty estimates returned by the evidential D-MPNN, and the set of 50 molecules was subsequently filtered according to varying confidence percentiles. Specifically, for a given confidence threshold p , molecules with estimated confidences below the associated p^{th} percentile are removed from the list of top 50 molecules, with p ranging from the 50th to 100th percentiles of greatest predictive confidence. The experimental hit rate for both the initial set of 50 molecules (i.e., no confidence filtering) and each filtered set was determined using empirically determined OD₆₀₀ values reported in Stokes et al. [43]. The hit rate was defined as the proportion of molecules in each candidate set with OD₆₀₀ < 0.2 (as previously reported [43]) relative to the total number of molecules in that candidate set. This screening procedure was also used for the dropout and ensemble-based baseline methods, following training and testing on the Stokes’ training and Broad discovery datasets, respectively.

7.7 Additional Results

| Cutoff | Delaney | | | Freesolv | | | Lipo | | | QM7 ($\times 10^2$) | | |
|--------|------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|-----------------------|------------------------|------------------------|
| | Dropout | Ensemble | Evidence | Dropout | Ensemble | Evidence | Dropout | Ensemble | Evidence | Dropout | Ensemble | Evidence |
| 0.0 | 0.68 \pm 0.02 | 0.65 \pm 0.03 | 0.66 \pm 0.02 | 1.00 \pm 0.06 | 0.94 \pm 0.06 | 0.96 \pm 0.07 | 0.55 \pm 0.01 | 0.53 \pm 0.02 | 0.55 \pm 0.02 | 1.18 \pm 0.02 | 1.12 \pm 0.02 | 1.15 \pm 0.03 |
| 0.5 | 0.62 \pm 0.03 | 0.55 \pm 0.03 | 0.44 \pm 0.01 | 0.79 \pm 0.07 | 0.45 \pm 0.04 | 0.42 \pm 0.04 | 0.52 \pm 0.01 | 0.40 \pm 0.01 | 0.50 \pm 0.01 | 0.88 \pm 0.06 | 0.88 \pm 0.06 | 0.39 \pm 0.03 |
| 0.75 | 0.59 \pm 0.03 | 0.50 \pm 0.05 | 0.35 \pm 0.02 | 0.85 \pm 0.12 | 0.41 \pm 0.05 | 0.36 \pm 0.04 | 0.50 \pm 0.02 | 0.38 \pm 0.02 | 0.51 \pm 0.02 | 0.65 \pm 0.03 | 0.81 \pm 0.06 | 0.23 \pm 0.04 |
| 0.90 | 0.55 \pm 0.03 | 0.51 \pm 0.09 | 0.28 \pm 0.02 | 0.66 \pm 0.20 | 0.40 \pm 0.06 | 0.35 \pm 0.08 | 0.46 \pm 0.03 | 0.38 \pm 0.02 | 0.53 \pm 0.03 | 0.69 \pm 0.05 | 0.71 \pm 0.11 | 0.10 \pm 0.04 |
| 0.95 | 0.53 \pm 0.06 | 0.45 \pm 0.06 | 0.22 \pm 0.02 | 0.75 \pm 0.30 | 0.27 \pm 0.04 | 0.38 \pm 0.12 | 0.49 \pm 0.04 | 0.36 \pm 0.03 | 0.50 \pm 0.04 | 0.73 \pm 0.08 | 0.69 \pm 0.11 | 0.10 \pm 0.04 |

| Cutoff | Enamine D-MPNN | | | QM9 D-MPNN | | | QM9 Atomistic ($\times 10^{-2}$) | | |
|--------|-----------------|------------------------|------------------------|-----------------|------------------------|------------------------|------------------------------------|------------------------|--|
| | Dropout | Ensemble | Evidence | Dropout | Ensemble | Evidence | Ensemble | Evidence | |
| 0.0 | 3.40 \pm 0.12 | 4.47 \pm 0.18 | 5.60 \pm 0.20 | 0.35 \pm 0.00 | 0.33 \pm 0.00 | 0.35 \pm 0.00 | 2.04 \pm 0.03 | 2.98 \pm 0.08 | |
| 0.5 | 3.64 \pm 0.05 | 2.12 \pm 0.02 | 1.55 \pm 0.12 | 0.33 \pm 0.00 | 0.32 \pm 0.00 | 0.30 \pm 0.00 | 1.45 \pm 0.02 | 1.52 \pm 0.02 | |
| 0.75 | 3.42 \pm 0.04 | 1.94 \pm 0.04 | 1.04 \pm 0.13 | 0.33 \pm 0.00 | 0.32 \pm 0.00 | 0.28 \pm 0.00 | 1.36 \pm 0.02 | 1.33 \pm 0.02 | |
| 0.90 | 3.30 \pm 0.06 | 1.80 \pm 0.03 | 0.63 \pm 0.12 | 0.32 \pm 0.00 | 0.32 \pm 0.01 | 0.27 \pm 0.00 | 1.31 \pm 0.03 | 1.18 \pm 0.03 | |
| 0.95 | 3.26 \pm 0.05 | 1.79 \pm 0.05 | 0.42 \pm 0.01 | 0.33 \pm 0.01 | 0.32 \pm 0.01 | 0.26 \pm 0.01 | 1.29 \pm 0.03 | 1.12 \pm 0.03 | |

| Cutoff | Clearance | | | LD50 | | | PPBR | | |
|--------|-------------------------|-------------------------|-------------------------|------------------------|-------------------------|-------------------------|-------------------------|-------------------------|-------------------------|
| | Dropout | Ensemble | Evidence | Dropout | Ensemble | Evidence | Dropout | Ensemble | Evidence |
| 0.0 | 47.45 \pm 0.90 | 44.05 \pm 0.68 | 5.60 \pm 0.20 | 0.58 \pm 0.01 | 0.56 \pm 0.05 | 0.56 \pm 0.01 | 11.67 \pm 0.42 | 11.17 \pm 0.40 | 11.35 \pm 0.25 |
| 0.5 | 33.20 \pm 2.12 | 38.81 \pm 1.60 | 37.68 \pm 2.35 | 0.50 \pm 0.01 | 00.47 \pm 0.01 | 0.49 \pm 0.01 | 7.57 \pm 0.68 | 6.42 \pm 0.62 | 6.83 \pm 0.59 |
| 0.75 | 30.33 \pm 2.94 | 37.23 \pm 2.94 | 32.02 \pm 3.60 | 0.48 \pm 0.011 | 0.44 \pm 0.01 | 0.45 \pm 0.02 | 6.47 \pm 0.87 | 4.90 \pm 0.71 | 4.90 \pm 0.86 |
| 0.90 | 27.71 \pm 3.62 | 38.87 \pm 5.50 | 29.75 \pm 5.64 | 0.47 \pm 0.02 | 0.46 \pm 0.018 | 0.44 \pm 0.03 | 5.99 \pm 1.12 | 4.97 \pm 0.96 | 2.88 \pm 0.87 |
| 0.95 | 27.94 \pm 6.13 | 34.18 \pm 5.96 | 25.84 \pm 7.29 | 0.47 \pm 0.03 | 0.47 \pm 0.04 | 0.45 \pm 0.038 | 5.26 \pm 1.23 | 4.59 \pm 1.24 | 0.91 \pm 0.16 |

Table 7.2: **Extended model error at various confidence percentile cutoffs including TDC lower-N data.** For a given confidence percentile cutoff, top performing methods based on prediction standard error of the mean (\pm s.e.m.) are bolded. A cutoff of 0.95 indicates that only the top 5% most confident predictions are considered. Full confidence plots for all datasets are shown in Figure 7.3 and Figures 7.7, 7.10, 7.11. Mean \pm s.e.m. (RMSE for all D-MPNN models, MAE for atomistic); $n = 10$ independent trials for lower-N datasets, $n = 5$ independent trials for higher-N datasets.

| Ligands Explored | Evidence vs. Dropout | | Ensemble vs. Evidence | | Ensemble vs. Dropout | |
|------------------|----------------------|----------|-----------------------|----------|----------------------|----------|
| | FC | p-value | FC | p-value | FC | p-value |
| 550 | 0.959 \pm 0.56 | 0.003669 | 0.926 \pm 0.33 | 0.195383 | 0.845 \pm 0.60 | 0.477471 |
| 775 | 1.052 \pm 0.17 | 0.259450 | 1.144 \pm 0.17 | 0.016443 | 1.196 \pm 0.12 | 0.001081 |
| 1000 | 1.039 \pm 0.11 | 0.242781 | 1.104 \pm 0.11 | 0.013200 | 1.142 \pm 0.09 | 0.000964 |
| 1225 | 1.057 \pm 0.05 | 0.009547 | 1.073 \pm 0.08 | 0.013934 | 1.131 \pm 0.06 | 0.000046 |
| 1450 | 1.055 \pm 0.05 | 0.007914 | 1.065 \pm 0.06 | 0.007536 | 1.122 \pm 0.05 | 0.000007 |
| 1675 | 1.056 \pm 0.05 | 0.006281 | 1.045 \pm 0.04 | 0.003669 | 1.104 \pm 0.04 | 0.000010 |

Table 7.3: **Statistical significance tests for Enamine 50k Bayesian optimization.** Pairwise comparison of uncertainty quantification methods for upper confidence bound (UCB) acquisition in Bayesian optimization on Enamine 50k data (Figure 7.5D). Fold changes (FC; mean \pm s.d.) reflect the fold change between the mean percentage of top-500 scores found for the first method listed relative to the second method listed. Significance values reflect the result of two-tailed unpaired t-tests over $n = 10$ independent trials.

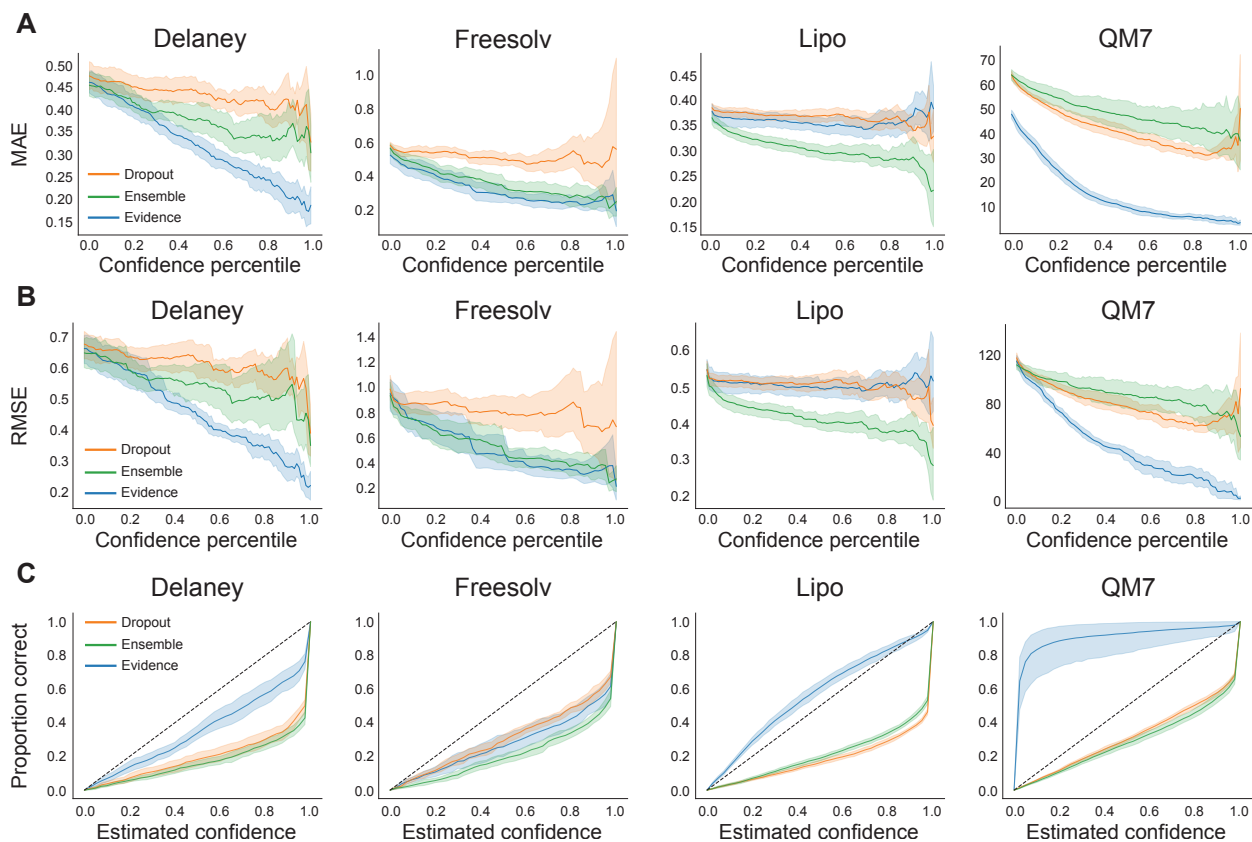


Figure 7.7: **Uncertainty benchmarking and calibration for lower-N datasets.** (A, B) Prediction error, measured as MAE (A) or RMSE (B), at different confidence percentile cutoffs for D-MPNNs evaluated on each of the 2D lower-N datasets. (C) Estimated confidence (cumulative probability) against the observed proportion correct for an evidential D-MPNN evaluated on each of the 2D lower-N datasets, with regularization parameter $\lambda = 0.2$. Mean \pm 95% c.i., $n = 10$ independent trials.

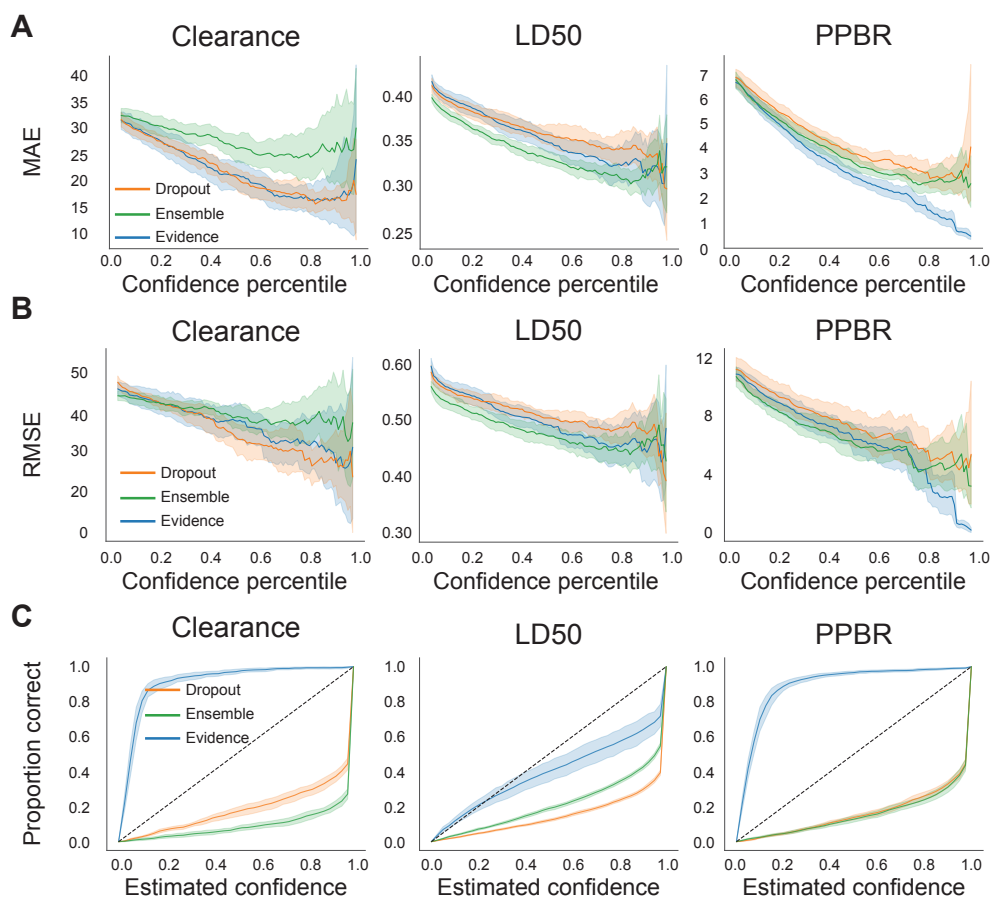


Figure 7.8: **Uncertainty benchmarking and calibration for additional lower-N Therapeutics Data Commons datasets.** (A, B) Prediction error, measured as MAE (A) or RMSE (B), at different confidence percentile cutoffs for D-MPNNs evaluated on each of the 2D lower-N datasets. (C) Estimated confidence (cumulative probability) against the observed proportion correct for an evidential D-MPNN evaluated on each of the 2D lower-N datasets, with regularization parameter $\lambda = 0.2$. Mean \pm 95% c.i., $n = 10$ independent trials.

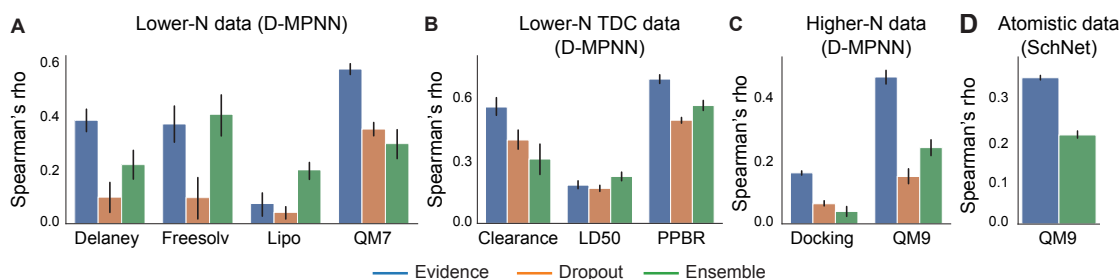


Figure 7.9: **Spearman rank correlation between error and uncertainty.** Spearman rank correlation coefficient between the estimated uncertainty and the absolute error for each point across lower-N (**A**), additional lower-N Therapeutics Data Commons (TDC) (**B**), higher-N (**C**), and atomistic (**D**) datasets. Mean \pm 95% c.i., $n = 10$ independent trials for lower-M and $n = 5$ independent trials for atomistic and higher-N datasets. For QM9 in the higher-N D-MPNN setting (C), standard error bars are computed across all tasks and independent trials ($n = 60$).

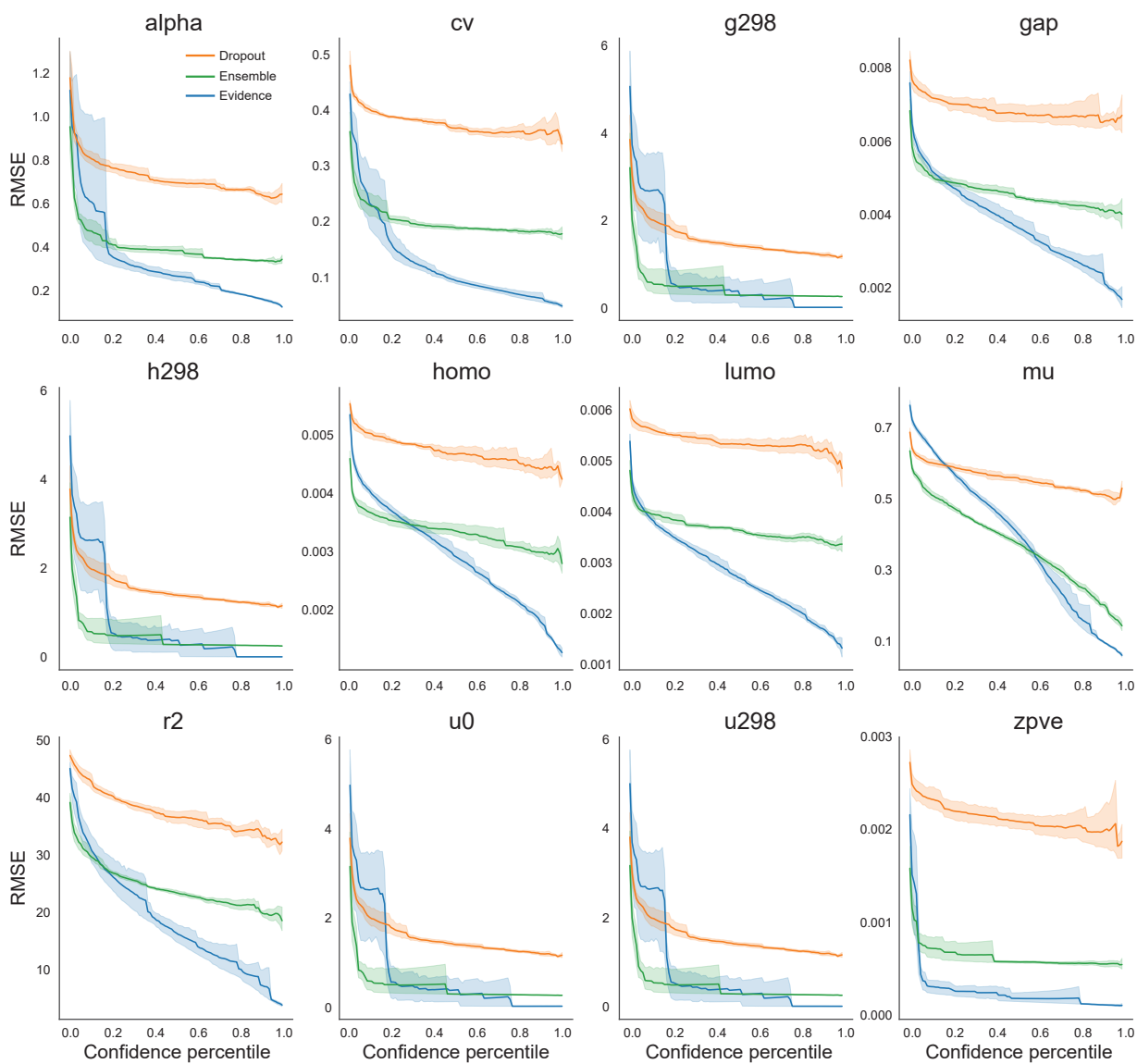


Figure 7.10: **Task-specific cutoffs for QM9 dataset.** RMSE at different confidence percentile cutoffs for D-MPNNs evaluated on each of the individual tasks from the QM9 dataset. Mean \pm 95% c.i., $n = 5$ independent trials.

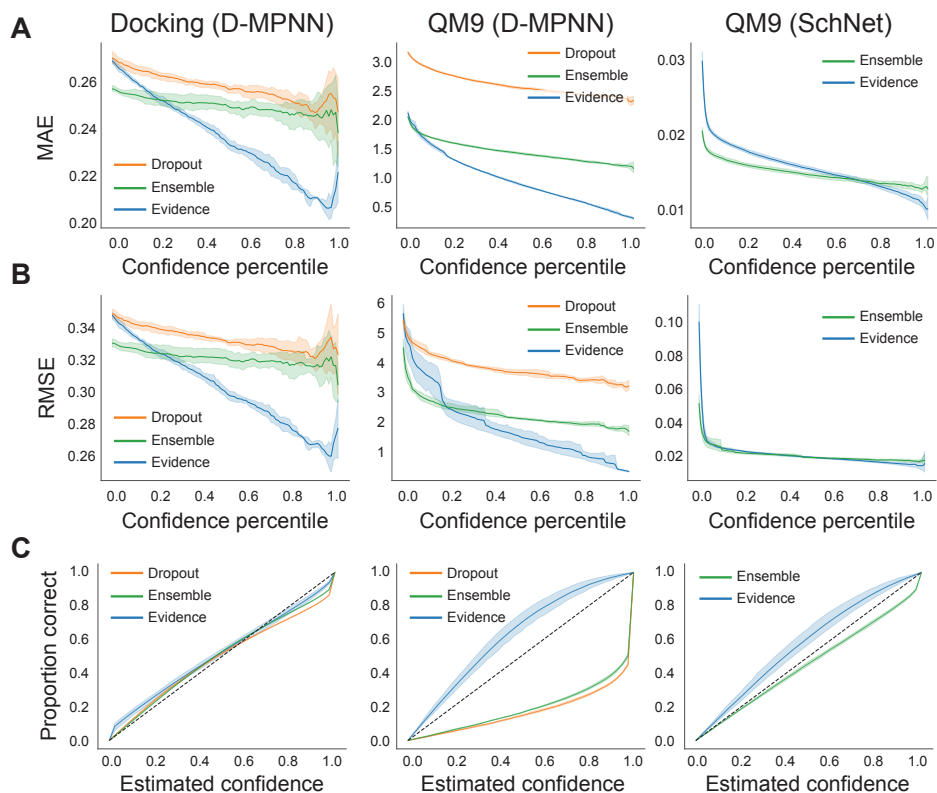


Figure 7.11: **Uncertainty benchmarking and calibration for higher-N 2D and 3D datasets.** (A, B) Prediction error, measured as MAE (A) or RMSE (B), at different confidence percentile cutoffs for models evaluated on the higher-N 2D and 3D datasets tested. (C) Estimated confidence (cumulative probability) against the observed proportion correct for the higher-N 2D and 3D datasets tested. Mean \pm 95% c.i., $n = 5$ independent trials.

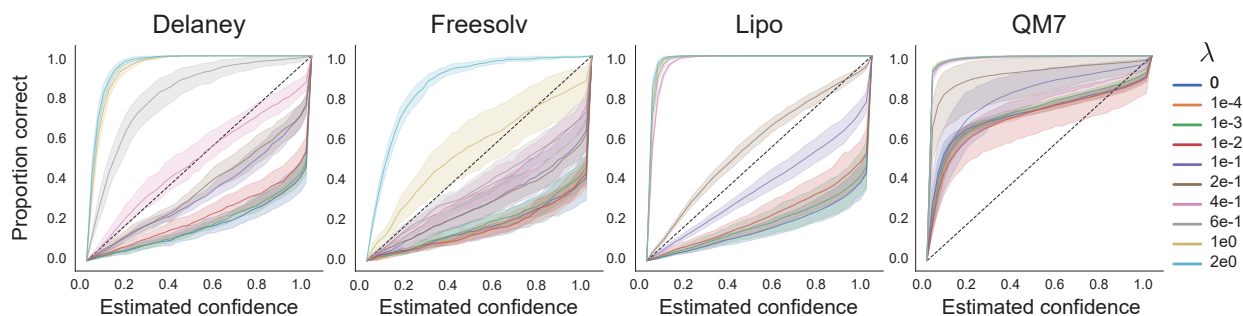


Figure 7.12: **Effect of λ on uncertainty calibration.** Evidential D-MPNNs are trained with different regularization coefficients λ on each of the lower-N datasets. Estimated confidence (cumulative probability) against the observed proportion correct is computed and plotted across different λ . Mean \pm 95% c.i., $n = 10$ independent trials.

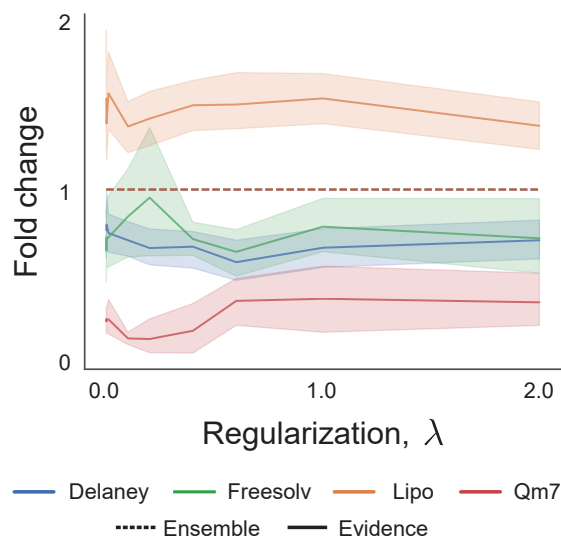


Figure 7.13: **Cutoff RMSE is robust to small λ shifts.** Evidential D-MPNNs are trained with different regularization coefficients λ on each of the lower-N datasets. RMSE is computed at the 10% confidence cutoff for each model. All computed errors are reported as the fold change of the top 10% RMSE for the evidential method relative to that of the ensemble baseline. Mean \pm 95% c.i., $n = 10$ independent trials.

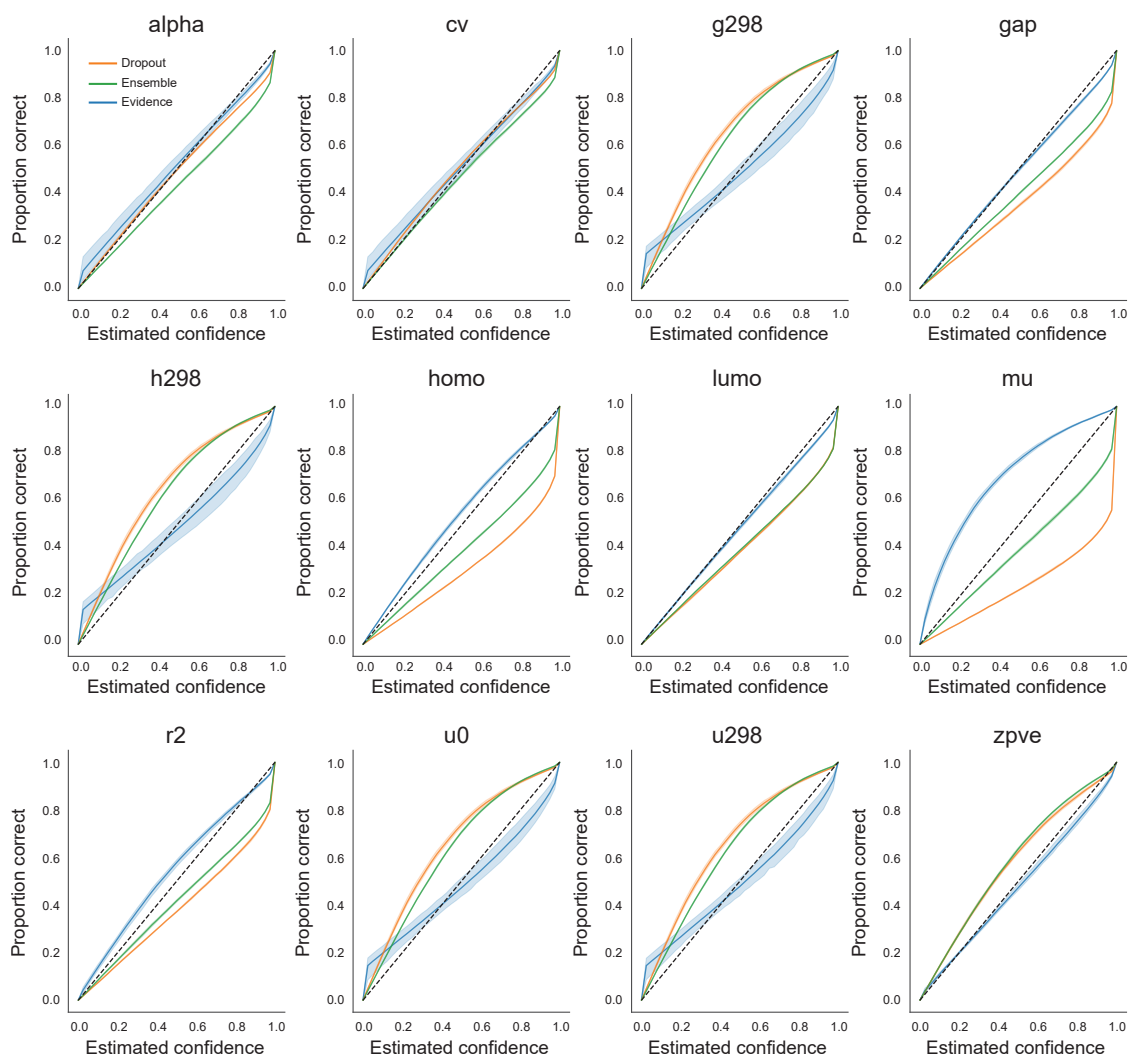


Figure 7.14: **Task-specific calibration for QM9 dataset.** Estimated confidence (cumulative probability) against the observed proportion correct is computed for an evidential D-MPNN evaluated on QM9 and then broken down into task-specific plots. Mean \pm 95% c.i., $n = 5$ independent trials.

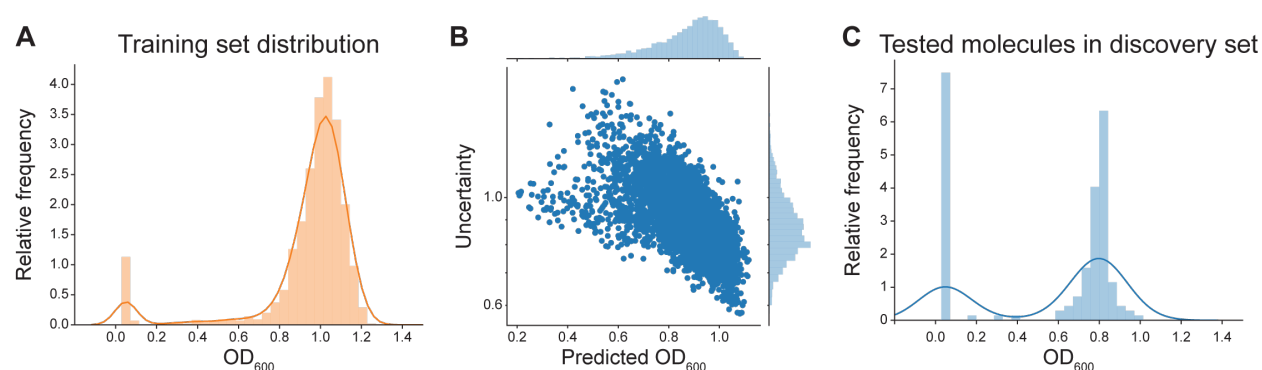


Figure 7.15: **Antibiotic discovery datasets and uncertainty predictions.** (A) Distribution of OD_{600} values for the training dataset of small molecules and their *in vitro* growth inhibitory activity against *E. coli*, as originally measured by Stokes et al [43]. Lower OD_{600} values indicate less *E. coli* growth and hence correspond to greater antibiotic activity. (B) Distribution of predicted OD_{600} values and evidential uncertainties for molecules in the Broad Drug Repurposing Hub discovery dataset. (C) Distribution of empirically determined OD_{600} for the subset of the discovery set (162 out of 6,111 total molecules) that was experimentally tested for *in vitro* growth inhibitory activity against *E. coli*.

Chapter 8

Conclusion

8.1 Contributions

This thesis introduces a set of practical and domain-tailored neural network models that can be used to discover new metabolites and interpret mass spectrometry data. Collectively, this is a step forward in an exciting journey to neuralize the small molecule annotation and discovery pipeline, a process that has been especially impactful in the analogous field of protein structure prediction.

Chapter 2 first focused on the task of predicting structural properties of target molecules, using no prior information of the molecule’s identity but its mass spectrum signature. This chapter introduced a formula-based representation of spectra as input to a Molecular Formula Transformer in a model termed MIST and showed how this could be utilized immediately to analyze mass spectrum data from IBD patient cohorts. All models and benchmarking comparisons are open-source to democratize progress. Chapter 3 extended this neural network architecture to neuralize an upstream component of the structure prediction pipeline, molecular formula prediction of the MS1 precursor mass, in a model MIST-CF.

The following two chapters inverted our modeling efforts to address the molecule-to-spectrum task using the same insights that drove improvements in performance within MIST. Chapter 4 presented SCARF, a model that utilizes a novel neural network prefix-tree decoding strategy to *generate* spectra at the level of molecular formulae, reasoning that formula-based representations had resulted in large performance boosts in the MIST and MIST-CF models.

Surprisingly, while SCARF led to high performance in spectrum prediction accuracy, the predictions were only modestly beneficial in the task of retrieval; that is, predicted spectra were not properly differentiated between highly similar isomer molecules. Chapter 5 improved upon this shortcoming by introducing ICEBERG, a model to predict spectra at the molecule-substructure representation level.

Chapter 6 and Chapter 7 described two orthogonal efforts related to downstream metabolite interactions and property prediction. Chapter 6 showcased how to predict compatibility between enzymes and substrates using learned representations and high throughput screening data. Chapter 7 incorporated evidential uncertainty into neural network predictions for small molecule properties.

8.2 Future Work

8.2.1 Building deployment-ready systems

The prospect of utilizing the models described in this thesis to effectively discover new biologically-derived metabolites and molecules from mass spectra is tantalizing. Further productizing the models and tasks developed by this thesis will be necessary to ensure prime time deployment and widespread use.

First, there is the practical component of packaging these models together into a system that is ready for consumption by end users with little technical computational background. I envision an effective system as accepting an input MGF file of unannotated tandem mass spectra, using MIST-CF to predict molecular formula for the precursor molecule, pipelining these output candidates into MIST to predict structural molecule fingerprints, querying the predicted fingerprints against a database of molecule candidates, and utilizing ICEBERG to re-rank the top 50 most promising candidates produced by MIST for increased accuracy. The leading academic tool that is most widely adopted may not be the most accurate, but will instead have the most optimal user experience.

Within this system, each individual component can be retrained with additional data and extended for improved performance across wider breadths of chemical space. For instance, all models can be extended and retrained to include negative mode ionization data and also data from several different adduct types, as users are increasingly interested in negative-mode acquisitions for chemical classes such as lipids. Such details were omitted in initial model development to focus purely on the task and method development, without conflating additional dependent variables related to data quality. Each individual model can also benefit from small adjustments accordingly: MIST-CF could be extended to utilize isotopic information to constrain the elements of plausible molecular formulae that are ranked, which is currently a user input feature; ICEBERG could be re-trained to make predictions at multiple collision energies (assuming this is provided with all training spectra); and MIST could consider mass spectrometer instrument types as additional input features and co-variates.

Beyond packaging the various models, there are several avenues for method development that could improve model accuracy. Most excitingly, I envision ways to combine the inverse spectrum-to-molecule and molecule-to-spectrum methods. I view the molecule-to-spectrum forward methods such as ICEBERG as more accurate but less efficient at exploring the full space of molecules present, as there are cost constraints at inference. I am optimistic that MIST can be retrained on predictions produced by ICEBERG to enable the accuracy of ICEBERG with the speed of MIST, a step I have already taken in subsequent unpublished, open-source work; MIST (and inverse models generally) can be a surrogate for high accuracy forward models that predict spectra from molecules. More general strategies and ideas for combining the forward and inverse directions of modeling into a single joint training process could be instrumental in unlocking new accuracy and even utilizing unlabeled spectra data, as has been done in unsupervised neural machine translation [300].

Orthogonal information and measurements can also be integrated into the models described to develop hybrid approaches. For instance, utilizing molecular networking data [69], NMR structural data, or even clues from paired genomic data [301] could all help constrain

the methods described above. Functional spectra readouts such as infrared spectra have already shown promise in improving the accuracy of molecule retrieval [302].

All the future research directions presented thus far concern themselves with improving models that *rank* candidate molecular structures. *De novo* molecular generation will be critical in going beyond the limits of current molecule retrieval databases and would complement improvements in candidate ranking or scoring. Many complex molecules that we detect analytically yet cannot annotate may be structurally novel (i.e., not found in large databases such as PubChem). Combining *de novo* molecule generation methods [180] with inverse and forward ranking methods described above could help overcome this limitation, with several methods already beginning these efforts [60], [61], [88], [303]. My belief throughout this thesis has been that ranking was more rate-limiting than candidate generation. While I still believe this, we may soon reach an inflection point where generative models of molecules are necessary to unlock higher structure elucidation accuracy.

8.2.2 Toward new discovery paradigms

Collectively, I am certain that these efforts will push the bounds of small molecule metabolite and natural product discovery to exciting new heights. This will be especially true with instrumentation improvements. Due to these multi-pronged advances, I expect metabolomics and lipidomics generally to emerge as more useful and widely adopted techniques.

These tools in turn will unlock new biomarkers, diagnostics, and even targets for intervening in human disease [67]. Toward this latter direction, improving the ease with which we can conduct quantitation, not just identification, of metabolites will be key and will provide exciting modeling avenues. One question, particularly in the case of therapeutics, is how to identify not only correlative biomarkers but rather causal metabolites. Follow up biochemical experiments to perturb or understand the mechanism by which a metabolite of interest acts can be challenging. It is not trivial to determine the pathways, enzymes, and intermediates that may have produced a target metabolite, nor the downstream receptors with which a metabolite may engage. Nevertheless, metabolite targets for therapeutic intervention are an exciting prospect that I expect to become more prevalent in the coming years.

More broadly than new biomarkers and targets, metabolomics and mass spectrometry have the potential to form the basis of novel molecular discovery platforms. The key challenge to deploying AI/ML in drug discovery is not new model development, as small molecule representation learning is a mature field, with many high quality models, such as graph neural networks, available to predict scalar properties from molecular structures [42]. Instead, the bottleneck is the dearth of high quality and streamlined data. Richard Sutton suggested with his so-called “bitter lesson” that scaling computation is a more promising direction than crafting methods with human intuition and biases [304]; similarly in the field of AI for biology, scaling data collection is more promising than crafting more clever algorithms and models. New open source data aggregation projects such as the Therapeutics Data Commons [183] have emerged to address this, alongside high throughput measurement platforms such as DNA-encoded libraries [305], [306] and phenotypic image-based screening [307], [308].

As a high throughput assay for measuring small molecules, mass spectrometry is well suited to this task of assaying small molecule functions in parallel in what is termed “functional metabolomics.” Already, companies are using mass spectrometry to comb through

natural products mined from the environments at record speeds, going beyond traditionally time intensive fractionation assays to effectively prioritize chemical matter. Chemoproteomics platforms utilize mass spectrometry to identify protein sites amenable to covalent ligands [309]. Next generation platforms such as equilibrium dialysis [310], affinity selection-mass spectrometry [311], and native metabolomics [312] all hold similar potential to unravel non-covalent protein-metabolite interactions. Analyzing the high dimensional, often noisy, and complex data from these various measurements in order to produce large datasets of biochemical data will be an important step toward AI enabled drug discovery.

Artificial intelligence and machine learning technologies are advancing at unprecedented speeds. While computational methods are not silver bullets, there has been no better time to leverage these tools for scientific discovery and life sciences innovation. It is tempting to pursue computational advances on already available data and well-established modeling tasks. Yet, I anticipate the most important breakthroughs will come when we exercise creativity to jointly invent new measurement platforms and modeling tasks, focusing on the exact data types and models most aligned with our ultimate scientific aims.

References

- [1] E. A. Underwood, “Boerhaave after three hundred years.,” *British Medical Journal*, vol. 4, no. 5634, p. 820, 1968.
- [2] J. Shendure and H. Ji, “Next-generation DNA sequencing,” *Nature Biotechnology*, vol. 26, no. 10, pp. 1135–1145, 2008.
- [3] F. Ozsolak and P. M. Milos, “RNA sequencing: advances, challenges and opportunities,” *Nature Reviews Genetics*, vol. 12, no. 2, pp. 87–98, 2011.
- [4] W. C. Cho, “Proteomics Technologies and Challenges,” *Genomics, Proteomics & Bioinformatics*, vol. 5, no. 2, pp. 77–85, 2007.
- [5] M. L. Metzker, “Sequencing technologies — the next generation,” *Nature Reviews Genetics*, vol. 11, no. 1, pp. 31–46, 2010.
- [6] D. Wishart, “Systems biology resources arising from the human metabolome project,” in *Genetics Meets Metabolomics: from experiment to systems biology*, Springer, 2012, pp. 157–175.
- [7] C. Jakobs, S. Schweitzer, and B. Dorland, “Galactitol in galactosemia,” *European Journal of Pediatrics*, vol. 154, S50–S52, 1995.
- [8] L. Dang, D. W. White, S. Gross, B. D. Bennett, M. A. Bittinger, E. M. Driggers, V. R. Fantin, H. G. Jang, S. Jin, M. C. Keenan, K. M. Marks, R. M. Prins, P. S. Ward, K. E. Yen, L. M. Liao, J. D. Rabinowitz, L. C. Cantley, C. B. Thompson, M. G. Vander Heiden, and S. M. Su, “Cancer-associated IDH1 mutations produce 2-hydroxyglutarate,” *Nature*, vol. 462, no. 7274, pp. 739–744, 2009.
- [9] W. Liu, Y. Wang, L. H. Bozi, P. D. Fischer, M. P. Jedrychowski, H. Xiao, T. Wu, N. Darabedian, X. He, E. L. Mills, *et al.*, “Lactate regulates cell cycle by remodelling the anaphase promoting complex,” *Nature*, vol. 616, no. 7958, pp. 790–797, 2023.
- [10] A. L. Demain and A. Fang, “The natural functions of secondary metabolites,” *History of Modern Biotechnology I*, pp. 1–39, 2000.
- [11] L. Yang, K.-S. Wen, X. Ruan, Y.-X. Zhao, F. Wei, and Q. Wang, “Response of plant secondary metabolites to environmental factors,” *Molecules*, vol. 23, no. 4, p. 762, 2018.
- [12] N. P. Keller, “Fungal secondary metabolism: regulation, function and drug discovery,” *Nature Reviews Microbiology*, vol. 17, no. 3, pp. 167–180, 2019.

- [13] R. A. Quinn, A. V. Melnik, A. Vrbanac, T. Fu, K. A. Patras, M. P. Christy, Z. Boudai, P. Belda-Ferre, A. Tripathi, L. K. Chung, M. Downes, R. D. Welch, M. Quinn, G. Humphrey, M. Panitchpakdi, K. C. Weldon, A. Aksenov, R. da Silva, J. Avila-Pacheco, C. Clish, S. Bae, H. Mallick, E. A. Franzosa, J. Lloyd-Price, R. Bussell, T. Thron, A. T. Nelson, M. Wang, E. Leszczynski, F. Vargas, J. M. Gauglitz, M. J. Meehan, E. Gentry, T. D. Arthur, A. C. Komor, O. Poulsen, B. S. Boland, J. T. Chang, W. J. Sandborn, M. Lim, N. Garg, J. C. Lumeng, R. J. Xavier, B. I. Kazmierczak, R. Jain, M. Egan, K. E. Rhee, D. Ferguson, M. Raffatellu, H. Vlamakis, G. G. Haddad, D. Siegel, C. Huttenhower, S. K. Mazmanian, R. M. Evans, V. Nizet, R. Knight, and P. C. Dorrestein, “Global chemical effects of the microbiome include new bile-acid conjugations,” *Nature*, vol. 579, no. 7797, pp. 123–129, 2020.
- [14] Y. Sato, K. Atarashi, D. R. Plichta, Y. Arai, S. Sasajima, S. M. Kearney, W. Suda, K. Takeshita, T. Sasaki, and S. Okamoto, “Novel bile acid biosynthetic pathways are enriched in the microbiome of centenarians,” *Nature*, vol. 599, no. 7885, pp. 458–464, 2021.
- [15] Y. Cao, J. Oh, M. Xue, W. J. Huh, J. Wang, J. A. Gonzalez-Hernandez, T. A. Rice, A. L. Martin, D. Song, J. M. Crawford, S. B. Herzon, and N. W. Palm, “Commensal microbiota from patients with inflammatory bowel disease produce genotoxic metabolites,” *Science*, vol. 378, no. 6618, eabm3233, 2022.
- [16] E. Y. Hsiao, S. W. McBride, S. Hsien, G. Sharon, E. R. Hyde, T. McCue, J. A. Codelli, J. Chow, S. E. Reisman, J. F. Petrosino, *et al.*, “The microbiota modulates gut physiology and behavioral abnormalities associated with autism,” *Cell*, vol. 155, no. 7, pp. 1451–1463, 2013.
- [17] A. L. Harvey, R. Edrada-Ebel, and R. J. Quinn, “The re-emergence of natural products for drug discovery in the genomics era,” *Nature Reviews Drug Discovery*, vol. 14, no. 2, pp. 111–129, 2015.
- [18] J. G. Bundy, M. P. Davey, and M. R. Viant, “Environmental metabolomics: a critical review and future perspectives,” *Metabolomics*, vol. 5, no. 1, pp. 3–21, 2009.
- [19] Z. Tian, H. Zhao, K. T. Peter, M. Gonzalez, J. Wetzell, C. Wu, X. Hu, J. Prat, E. Mudrock, R. Hettinger, A. E. Cortina, R. G. Biswas, F. V. C. Kock, R. Soong, A. Jenne, B. Du, F. Hou, H. He, R. Lundeen, A. Gilbreath, R. Sutton, N. L. Scholz, J. W. Davis, M. C. Dodd, A. Simpson, J. K. McIntyre, and E. P. Kolodziej, “A ubiquitous tire rubber-derived chemical induces acute mortality in coho salmon,” *Science*, vol. 371, no. 6525, pp. 185–189, 2021.
- [20] B. P. Bowen and T. R. Northen, “Dealing with the unknown: metabolomics and metabolite atlases,” *Journal of the American Society for Mass Spectrometry*, vol. 21, pp. 1471–1476, 2010.
- [21] G. J. Patti, O. Yanes, and G. Siuzdak, “Metabolomics: the apogee of the omics trilogy,” *Nature Reviews Molecular Cell Biology*, vol. 13, no. 4, pp. 263–269, 2012.
- [22] M. Beattie and O. A. Jones, “Rate of Advancement of Detection Limits in Mass Spectrometry: Is there a Moore’s Law of Mass Spec?” *Mass Spectrometry*, vol. 12, no. 1, A0118–A0118, 2023.

- [23] L. Perez de Souza, S. Alseekh, F. Scossa, and A. R. Fernie, “Ultra-high-performance liquid chromatography high-resolution mass spectrometry variants for metabolomics research,” *Nature Methods*, vol. 18, no. 7, pp. 733–746, 2021.
- [24] M. Giera, O. Yanes, and G. Siuzdak, “Metabolite discovery: Biochemistry’s scientific driver,” *Cell Metabolism*, vol. 34, no. 1, pp. 21–34, 2022.
- [25] S. Alseekh, A. Aharoni, Y. Brotman, K. Contrepois, J. D’Auria, J. Ewald, J. C. Ewald, P. D. Fraser, P. Giavalisco, R. D. Hall, M. Heinemann, H. Link, J. Luo, S. Neumann, J. Nielsen, L. Perez De Souza, K. Saito, U. Sauer, F. C. Schroeder, S. Schuster, G. Siuzdak, A. Skirycz, L. W. Sumner, M. P. Snyder, H. Tang, T. Tohge, Y. Wang, W. Wen, S. Wu, G. Xu, N. Zamboni, and A. R. Fernie, “Mass spectrometry-based metabolomics: a guide for annotation, quantification and best reporting practices,” *Nature Methods*, vol. 18, no. 7, pp. 747–756, 2021.
- [26] D. S. Wishart, Y. D. Feunang, A. Marcu, A. C. Guo, K. Liang, R. Vázquez-Fresno, T. Sajed, D. Johnson, C. Li, N. Karu, Z. Sayeeda, E. Lo, N. Assempour, M. Berjanskii, S. Singhal, D. Arndt, Y. Liang, H. Badran, J. Grant, A. Serra-Cayuela, Y. Liu, R. Mandal, V. Neveu, A. Pon, C. Knox, M. Wilson, C. Manach, and A. Scalbert, “HMDB 4.0: the human metabolome database for 2018,” *Nucleic Acids Research*, vol. 46, no. D1, pp. D608–D617, 2018.
- [27] S. Alseekh and A. R. Fernie, “Metabolomics 20 years on: what have we learned and what hurdles remain?” *The Plant Journal*, vol. 94, no. 6, pp. 933–942, 2018.
- [28] M. Wang, J. J. Carver, V. V. Phelan, L. M. Sanchez, N. Garg, Y. Peng, D. D. Nguyen, J. Watrous, C. A. Kapon, and T. Luzzatto-Knaan, “Sharing and community curation of mass spectrometry data with Global Natural Products Social Molecular Networking,” *Nature Biotechnology*, vol. 34, no. 8, pp. 828–837, 2016.
- [29] W. Bittremieux, M. Wang, and P. C. Dorrestein, “The critical role that spectral libraries play in capturing the metabolomics community knowledge,” *Metabolomics*, vol. 18, no. 12, p. 94, 2022.
- [30] K. He, X. Zhang, S. Ren, and J. Sun, “Deep Residual Learning for Image Recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [31] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “ImageNet Classification with Deep Convolutional Neural Networks,” in *Advances in Neural Information Processing Systems 25*, 2012.
- [32] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding,” *arXiv preprint arXiv:1810.04805*, 2018.
- [33] D. Bahdanau, K. Cho, and Y. Bengio, “Neural Machine Translation by Jointly Learning to Align and Translate,” in *International Conference on Learning Representations*, 2015.
- [34] L. Jiao and J. Zhao, “A survey on the new generation of deep learning in image processing,” *IEEE Access*, vol. 7, pp. 172 231–172 263, 2019.

- [35] J. Jumper, R. Evans, A. Pritzel, T. Green, M. Figurnov, O. Ronneberger, K. Tunyasuvunakool, R. Bates, A. Žídek, A. Potapenko, *et al.*, “Highly accurate protein structure prediction with AlphaFold,” *Nature*, vol. 596, no. 7873, pp. 583–589, 2021.
- [36] Z. Lin, H. Akin, R. Rao, B. Hie, Z. Zhu, W. Lu, N. Smetanin, R. Verkuil, O. Kabeli, Y. Shmueli, *et al.*, “Evolutionary-scale prediction of atomic-level protein structure with a language model,” *Science*, vol. 379, no. 6637, pp. 1123–1130, 2023.
- [37] M. Baek, F. DiMaio, I. Anishchenko, J. Dauparas, S. Ovchinnikov, G. R. Lee, J. Wang, Q. Cong, L. N. Kinch, R. D. Schaeffer, *et al.*, “Accurate prediction of protein structures and interactions using a three-track neural network,” *Science*, vol. 373, no. 6557, pp. 871–876, 2021.
- [38] R. Bommasani, D. A. Hudson, E. Adeli, R. Altman, S. Arora, S. von Arx, M. S. Bernstein, J. Bohg, A. Bosselut, E. Brunskill, *et al.*, “On the opportunities and risks of foundation models,” *arXiv preprint arXiv:2108.07258*, 2021.
- [39] H. Wang, T. Fu, Y. Du, W. Gao, K. Huang, Z. Liu, P. Chandak, S. Liu, P. Van Katwyk, A. Deac, *et al.*, “Scientific discovery in the age of artificial intelligence,” *Nature*, vol. 620, no. 7972, pp. 47–60, 2023.
- [40] M. H. Segler, M. Preuss, and M. P. Waller, “Planning chemical syntheses with deep neural networks and symbolic AI,” *Nature*, vol. 555, no. 7698, pp. 604–610, 2018.
- [41] C. W. Coley, D. A. Thomas III, J. A. Lummiss, J. N. Jaworski, C. P. Breen, V. Schultz, T. Hart, J. S. Fishman, L. Rogers, H. Gao, *et al.*, “A robotic platform for flow synthesis of organic compounds informed by AI planning,” *Science*, vol. 365, no. 6453, eaax1566, 2019.
- [42] K. Yang, K. Swanson, W. Jin, C. Coley, P. Eiden, H. Gao, A. Guzman-Perez, T. Hopper, B. Kelley, M. Mathea, *et al.*, “Analyzing Learned Molecular Representations for Property Prediction,” *Journal of Chemical Information and Modeling*, vol. 59, no. 8, pp. 3370–3388, 2019.
- [43] J. M. Stokes, K. Yang, K. Swanson, W. Jin, A. Cubillos-Ruiz, N. M. Donghia, C. R. MacNair, S. French, L. A. Carfrae, Z. Bloom-Ackerman, *et al.*, “A Deep Learning Approach to Antibiotic Discovery,” *Cell*, vol. 180, no. 4, pp. 688–702, 2020.
- [44] D. H. Bryant, A. Bashir, S. Sinai, N. K. Jain, P. J. Ogden, P. F. Riley, G. M. Church, L. J. Colwell, and E. D. Kelsic, “Deep diversification of an AAV capsid protein by machine learning,” *Nature Biotechnology*, vol. 39, no. 6, pp. 691–696, 2021.
- [45] A. Yala, P. G. Mikhael, F. Strand, G. Lin, K. Smith, Y.-L. Wan, L. Lamb, K. Hughes, C. Lehman, and R. Barzilay, “Toward robust mammography-based models for breast cancer risk,” *Science Translational Medicine*, vol. 13, no. 578, eaba4373, 2021.
- [46] S. Ovchinnikov and P.-S. Huang, “Structure-based protein design with deep learning,” *Current Opinion in Chemical Biology*, vol. 65, pp. 136–144, 2021.
- [47] M. AlQuraishi, “Machine learning in protein structure prediction,” *Current Opinion in Chemical Biology*, vol. 65, pp. 1–8, 2021.

- [48] A. Kryshchak, T. Schwede, M. Topf, K. Fidelis, and J. Moult, “Critical assessment of methods of protein structure prediction (CASP)—Round XIV,” *Proteins: Structure, Function, and Bioinformatics*, vol. 89, no. 12, pp. 1607–1617, 2021.
- [49] S. K. Burley, H. M. Berman, G. J. Kleywegt, J. L. Markley, H. Nakamura, and S. Velankar, “Protein Data Bank (PDB): the single global macromolecular structure archive,” *Protein Crystallography*, pp. 627–641, 2017.
- [50] K. Dührkop, M. Fleischauer, M. Ludwig, A. A. Aksenov, A. V. Melnik, M. Meusel, P. C. Dorrestein, J. Rousu, and S. Böcker, “SIRIUS 4: a rapid tool for turning tandem mass spectra into metabolite structure information,” *Nature Methods*, vol. 16, no. 4, pp. 299–302, 2019.
- [51] F. Wang, J. Liigand, S. Tian, D. Arndt, R. Greiner, and D. S. Wishart, “CFM-ID 4.0: More Accurate ESI-MS/MS Spectral Prediction and Compound Identification,” *Analytical Chemistry*, vol. 93, no. 34, pp. 11 692–11 700, 2021.
- [52] S. Xing, S. Shen, B. Xu, X. Li, and T. Huan, “BUDDY: molecular formula discovery via bottom-up MS/MS interrogation,” *Nature Methods*, vol. 20, pp. 881–890, 2023.
- [53] K. Dührkop, H. Shen, M. Meusel, J. Rousu, and S. Böcker, “Searching molecular structure databases with tandem mass spectra using CSI:FingerID,” *Proceedings of the National Academy of Sciences*, vol. 112, no. 41, pp. 12 580–12 585, 2015.
- [54] S. Böcker, M. C. Letzel, Z. Lipták, and A. Pervukhin, “SIRIUS: decomposing isotope patterns for metabolite identification,” *Bioinformatics*, vol. 25, no. 2, pp. 218–224, 2009.
- [55] G. Voronov, A. Frandsen, B. Bargh, D. Healey, R. Lightheart, T. Kind, P. Dorrestein, V. Colluru, and T. Butler, “MS2Prop: A machine learning model that directly predicts chemical properties from mass spectrometry data for novel compounds,” *BioRxiv*, 2022.
- [56] Z. Fan, A. Alley, K. Ghaffari, and H. W. Ransom, “MetFID: artificial neural network-based compound fingerprint prediction for metabolite annotation,” *Metabolomics*, vol. 16, no. 10, p. 104, 2020.
- [57] K. Dührkop, L.-F. Nothias, M. Fleischauer, R. Reher, M. Ludwig, M. A. Hoffmann, D. Petras, W. H. Gerwick, J. Rousu, and P. C. Dorrestein, “Systematic classification of unknown metabolites using high-resolution fragmentation mass spectra,” *Nature Biotechnology*, vol. 39, no. 4, pp. 462–471, 2021.
- [58] S. Kim, J. Chen, T. Cheng, A. Gindulyte, J. He, S. He, Q. Li, B. A. Shoemaker, P. A. Thiessen, and B. Yu, “PubChem 2019 update: improved access to chemical data,” *Nucleic Acids Research*, vol. 47, no. D1, pp. D1102–D1109, 2019.
- [59] D. S. Wishart, A. Guo, E. Oler, F. Wang, A. Anjum, H. Peters, R. Dizon, Z. Sayeeda, S. Tian, B. L. Lee, *et al.*, “HMDB 5.0: the Human Metabolome Database for 2022,” *Nucleic Acids Research*, vol. 50, no. D1, pp. D622–D631, 2022.
- [60] T. Butler, A. Frandsen, R. Lightheart, B. Bargh, J. Taylor, T. Bollerman, T. Kerby, K. West, G. Voronov, K. Moon, *et al.*, “MS2Mol: A transformer model for illuminating dark chemical space from mass spectra,” *ChemRxiv*, 2023.

- [61] M. A. Stravs, K. Dührkop, S. Böcker, and N. Zamboni, “MSNovelist: De novo structure generation from mass spectra,” *Nature Methods*, vol. 19, no. 7, pp. 865–870, 2022.
- [62] W. Gao, T. Fu, J. Sun, and C. Coley, “Sample Efficiency Matters: A Benchmark for Practical Molecular Optimization,” in *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2022.
- [63] NIST, *Tandem Mass Spectral Library*, 2020.
- [64] S. Goldman, J. Wohlwend, M. Stražar, G. Haroush, R. J. Xavier, and C. W. Coley, “Annotating metabolite mass spectra with domain-inspired chemical formula transformers,” *Nature Machine Intelligence*, vol. 5, no. 9, pp. 965–979, 2023.
- [65] W. Xu, H. Yang, Y. Liu, Y. Yang, P. Wang, S.-H. Kim, S. Ito, C. Yang, P. Wang, and M.-T. Xiao, “Oncometabolite 2-hydroxyglutarate is a competitive inhibitor of α -ketoglutarate-dependent dioxygenases,” *Cancer Cell*, vol. 19, no. 1, pp. 17–30, 2011.
- [66] M. P. Torrens-Spence, A. Bobokalonova, V. Carballo, C. M. Glinkerman, T. Pluskal, A. Shen, and J.-K. Weng, “PBS3 and EPS1 complete salicylic acid biosynthesis from isochorismate in *Arabidopsis*,” *Molecular Plant*, vol. 12, no. 12, pp. 1577–1586, 2019.
- [67] D. S. Wishart, “Metabolomics for Investigating Physiological and Pathophysiological Processes,” *Physiological Reviews*, vol. 99, no. 4, pp. 1819–1875, 2019.
- [68] S. Neumann and S. Böcker, “Computational mass spectrometry for metabolomics: identification of metabolites and small molecules,” *Analytical and Bioanalytical Chemistry*, vol. 398, no. 7, pp. 2779–2788, 2010.
- [69] L.-F. Nothias, D. Petras, R. Schmid, K. Dührkop, J. Rainer, A. Sarvepalli, I. Protsyuk, M. Ernst, H. Tsugawa, M. Fleischauer, F. Aicheler, A. A. Aksenov, O. Alka, P.-M. Allard, A. Barsch, X. Cachet, A. M. Caraballo-Rodriguez, R. R. Da Silva, T. Dang, N. Garg, J. M. Gauglitz, A. Gurevich, G. Isaac, A. K. Jarmusch, Z. Kameník, K. B. Kang, N. Kessler, I. Koester, A. Korf, A. Le Gouellec, M. Ludwig, C. Martin H., L.-I. McCall, J. McSayles, S. W. Meyer, H. Mohimani, M. Morsy, O. Moyne, S. Neumann, H. Neuweger, N. H. Nguyen, M. Nothias-Esposito, J. Paolini, V. V. Phelan, T. Pluskal, R. A. Quinn, S. Rogers, B. Shrestha, A. Tripathi, J. J. J. van der Hooft, F. Vargas, K. C. Weldon, M. Witting, H. Yang, Z. Zhang, F. Zubeil, O. Kohlbacher, S. Böcker, T. Alexandrov, N. Bandeira, M. Wang, and P. C. Dorrestein, “Feature-based molecular networking in the GNPS analysis environment,” *Nature Methods*, vol. 17, no. 9, pp. 905–908, 2020.
- [70] J. Watrous, P. Roach, T. Alexandrov, B. S. Heath, J. Y. Yang, R. D. Kersten, M. van der Voort, K. Pogliano, H. Gross, and J. M. Raaijmakers, “Mass spectral molecular networking of living microbial colonies,” *Proceedings of the National Academy of Sciences*, vol. 109, no. 26, E1743–E1752, 2012.
- [71] S. Wolf, S. Schmidt, M. Müller-Hannemann, and S. Neumann, “In silico fragmentation for computer assisted identification of metabolite mass spectra,” *BMC Bioinformatics*, vol. 11, no. 1, p. 148, 2010.

- [72] C. Ruttkies, E. L. Schymanski, S. Wolf, J. Hollender, and S. Neumann, “MetFrag re-launched: incorporating strategies beyond in silico fragmentation,” *Journal of Cheminformatics*, vol. 8, no. 1, pp. 1–16, 2016.
- [73] F. Allen, R. Greiner, and D. Wishart, “Competitive fragmentation modeling of ESI-MS/MS spectra for putative metabolite identification,” *Metabolomics*, vol. 11, no. 1, pp. 98–110, 2015.
- [74] H. Shen, N. Zamboni, M. Heinonen, and J. Rousu, “Metabolite identification through machine learning—tackling CASMI challenge using FingerID,” *Metabolites*, vol. 3, no. 2, pp. 484–505, 2013.
- [75] CASMI, *Critical Assessment of Small Molecule Identification*.
- [76] E. L. Schymanski, C. Ruttkies, M. Krauss, C. Brouard, T. Kind, K. Dührkop, F. Allen, A. Vaniya, D. Verdegem, and S. Böcker, “Critical assessment of small molecule identification 2016: automated methods,” *Journal of Cheminformatics*, vol. 9, no. 1, pp. 1–21, 2017.
- [77] S. Böcker and K. Dührkop, “Fragmentation trees reloaded,” *Journal of Cheminformatics*, vol. 8, no. 1, pp. 1–26, 2016.
- [78] G. Hjörleifsson Eldjárn, A. Ramsay, J. J. Van Der Hooft, K. R. Duncan, S. Soldatou, J. Rousu, R. Daly, J. Wandy, and S. Rogers, “Ranking microbial metabolomic and genomic links in the NPLinker framework using complementary scoring functions,” *PLoS Computational Biology*, vol. 17, no. 5, e1008920, 2021.
- [79] M. A. Hoffmann, L.-F. Nothias, M. Ludwig, M. Fleischauer, E. C. Gentry, M. Witting, P. C. Dorrestein, K. Dührkop, and S. Böcker, “High-confidence structural annotation of metabolites absent from spectral libraries,” *Nature Biotechnology*, vol. 40, no. 3, pp. 411–421, 2021.
- [80] A. Tripathi, Y. Vázquez-Baeza, J. M. Gauglitz, M. Wang, K. Dührkop, M. Nothias-Esposito, D. D. Acharya, M. Ernst, J. J. van der Hooft, and Q. Zhu, “Chemically informed analyses of metabolomics mass spectrometry data with Qemistree,” *Nature Chemical Biology*, vol. 17, no. 2, pp. 146–151, 2021.
- [81] F. Huber, L. Ridder, S. Verhoeven, J. H. Spaaks, F. Diblen, S. Rogers, and J. J. J. v. d. Hooft, “Spec2Vec: Improved mass spectral similarity scoring through learning of structural relationships,” *PLoS Computational Biology*, vol. 17, no. 2, e1008724, 2021.
- [82] F. Huber, S. van der Burg, J. J. van der Hooft, and L. Ridder, “MS2DeepScore: a novel deep learning similarity measure to compare tandem mass spectra,” *Journal of Cheminformatics*, vol. 13, no. 1, pp. 1–14, 2021.
- [83] G. Voronov, R. Lighthead, J. Davison, C. A. Krettler, D. Healey, and T. Butler, “Multi-scale Sinusoidal Embeddings Enable Learning on High Resolution Mass Spectrometry Data,” *arXiv preprint arXiv:2207.02980*, 2022.
- [84] J. N. Wei, D. Belanger, R. P. Adams, and D. Sculley, “Rapid Prediction of Electron-Ionization Mass Spectrometry Using Neural Networks,” *ACS Central Science*, vol. 5, no. 4, pp. 700–708, 2019.

- [85] X. Li, H. Zhu, L.-p. Liu, and S. Hassoun, "Ensemble Spectral Prediction (ESP) Model for Metabolite Annotation," *arXiv preprint arXiv:2203.13783*, 2022.
- [86] A. Young, B. Wang, and H. Röst, "MassFormer: Tandem Mass Spectrum Prediction with Graph Transformers," *arXiv preprint arXiv:2111.04824*, 2021.
- [87] A. D. Shrivastava, N. Swainston, S. Samanta, I. Roberts, M. Wright Muelas, and D. B. Kell, "MassGenie: A transformer-based deep learning method for identifying small molecules from their mass spectra," *Biomolecules*, vol. 11, no. 12, p. 1793, 2021.
- [88] E. E. Litsa, V. Chenthamarakshan, P. Das, and L. E. Kaviraki, "An end-to-end deep learning framework for translating mass spectra to de-novo molecules," *Communications Chemistry*, vol. 6, no. 1, p. 132, 2023.
- [89] D. Weininger, "SMILES, a chemical language and information system. 1. Introduction to methodology and encoding rules," *Journal of Chemical Information and Computer Sciences*, vol. 28, no. 1, pp. 31–36, 1988.
- [90] K. Dührkop, "Deep kernel learning improves molecular fingerprint prediction from tandem mass spectra," *Bioinformatics*, vol. 38, no. Supplement_1, pp. i342–i349, 2022.
- [91] J. Lee, Y. Lee, J. Kim, A. Kosior, S. Choi, and Y. W. Teh, "Set Transformer: A Framework for Attention-based Permutation-Invariant Neural Networks," in *Proceedings of the 36th International Conference on Machine Learning*, 2019, pp. 3744–3753.
- [92] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention Is All You Need," in *Advances in Neural Information Processing Systems 30*, 2017, pp. 5998–6008.
- [93] A. Aisporna, H. P. Benton, A. Chen, R. J. Derks, J. M. Galano, M. Giera, and G. Siuzdak, "Neutral loss mass spectral data enhances molecular similarity analysis in METLIN," *Journal of the American Society for Mass Spectrometry*, vol. 33, no. 3, pp. 530–534, 2022.
- [94] T. Karras, T. Aila, S. Laine, and J. Lehtinen, "Progressive Growing of GANs for Improved Quality, Stability, and Variation," in *International Conference on Learning Representations*, 2018.
- [95] L. Ridder, J. J. van der Hooft, and S. Verhoeven, "Automatic compound annotation from mass spectrometry data using MAGMa," *Mass Spectrometry*, vol. 3, no. Spec Iss 2, S0033, 2014.
- [96] Q. Xie, M.-T. Luong, E. Hovy, and Q. V. Le, "Self-Training With Noisy Student Improves ImageNet Classification," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 10 687–10 698.
- [97] *MassBank of North America*.
- [98] M. Ludwig, K. Dührkop, and S. Böcker, "Bayesian networks for mass spectrometric metabolite identification via molecular fingerprints," *Bioinformatics*, vol. 34, no. 13, pp. i333–i340, 2018.

- [99] A. v. d. Oord, Y. Li, and O. Vinyals, “Representation Learning with Contrastive Predictive Coding,” *arXiv preprint arXiv:1807.03748*, 2018.
- [100] F. Huber, S. Verhoeven, C. Meijer, H. Spreeuw, E. M. V. Castilla, C. Geng, J. J. j. van der Hooft, S. Rogers, A. Belloum, F. Diblen, and J. H. Spaaks, “matchms - processing and similarity evaluation of mass spectrometry data,” *Journal of Open Source Software*, vol. 5, p. 2411, 2020.
- [101] H. W. Kim, M. Wang, C. A. Leber, L.-F. Nothias, R. Reher, K. B. Kang, J. J. Van Der Hooft, P. C. Dorrestein, W. H. Gerwick, and G. W. Cottrell, “NPCClassifier: a deep neural network-based structural classification tool for natural products,” *Journal of Natural Products*, vol. 84, no. 11, pp. 2795–2807, 2021.
- [102] R. H. Mills, P. S. Dulai, Y. Vázquez-Baeza, C. Saucedo, N. Daniel, R. R. Gerner, L. E. Batachari, M. Malfavon, Q. Zhu, and K. Weldon, “Multi-omics analyses of the ulcerative colitis gut microbiome link *Bacteroides vulgatus* proteases with disease severity,” *Nature Microbiology*, vol. 7, no. 2, pp. 262–276, 2022.
- [103] E. A. Franzosa, A. Sirota-Madi, J. Avila-Pacheco, N. Fornelos, H. J. Haiser, S. Reinker, T. Vatanen, A. B. Hall, H. Mallick, and L. J. McIver, “Gut microbiome structure and metabolic activity in inflammatory bowel disease,” *Nature Microbiology*, vol. 4, no. 2, pp. 293–305, 2019.
- [104] M. Schirmer, L. Denson, H. Vlamakis, E. A. Franzosa, S. Thomas, N. M. Gotman, P. Rufo, S. S. Baker, C. Sauer, and J. Markowitz, “Compositional and temporal changes in the gut microbiome of pediatric ulcerative colitis patients are linked to disease course,” *Cell Host & Microbe*, vol. 24, no. 4, pp. 600–610. e4, 2018.
- [105] D. F. Rojas-Tapias, E. M. Brown, E. R. Temple, M. A. Onyekaba, A. M. T. Mohamed, K. Duncan, M. Schirmer, R. L. Walker, T. Mayassi, K. A. Pierce, J. Ávila-Pacheco, C. B. Clish, H. Vlamakis, and R. J. Xavier, “Inflammation-associated nitrate facilitates ectopic colonization of oral bacterium *Veillonella parvula* in the intestine,” *Nature Microbiology*, vol. 7, no. 10, pp. 1673–1685, 2022.
- [106] G. A. Bezerra, Y. Ohara-Nemoto, I. Cornaciu, S. Fedosyuk, G. Hoffmann, A. Round, J. A. Márquez, T. K. Nemoto, and K. Djinić-Carugo, “Bacterial protease uses distinct thermodynamic signatures for substrate recognition,” *Scientific Reports*, vol. 7, no. 1, p. 2848, 2017.
- [107] M. Wlodarska, C. Luo, R. Kolde, E. d’Hennezel, J. W. Annand, C. E. Heim, P. Krastel, E. K. Schmitt, A. S. Omar, E. A. Creasey, A. L. Garner, S. Mohammadi, D. J. O’Connell, S. Abubucker, T. D. Arthur, E. A. Franzosa, C. Huttenhower, L. O. Murphy, H. J. Haiser, H. Vlamakis, J. A. Porter, and R. J. Xavier, “Indoleacrylic Acid Produced by Commensal *Peptostreptococcus* Species Suppresses Inflammation,” *Cell Host & Microbe*, vol. 22, no. 1, pp. 25–37.e6, 2017.
- [108] E. L. Schymanski and S. Neumann, “The critical assessment of small molecule identification (CASMI): challenges and solutions,” *Metabolites*, vol. 3, no. 3, pp. 517–538, 2013.
- [109] G. Landrum, “RDKit: A software suite for cheminformatics, computational chemistry, and predictive modeling,” 2013.

- [110] T. Malisiewicz, A. Gupta, and A. A. Efros, “Ensemble of exemplar-svms for object detection and beyond,” in *International Conference on Computer Vision*, IEEE, 2011, pp. 89–96.
- [111] M. Ludwig, L.-F. Nothias, K. Dührkop, I. Koester, M. Fleischauer, M. A. Hoffmann, D. Petras, F. Vargas, M. Morsy, and L. Aluwihare, “Database-independent molecular formula annotation using Gibbs sampling through ZODIAC,” *Nature Machine Intelligence*, vol. 2, no. 10, pp. 629–641, 2020.
- [112] Z. Tu and C. W. Coley, “Permutation Invariant Graph-to-Sequence Model for Template-Free Retrosynthesis and Reaction Prediction,” *Journal of Chemical Information and Modeling*, vol. 62, no. 15, pp. 3503–3513, 2022.
- [113] Z. Dai, Z. Yang, Y. Yang, J. Carbonell, Q. V. Le, and R. Salakhutdinov, “Transformer-XL: Attentive Language Models Beyond a Fixed-Length Context,” *arXiv preprint arXiv:1901.02860*, 2019.
- [114] D. Rogers and M. Hahn, “Extended-Connectivity Fingerprints,” *Journal of Chemical Information and Modeling*, vol. 50, no. 5, pp. 742–754, 2010.
- [115] M. Gutmann and A. Hyvärinen, “Noise-contrastive estimation: A new estimation principle for unnormalized statistical models,” in *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, JMLR Workshop and Conference Proceedings, 2010, pp. 297–304.
- [116] L. Liu, H. Jiang, P. He, W. Chen, X. Liu, J. Gao, and J. Han, “On the Variance of the Adaptive Learning Rate and Beyond,” in *International Conference on Learning Representations*, 2020.
- [117] D. S. Wishart, D. Tzur, C. Knox, R. Eisner, A. C. Guo, N. Young, D. Cheng, K. Jewell, D. Arndt, and S. Sawhney, “HMDB: the Human Metabolome Database,” *Nucleic Acids Research*, vol. 35, no. suppl_1, pp. D521–D526, 2007.
- [118] Y. Shinbo, Y. Nakamura, M. Altaf-Ul-Amin, H. Asahi, K. Kurokawa, M. Arita, K. Saito, D. Ohta, D. Shibata, and S. Kanaya, “KNAPSAcK: a comprehensive species-metabolite relationship database,” *Plant Metabolomics*, pp. 165–181, 2006.
- [119] M. Kanehisa, S. Goto, S. Kawashima, and A. Nakaya, “The KEGG databases at GenomeNet,” *Nucleic Acids Research*, vol. 30, no. 1, pp. 42–46, 2002.
- [120] R. Schmid, S. Heuckeroth, A. Korf, A. Smirnov, O. Myers, T. S. Dyrland, R. Bushuiev, K. J. Murray, N. Hoffmann, M. Lu, A. Sarvepalli, Z. Zhang, M. Fleischauer, K. Dührkop, M. Wesner, S. J. Hoogstra, E. Rudt, O. Mokshyna, C. Brungs, K. Ponomarov, L. Mutabdžija, T. Damiani, C. J. Pudney, M. Earll, P. O. Helmer, T. R. Fallon, T. Schulze, A. Rivas-Ubach, A. Bilbao, H. Richter, L.-F. Nothias, M. Wang, M. Orešič, J.-K. Weng, S. Böcker, A. Jeibmann, H. Hayen, U. Karst, P. C. Dorrestein, D. Petras, X. Du, and T. Pluskal, “Integrative analysis of multimodal mass spectrometry data in MZmine 3,” *Nature Biotechnology*, vol. 41, no. 4, pp. 447–449, 2023.
- [121] L. McInnes, J. Healy, and J. Melville, “UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction,” *arXiv preprint arXiv:1802.03426*, 2018.

- [122] J. Bergstra, R. Bardenet, Y. Bengio, and B. Kégl, “Algorithms for Hyper-Parameter Optimization,” in *Advances in Neural Information Processing Systems 24*, 2011.
- [123] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama, “Optuna: A Next-generation Hyperparameter Optimization Framework,” in *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, 2019, pp. 2623–2631.
- [124] R. Liaw, E. Liang, R. Nishihara, P. Moritz, J. E. Gonzalez, and I. Stoica, “Tune: A Research Platform for Distributed Model Selection and Training,” in *ICML AutoML Workshop*, 2018.
- [125] S. Goldman, J. Xin, J. Provenzano, and C. W. Coley, “MIST-CF: Chemical Formula Inference from Tandem Mass Spectra,” *Journal of Chemical Information and Modeling*, 2023.
- [126] T. Pluskal, M. P. Torrens-Spence, T. R. Fallon, A. De Abreu, C. H. Shi, and J.-K. Weng, “The biosynthetic origin of psychoactive kavalactones in kava,” *Nature Plants*, vol. 5, no. 8, pp. 867–878, 2019.
- [127] D. Paik, L. Yao, Y. Zhang, S. Bae, G. D. D’Agostino, M. Zhang, E. Kim, E. A. Franzosa, J. Avila-Pacheco, J. E. Bisanz, C. K. Rakowski, H. Vlamakis, R. J. Xavier, P. J. Turnbaugh, R. S. Longman, M. R. Krout, C. B. Clish, F. Rastinejad, C. Huttenhower, J. R. Huh, and A. S. Devlin, “Human gut bacteria produce TH17-modulating bile acid metabolites,” *Nature*, vol. 603, no. 7903, pp. 907–912, 2022.
- [128] S. Kim, P. A. Thiessen, E. E. Bolton, J. Chen, G. Fu, A. Gindulyte, L. Han, J. He, S. He, and B. A. Shoemaker, “PubChem Substance and Compound databases,” *Nucleic Acids Research*, vol. 44, no. D1, pp. D1202–D1213, 2016.
- [129] T. Pluskal, T. Uehara, and M. Yanagida, “Highly accurate chemical formula prediction tool utilizing high-resolution mass spectra, MS/MS fragmentation, heuristic rules, and isotope pattern matching,” *Analytical Chemistry*, vol. 84, no. 10, pp. 4396–4403, 2012.
- [130] S. Goldman, J. Bradshaw, J. Xin, and C. W. Coley, “Prefix-Tree Decoding for Predicting Mass Spectra from Molecules,” in *Advances in Neural Information Processing Systems 37*, 2023.
- [131] E. Pretsch, P. Bühlmann, C. Affolter, E. Pretsch, P. Bhuhlmann, and C. Affolter, *Structure determination of organic compounds*. Springer, 2000.
- [132] N. Zhang, R. Aebersold, and B. Schwikowski, “ProbID: A probabilistic algorithm to identify peptides through sequence database searching using tandem mass spectral data,” *Proteomics*, vol. 2, no. 10, pp. 1406–1412, 2002.
- [133] Y. LeCun, S. Chopra, R. Hadsell, M. Ranzato, and F. Huang, “A Tutorial on Energy-Based Learning,” 2006.
- [134] M. H. Lin, Z. Tu, and C. W. Coley, “Improving the performance of models for one-step retrosynthesis through re-ranking,” *Journal of Cheminformatics*, vol. 14, no. 1, pp. 1–13, 2022.
- [135] R. Sun, H. Dai, L. Li, S. Kearnes, and B. Dai, “Energy-based View of Retrosynthesis,” *arXiv preprint arXiv:2007.13437*, 2020.

- [136] Y. Du, J. Meier, J. Ma, R. Fergus, and A. Rives, “Energy-based models for atomic-resolution protein conformations,” in *International Conference on Learning Representations*, 2020.
- [137] C. Ying, T. Cai, S. Luo, S. Zheng, G. Ke, D. He, Y. Shen, and T.-Y. Liu, “Do Transformers Really Perform Bad for Graph Representation?” In *Advances in Neural Information Processing Systems 34*, 2021, pp. 28 877–28 888.
- [138] C. R. Harris, K. J. Millman, S. J. van der Walt, R. Gommers, P. Virtanen, D. Cournapeau, E. Wieser, J. Taylor, S. Berg, N. J. Smith, R. Kern, M. Picus, S. Hoyer, M. H. van Kerkwijk, M. Brett, A. Haldane, J. Fernández del Río, M. Wiebe, P. Peterson, P. Gérard-Marchant, K. Sheppard, T. Reddy, W. Weckesser, H. Abbasi, C. Gohlke, and T. E. Oliphant, “Array programming with NumPy,” *Nature*, vol. 585, no. 7825, pp. 357–362, 2020.
- [139] M. Tancik, P. P. Srinivasan, B. Mildenhall, S. Fridovich-Keil, N. Raghavan, U. Singhal, R. Ramamoorthi, J. T. Barron, and R. Ng, “Fourier Features Let Networks Learn High Frequency Functions in Low Dimensional Domains,” in *Advances in Neural Information Processing Systems 33*, 2020, pp. 7537–7547.
- [140] W. Falcon and The PyTorch Lightning team, *PyTorch Lightning*, 2019.
- [141] D. P. Kingma and J. Ba, “Adam: A Method for Stochastic Optimization,” in *International Conference on Learning Representations*, 2015.
- [142] S. Bocker and Z. Lipták, “A fast and simple algorithm for the Money Changing Problem,” *Algorithmica*, vol. 48, no. 4, pp. 413–432, 2007.
- [143] K. Dührkop, M. Ludwig, M. Meusel, and S. Böcker, “Faster mass decomposition,” in *Algorithms in Bioinformatics: 13th International Workshop*, 2013, pp. 45–58.
- [144] R. Schmid, D. Petras, L.-F. Nothias, M. Wang, A. T. Aron, A. Jagels, H. Tsugawa, J. Rainer, M. Garcia-Aloy, K. Dührkop, A. Korf, T. Pluskal, Z. Kameník, A. K. Jarmusch, A. M. Caraballo-Rodríguez, K. C. Weldon, M. Nothias-Esposito, A. A. Ak-senov, A. Bauermeister, A. Albarracin Orio, C. O. Grundmann, F. Vargas, I. Koester, J. M. Gauglitz, E. C. Gentry, Y. Hövelmann, S. A. Kalinina, M. A. Pendergraft, M. Panitchpakdi, R. Tehan, A. Le Gouellec, G. Aleti, H. Mannocho Russo, B. Arndt, F. Hübner, H. Hayen, H. Zhi, M. Raffatellu, K. A. Prather, L. I. Aluwihare, S. Böcker, K. L. McPhail, H.-U. Humpf, U. Karst, and P. C. Dorrestein, “Ion identity molecular networking for mass spectrometry-based metabolomics in the GNPS environment,” *Nature Communications*, vol. 12, no. 1, p. 3832, 2021.
- [145] S. Goldman, J. Li, and C. W. Coley, “Generating Molecular Fragmentation Graphs with Autoregressive Neural Networks,” *arXiv preprint arXiv:2304.13136*, 2023.
- [146] A. Bauermeister, H. Mannocho-Russo, L. V. Costa-Lotufu, A. K. Jarmusch, and P. C. Dorrestein, “Mass spectrometry-based metabolomics in microbiome investigations,” *Nature Reviews Microbiology*, vol. 20, no. 3, pp. 143–160, 2022.
- [147] J. Gasteiger, W. Hanebeck, and K. P. Schulz, “Prediction of Mass Spectra from Structural Information,” *Journal of Chemical Information and Computer Sciences*, vol. 32, no. 4, pp. 264–271, 1992.

- [148] D. P. Demarque, A. E. Crotti, R. Vessecchi, J. L. Lopes, and N. P. Lopes, "Fragmentation reactions using electrospray ionization mass spectrometry: an important tool for the structural elucidation and characterization of synthetic and natural products," *Natural Product Reports*, vol. 33, no. 3, pp. 432–455, 2016.
- [149] J. Chen, Y. Chen, Y. Jiang, H. Fu, B. Xin, and Y.-F. Zhao, "Rearrangement of P-N to P-O bonds in mass spectra of N-diisopropoxyphosphoryl amino acids/alcohols," *Rapid Communications in Mass Spectrometry*, vol. 15, no. 20, pp. 1936–1940, 2001.
- [150] H. Zhu, L. Liu, and S. Hassoun, "Using Graph Neural Networks for Mass Spectrometry Prediction," in *Machine Learning for Molecules Workshop at NeurIPS*, 2020.
- [151] F. Hufsky, K. Scheubert, and S. Böcker, "Computational mass spectrometry for small-molecule fragmentation," *TrAC Trends in Analytical Chemistry*, vol. 53, pp. 41–48, 2014.
- [152] M. Murphy, S. Jegelka, E. Fraenkel, T. Kind, D. Healey, and T. Butler, "Efficiently predicting high resolution mass spectra with graph neural networks," in *Proceedings of the 40th International Conference on Machine Learning*, 2023.
- [153] R. L. Zhu and E. Jonas, "Rapid Approximate Subset-Based Spectra Prediction for Electron Ionization–Mass Spectrometry," *Analytical Chemistry*, vol. 95, no. 5, pp. 2653–2663, 2023.
- [154] W. Bittremieux, R. Schmid, F. Huber, J. J. J. van der Hooft, M. Wang, and P. C. Dorrestein, "Comparison of Cosine, Modified Cosine, and Neutral Loss Based Spectrum Alignment For Discovery of Structurally Related Molecules," *Journal of the American Society for Mass Spectrometry*, vol. 33, no. 9, pp. 1733–1744, 2022.
- [155] O. Vinyals, S. Bengio, and M. Kudlur, "Order Matters: Sequence to sequence for sets," in *International Conference on Learning Representations*, 2016.
- [156] M. Zaheer, S. Kottur, S. Ravanbakhsh, B. Póczos, R. Salakhutdinov, and A. Smola, "Deep Sets," in *Advances in Neural Information Processing Systems 30*, 2017, pp. 3391–3401.
- [157] Y. Hong, S. Li, C. J. Welch, S. Tichy, Y. Ye, and H. Tang, "3DMolMS: prediction of tandem mass spectra from 3D molecular conformations," *Bioinformatics*, vol. 39, no. 6, btad354, 2023.
- [158] G. Landrum, "RDKit: Open-source cheminformatics software," 2016.
- [159] M. A. Hoffmann, F. Kretschmer, M. Ludwig, and S. Böcker, "Mad Hatter correctly annotates 98% of small molecule tandem mass spectra searching in PubChem," *Metabolites*, vol. 13, no. 3, p. 314, 2023.
- [160] B. Curry and D. E. Rumelhart, "MSnet: A neural network which classifies mass spectra," *Tetrahedron Computer Methodology*, vol. 3, no. 3-4, pp. 213–237, 1990.
- [161] B. Buchanan, G. Sutherland, and E. A. Feigenbaum, "Heuristic Dendral: A program for generating explanatory hypotheses," *Machine Intelligence*, vol. 4, pp. 209–254, 1969.

- [162] G. Sutherland, “Dendral-a computer program for generating and filtering chemical structures,” Stanford University, Department of Computer Science, Tech. Rep., 1967.
- [163] M. Yilmaz, W. Fondrie, W. Bittremieux, S. Oh, and W. S. Noble, “De novo mass spectrometry peptide sequencing with a transformer model,” in *Proceedings of the 39th International Conference on Machine Learning*, 2022, pp. 25 514–25 522.
- [164] O. Shouman, W. Gabriel, V.-G. Giurcoiu, V. Sternlicht, and M. Wilhelm, “PROSPECT: Labeled Tandem Mass Spectrometry Dataset for Machine Learning in Proteomics,” in *Advances in Neural Information Processing Systems 35*, 2022.
- [165] Y. Zhang, J. Hare, and A. Prugel-Bennett, “Deep Set Prediction Networks,” in *Advances in Neural Information Processing Systems 32*, 2019.
- [166] F. Locatello, D. Weissenborn, T. Unterthiner, A. Mahendran, G. Heigold, J. Uszkoreit, A. Dosovitskiy, and T. Kipf, “Object-centric learning with Slot Attention,” in *Advances in Neural Information Processing Systems 33*, 2020.
- [167] A. R. Kosioerek, H. Kim, and D. J. Rezende, “Conditional Set Generation with Transformers,” in *Workshop on Object-Oriented Learning at ICML*, 2020.
- [168] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [169] W. L. Hamilton, “Graph Representation Learning,” *Synthesis Lectures on Artificial Intelligence and Machine Learning*, vol. 14, no. 3, pp. 1–159, 2020.
- [170] M. M. Bronstein, J. Bruna, T. Cohen, and P. Veličković, “Geometric Deep Learning: Grids, Groups, Graphs, Geodesics, and Gauges,” *arXiv preprint arXiv:2104.13478*, 2021.
- [171] P. W. Battaglia, J. B. Hamrick, V. Bapst, A. Sanchez-Gonzalez, V. Zambaldi, M. Malinowski, A. Tacchetti, D. Raposo, A. Santoro, R. Faulkner, C. Gulcehre, F. Song, A. Ballard, J. Gilmer, G. Dahl, A. Vaswani, K. Allen, C. Nash, V. Langston, C. Dyer, N. Heess, D. Wierstra, P. Kohli, M. Botvinick, O. Vinyals, Y. Li, and R. Pascanu, “Relational inductive biases, deep learning, and graph networks,” *arXiv preprint arXiv:1806.01261*, 2018.
- [172] Y. Li, D. Tarlow, M. Brockschmidt, and R. Zemel, “Gated Graph Sequence Neural Networks,” in *International Conference on Learning Representations*, 2016.
- [173] M. Wang, D. Zheng, Z. Ye, Q. Gan, M. Li, X. Song, J. Zhou, C. Ma, L. Yu, Y. Gai, T. Xiao, T. He, G. Karypis, J. Li, and Z. Zhang, “Deep Graph Library: A Graph-Centric, Highly-Performant Package for Graph Neural Networks,” *arXiv preprint arXiv:1909.01315*, 2019.
- [174] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, *et al.*, “PyTorch: An Imperative Style, High-Performance Deep Learning Library,” in *Advances in Neural Information Processing Systems 32*, 2019.
- [175] M. Szeremeta, K. Pietrowska, A. Niemcunowicz-Janica, A. Kretowski, and M. Ci-borowski, “Applications of Metabolomics in Forensic Toxicology and Forensic Medicine,” *International Journal of Molecular Sciences*, vol. 22, no. 6, 2021.

- [176] P. Kirkpatrick and C. Ellis, "Chemical space," *Nature*, vol. 432, no. 7019, pp. 823–824, 2004.
- [177] B. E. Frewen, G. E. Merrihew, C. C. Wu, W. S. Noble, and M. J. MacCoss, "Analysis of Peptide MS/MS Spectra from Large-Scale Proteomics Experiments Using Spectrum Libraries," *Analytical Chemistry*, vol. 78, no. 16, pp. 5678–5684, 2006.
- [178] B. Weisfeiler and A. Leman, "The reduction of a graph to canonical form and the algebra which appears therein," *nti, Series*, vol. 2, no. 9, pp. 12–16, 1968.
- [179] J. Bradshaw, B. Paige, M. J. Kusner, M. Segler, and J. M. Hernández-Lobato, "Barking up the right tree: an approach to search over molecule synthesis DAGs," in *Advances in Neural Information Processing Systems 33*, 2020, pp. 6852–6866.
- [180] W. Gao, R. Mercado, and C. W. Coley, "Amortized Tree Generation for Bottom-up Synthesis Planning and Synthesizable Molecular Design," in *International Conference on Learning Representations*, 2021.
- [181] F. W. McLafferty, "Mass Spectrometric Analysis. Molecular Rearrangements," *Analytical Chemistry*, vol. 31, no. 1, pp. 82–87, 1959.
- [182] P. Ertl and A. Schuffenhauer, "Estimation of synthetic accessibility score of drug-like molecules based on molecular complexity and fragment contributions," *Journal of Cheminformatics*, vol. 1, pp. 1–11, 2009.
- [183] K. Huang, T. Fu, W. Gao, Y. Zhao, Y. Roohani, J. Leskovec, C. W. Coley, C. Xiao, J. Sun, and M. Zitnik, "Artificial intelligence foundation for therapeutic science," *Nature Chemical Biology*, vol. 18, no. 10, pp. 1033–1036, 2022.
- [184] J. A. Fine, A. A. Rajasekar, K. P. Jethava, and G. Chopra, "Spectral deep learning for prediction and prospective validation of functional groups," *Chemical Science*, vol. 11, no. 18, pp. 4618–4630, 2020.
- [185] J. R. Doroghazi, J. C. Albright, A. W. Goering, K.-S. Ju, R. R. Haines, K. A. Tchalukov, D. P. Labeda, N. L. Kelleher, and W. W. Metcalf, "A roadmap for natural product discovery based on large-scale genomics and metabolomics," *Nature Chemical Biology*, vol. 10, no. 11, pp. 963–968, 2014.
- [186] S. Goldman, R. Das, K. K. Yang, and C. W. Coley, "Machine learning modeling of family wide enzyme-substrate specificity screens," *PLoS Computational Biology*, vol. 18, no. 2, e1009853, 2022.
- [187] G.-M. Lin, R. Warden-Rothman, and C. A. Voigt, "Retrosynthetic design of metabolic pathways to chemicals not found in nature," *Current Opinion in Systems Biology*, vol. 14, pp. 82–107, 2019.
- [188] S. Wu, R. Snajdrova, J. C. Moore, K. Baldenius, and U. T. Bornscheuer, "Biocatalysis: Enzymatic synthesis for industrial applications," *Angewandte Chemie International Edition*, vol. 60, no. 1, pp. 88–119, 2021.
- [189] F. H. Arnold, "Directed evolution: bringing new chemistry to life," *Angewandte Chemie International Edition*, vol. 57, no. 16, pp. 4143–4148, 2018.

- [190] C. A. Voigt, “Synthetic biology 2020–2030: six commercially-available products that are changing our world,” *Nature Communications*, vol. 11, no. 1, pp. 1–6, 2020.
- [191] M. A. Huffman, A. Fryszkowska, O. Alvizo, M. Borra-Garske, K. R. Campos, P. N. Devine, D. Duan, J. H. Forstater, S. T. Grosser, and H. M. Halsey, “Design of an in vitro biocatalytic cascade for the manufacture of islatravir,” *Science*, vol. 366, no. 6470, pp. 1255–1259, 2019.
- [192] C. W. Coley, W. H. Green, and K. F. Jensen, “Machine Learning in Computer-Aided Synthesis Planning,” *Accounts of chemical research*, vol. 51, no. 5, pp. 1281–1289, 2018.
- [193] S. Zheng, T. Zeng, C. Li, B. Chen, C. W. Coley, Y. Yang, and R. Wu, “Deep learning driven biosynthetic pathways navigation for natural products with BioNavi-NP,” *Nature Communications*, vol. 13, no. 1, p. 3342, 2022.
- [194] W. Finnigan, L. J. Hepworth, S. L. Flitsch, and N. J. Turner, “RetroBioCat as a computer-aided synthesis planning tool for biocatalytic reactions and cascades,” *Nature Catalysis*, vol. 4, no. 2, pp. 98–104, 2021.
- [195] M. Koch, T. Duigou, and J.-L. Faulon, “Reinforcement Learning for Bioretrosynthesis,” *ACS Synthetic Biology*, vol. 9, no. 1, pp. 157–168, 2019.
- [196] P. Carbonell, J. Wong, N. Swainston, E. Takano, N. J. Turner, N. S. Scrutton, D. B. Kell, R. Breitling, and J.-L. Faulon, “Selenzyme: Enzyme selection tool for pathway design,” *Bioinformatics*, vol. 34, no. 12, pp. 2153–2154, 2018.
- [197] A. Cho, H. Yun, J. H. Park, S. Y. Lee, and S. Park, “Prediction of novel synthetic pathways for the production of desired chemicals,” *BMC Systems Biology*, vol. 4, no. 1, pp. 1–16, 2010.
- [198] K. Hult and P. Berglund, “Enzyme promiscuity: mechanism and applications,” *Trends in Biotechnology*, vol. 25, no. 5, pp. 231–238, 2007.
- [199] C. N. Bedbrook, K. K. Yang, J. E. Robinson, E. D. Mackey, V. Gradinaru, and F. H. Arnold, “Machine learning-guided channelrhodopsin engineering enables minimally invasive optogenetics,” *Nature Methods*, vol. 16, no. 11, pp. 1176–1184, 2019.
- [200] Z. Wu, S. J. Kan, R. D. Lewis, B. J. Wittmann, and F. H. Arnold, “Machine learning-assisted directed protein evolution with combinatorial libraries,” *Proceedings of the National Academy of Sciences*, vol. 116, no. 18, pp. 8852–8858, 2019.
- [201] R. J. Fox, S. C. Davis, E. C. Mundorff, L. M. Newman, V. Gavrilovic, S. K. Ma, L. M. Chung, C. Ching, S. Tam, S. Muley, *et al.*, “Improving catalytic function by ProSAR-driven enzyme evolution,” *Nature Biotechnology*, vol. 25, no. 3, pp. 338–344, 2007.
- [202] T. Bepler and B. Berger, “Learning protein sequence embeddings using information from structure,” in *International Conference on Learning Representations*, 2019.
- [203] T. Bepler and B. Berger, “Learning the protein language: Evolution, structure, and function,” *Cell Systems*, vol. 12, no. 6, pp. 654–669, 2021.

- [204] R. Rao, N. Bhattacharya, N. Thomas, Y. Duan, P. Chen, J. Canny, P. Abbeel, and Y. Song, “Evaluating Protein Transfer Learning with TAPE,” in *Advances in Neural Information Processing Systems 33*, 2019, pp. 9689–9701.
- [205] R. M. Rao, J. Meier, T. Sercu, S. Ovchinnikov, and A. Rives, “Transformer protein language models are unsupervised structure learners,” in *International Conference on Learning Representations*, 2020.
- [206] A. Rives, J. Meier, T. Sercu, S. Goyal, Z. Lin, J. Liu, D. Guo, M. Ott, C. L. Zitnick, J. Ma, and R. Fergus, “Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences,” *Proceedings of the National Academy of Sciences*, vol. 118, no. 15, e2016239118, 2021.
- [207] E. C. Alley, G. Khimulya, S. Biswas, M. AlQuraishi, and G. M. Church, “Unified rational protein engineering with sequence-based deep representation learning,” *Nature Methods*, vol. 16, no. 12, pp. 1315–1322, 2019.
- [208] C. Dallago, K. Schütze, M. Heinzinger, T. Olenyi, M. Littmann, A. X. Lu, K. K. Yang, S. Min, S. Yoon, J. T. Morton, *et al.*, “Learned embeddings from deep learning to visualize and predict protein sets,” *Current Protocols*, vol. 1, no. 5, e113, 2021.
- [209] S. Biswas, G. Khimulya, E. C. Alley, K. M. Esvelt, and G. M. Church, “Low-N protein engineering with data-efficient deep learning,” *Nature Methods*, vol. 18, no. 4, pp. 389–396, 2021.
- [210] J. C. Greenhalgh, S. A. Fahlberg, B. F. Pfeleger, and P. A. Romero, “Machine learning-guided acyl-ACP reductase engineering for improved in vivo fatty alcohol production,” *Nature Communications*, vol. 12, no. 1, p. 5825, 2021.
- [211] B. M. Bonk, Y. Tarasova, M. A. Hicks, B. Tidor, and K. L. Prather, “Rational design of thiolase substrate specificity for metabolic engineering applications,” *Biotechnology and Bioengineering*, vol. 115, no. 9, pp. 2167–2182, 2018.
- [212] R. C. de Melo-Minardi, K. Bastard, and F. Artiguenave, “Identification of subfamily-specific sites based on active sites modeling and clustering,” *Bioinformatics*, vol. 26, no. 24, pp. 3075–3082, 2010.
- [213] G. Rix, E. J. Watkins-Dulaney, P. J. Almhjell, C. E. Boville, F. H. Arnold, and C. C. Liu, “Scalable continuous evolution for the generation of diverse enzyme variants encompassing promiscuous activities,” *Nature Communications*, vol. 11, no. 1, p. 5644, 2020.
- [214] K. Chen and F. H. Arnold, “Enzyme Engineering for Nonaqueous Solvents: Random Mutagenesis to Enhance Activity of Subtilisin E in Polar Organic Media,” *Bio/Technology*, vol. 9, no. 11, pp. 1073–1077, 1991.
- [215] P. A. Romero and F. H. Arnold, “Exploring protein fitness landscapes by directed evolution,” *Nature Reviews Molecular Cell Biology*, vol. 10, no. 12, pp. 866–876, 2009.
- [216] K. Chen and F. H. Arnold, “Engineering cytochrome P450s for enantioselective cyclopropanation of internal alkynes,” *Journal of the American Chemical Society*, vol. 142, no. 15, pp. 6891–6895, 2020.

- [217] C. Corre and G. L. Challis, “New natural product biosynthetic chemistry discovered by genome mining,” *Natural Product Reports*, vol. 26, no. 8, pp. 977–986, 2009.
- [218] P.-Y. Colin, B. Kintsjes, F. Gielen, C. M. Miton, G. Fischer, M. F. Mohamed, M. Hyvönen, D. P. Morgavi, D. B. Janssen, and F. Hollfelder, “Ultra-high-throughput discovery of promiscuous enzymes by picodroplet functional metagenomics,” *Nature Communications*, vol. 6, no. 1, pp. 1–12, 2015.
- [219] B. F. Fisher, H. M. Snodgrass, K. A. Jones, M. C. Andorfer, and J. C. Lewis, “Site-Selective C–H Halogenation Using Flavin-Dependent Halogenases Identified via Family-Wide Activity Profiling,” *ACS Central Science*, vol. 5, no. 11, pp. 1844–1856, 2019.
- [220] J. Hon, S. Borko, J. Stourac, Z. Prokop, J. Zendulka, D. Bednar, T. Martinek, and J. Damborsky, “EnzymeMiner: automated mining of soluble enzymes with diverse structures, catalytic properties and stabilities,” *Nucleic Acids Research*, vol. 48, no. W1, W104–W109, 2020.
- [221] J. R. Marshall, P. Yao, S. L. Montgomery, J. D. Finnigan, T. W. Thorpe, R. B. Palmer, J. Mangas-Sanchez, R. A. M. Duncan, R. S. Heath, K. M. Graham, D. J. Cook, S. J. Charnock, and N. J. Turner, “Screening and characterization of a diverse panel of metagenomic imine reductases for biocatalytic reductive amination,” *Nature Chemistry*, vol. 13, no. 2, pp. 140–148, 2020.
- [222] E. E. Kempa, J. L. Galman, F. Parmeggiani, J. R. Marshall, J. Malassis, C. Q. Fontenelle, J.-B. Vendeville, B. Linclau, S. J. Charnock, and S. L. Flitsch, “Rapid Screening of Diverse Biotransformations for Enzyme Evolution,” *JACS Au*, vol. 1, no. 4, pp. 508–516, 2021.
- [223] M. I. Davis, J. P. Hunt, S. Herrgard, P. Ciceri, L. M. Wodicka, G. Pallares, M. Hocker, D. K. Treiber, and P. P. Zarrinkar, “Comprehensive analysis of kinase inhibitor selectivity,” *Nature Biotechnology*, vol. 29, no. 11, pp. 1046–1051, 2011.
- [224] N. J. Schauer, R. S. Magin, X. Liu, L. M. Doherty, and S. J. Buhrlage, “Advances in discovering deubiquitinating enzyme (DUB) inhibitors,” *Journal of Medicinal Chemistry*, vol. 63, no. 6, pp. 2731–2750, 2019.
- [225] A. Ernst, G. Avvakumov, J. Tong, Y. Fan, Y. Zhao, P. Alberts, A. Persaud, J. R. Walker, A.-M. Neculai, D. Neculai, A. Vorobyov, P. Garg, L. Beatty, P.-K. Chan, Y.-C. Juang, M.-C. Landry, C. Yeh, E. Zeqiraj, K. Karamboulas, A. Allali-Hassani, M. Vedadi, M. Tyers, J. Moffat, F. Sicheri, L. Pelletier, D. Durocher, B. Raught, D. Rotin, J. Yang, M. F. Moran, S. Dhe-Paganon, and S. S. Sidhu, “A Strategy for Modulation of Enzymes in the Ubiquitin System,” *Science*, vol. 339, no. 6119, pp. 590–595, 2013.
- [226] B. Hie, B. D. Bryson, and B. Berger, “Leveraging uncertainty in machine learning accelerates biological discovery and design,” *Cell Systems*, vol. 11, no. 5, pp. 461–477, 2020.
- [227] S. Li, F. Wan, H. Shu, T. Jiang, D. Zhao, and J. Zeng, “MONN: a multi-objective neural network for predicting compound-protein interactions and affinities,” *Cell Systems*, vol. 10, no. 4, pp. 308–322, 2020.

- [228] I. Schomburg, A. Chang, C. Ebeling, M. Gremse, C. Heldt, G. Huhn, and D. Schomburg, “BRENDA, the enzyme database: updates and major new developments,” *Nucleic Acids Research*, vol. 32, no. suppl_1, pp. D431–D433, 2004.
- [229] M. Yang, C. Fehl, K. V. Lees, E.-K. Lim, W. A. Offen, G. J. Davies, D. J. Bowles, M. G. Davidson, S. J. Roberts, and B. G. Davis, “Functional and informatics analysis enables glycosyltransferase activity prediction,” *Nature Chemical Biology*, vol. 14, no. 12, pp. 1109–1117, 2018.
- [230] S. L. Robinson, M. D. Smith, J. E. Richman, K. G. Aukema, and L. P. Wackett, “Machine learning-based prediction of activity and substrate specificity for OleA enzymes in the thiolase superfamily,” *Synthetic Biology*, vol. 5, no. 1, ysa004, 2020.
- [231] K. Bastard, A. A. T. Smith, C. Vergne-Vaxelaire, A. Perret, A. Zapparucha, R. De Melo-Minardi, A. Mariage, M. Boutard, A. Debard, C. Lechaplais, C. Pelle, V. Pellouin, N. Perchat, J.-L. Petit, A. Kreimeyer, C. Medigue, J. Weissenbach, F. Artiguenave, V. De Berardinis, D. Vallenet, and M. Salanoubat, “Revealing the hidden functional diversity of an enzyme family,” *Nature Chemical Biology*, vol. 10, no. 1, pp. 42–49, 2014.
- [232] M. Martínez-Martínez, C. Coscolín, G. Santiago, J. Chow, P. J. Stogios, R. Bargiela, C. Gertler, J. Navarro-Fernández, A. Bollinger, S. Thies, C. Méndez-García, A. Popovic, G. Brown, T. N. Chernikova, A. García-Moyano, G. E. K. Bjerga, P. Pérez-García, T. Hai, M. V. Del Pozo, R. Stokke, I. H. Steen, H. Cui, X. Xu, B. P. Nocek, M. Alcaide, M. Distaso, V. Mesa, A. I. Peláez, J. Sánchez, P. C. F. Buchholz, J. Pleiss, A. Fernández-Guerra, F. O. Glöckner, O. V. Golyshina, M. M. Yakimov, A. Savchenko, K.-E. Jaeger, A. F. Yakunin, W. R. Streit, P. N. Golyshin, V. Guallar, M. Ferrer, and The INMARE Consortium, “Determinants and Prediction of Esterase Substrate Promiscuity Patterns,” *ACS Chemical Biology*, vol. 13, no. 1, pp. 225–234, 2018.
- [233] H. Huang, C. Pandya, C. Liu, N. F. Al-Obaidi, M. Wang, L. Zheng, S. T. Keating, M. Aono, J. D. Love, B. Evans, R. D. Seidel, B. S. Hillerich, S. J. Garforth, S. C. Almo, P. S. Mariano, D. Dunaway-Mariano, K. N. Allen, and J. D. Farelli, “Panoramic view of a superfamily of phosphatases through substrate profiling,” *Proceedings of the National Academy of Sciences*, vol. 112, no. 16, E1974–E1983, 2015.
- [234] P. Kim, R. Winter, and D.-A. Clevert, “Deep Protein-Ligand Binding Prediction Using Unsupervised Learned Representations,” *International Journal of Molecular Sciences*, vol. 22, no. 23, p. 12 882, 2021.
- [235] M. Heinzinger, A. Elnaggar, Y. Wang, C. Dallago, D. Nechaev, F. Matthes, and B. Rost, “Modeling aspects of the language of life through transfer-learning protein sequences,” *BMC Bioinformatics*, vol. 20, no. 1, pp. 1–17, 2019.
- [236] W. Jin, R. Barzilay, and T. Jaakkola, “Junction Tree Variational Autoencoder for Molecular Graph Generation,” in *International Conference on Machine Learning*, PMLR, 2018, pp. 2323–2332.
- [237] H. L. Morgan, “The generation of a unique machine description for chemical structures—a technique developed at chemical abstracts service,” *Journal of Chemical Documentation*, vol. 5, no. 2, pp. 107–113, 1965.

- [238] K. Huang, T. Fu, W. Gao, Y. Zhao, Y. Roohani, J. Leskovec, C. W. Coley, C. Xiao, J. Sun, and M. Zitnik, “Therapeutics Data Commons: Machine Learning Datasets and Tasks for Drug Discovery and Development,” in *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*, 2021.
- [239] A. Shanehsazzadeh, D. Belanger, and D. Dohan, “Is Transfer Learning Necessary for Protein Landscape Prediction?” *arXiv preprint arXiv:2011.03443*, 2020.
- [240] T. Lu, A. X. Lu, and A. M. Moses, “Random Embeddings and Linear Regression can Predict Protein Function,” *arXiv preprint arXiv:2104.14661*, 2021.
- [241] N. S. Detlefsen, S. Hauberg, and W. Boomsma, “What is a meaningful representation of protein sequences?” *arXiv preprint arXiv:2012.02679*, 2020.
- [242] E. F. Pettersen, T. D. Goddard, C. C. Huang, G. S. Couch, D. M. Greenblatt, E. C. Meng, and T. E. Ferrin, “UCSF Chimera—a visualization system for exploratory research and analysis,” *Journal of Computational Chemistry*, vol. 25, no. 13, pp. 1605–1612, 2004.
- [243] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, and V. Dubourg, “Scikit-learn: Machine learning in Python,” *The Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.
- [244] F. Chollet *et al.*, *Keras*, <https://keras.io>, 2015.
- [245] P. J. Cock, T. Antao, J. T. Chang, B. A. Chapman, C. J. Cox, A. Dalke, I. Friedberg, T. Hamelryck, F. Kauff, and B. Wilczynski, “Biopython: freely available Python tools for computational molecular biology and bioinformatics,” *Bioinformatics*, vol. 25, no. 11, pp. 1422–1423, 2009.
- [246] M. Swain, “CIRpy-A Python interface for the Chemical Identifier Resolver (CIR),” *Matt Swain’s Blog*, 2012.
- [247] U. Consortium, “UniProt: a hub for protein information,” *Nucleic Acids Research*, vol. 43, no. D1, pp. D204–D212, 2015.
- [248] A. Bateman, L. Coin, R. Durbin, R. D. Finn, V. Hollich, S. Griffiths-Jones, A. Khanna, M. Marshall, S. Moxon, and E. L. Sonnhammer, “The Pfam protein families database,” *Nucleic Acids Research*, vol. 32, no. suppl_1, pp. D138–D141, 2004.
- [249] R. D. Finn, J. Clements, and S. R. Eddy, “HMMER web server: interactive sequence similarity searching,” *Nucleic Acids Research*, vol. 39, no. suppl_2, W29–W37, 2011.
- [250] A. P. Soleimany, A. Amini, S. Goldman, D. Rus, S. N. Bhatia, and C. W. Coley, “Evidential Deep Learning for Guided Molecular Property Prediction and Discovery,” *ACS Central Science*, vol. 7, no. 8, pp. 1356–1367, 2021.
- [251] A. Nigam, R. Pollice, M. F. Hurley, R. J. Hickman, M. Aldeghi, N. Yoshikawa, S. Chithrananda, V. A. Voelz, and A. Aspuru-Guzik, “Assigning Confidence to Molecular Property Prediction,” *Expert Opinion on Drug Discovery*, vol. 16, no. 9, pp. 1009–1023, 2021.
- [252] G. Lamb and B. Paige, “Bayesian Graph Neural Networks for Molecular Property Prediction,” *arXiv preprint arXiv:2012.02089*, 2020.

- [253] J. Jiménez-Luna, F. Grisoni, and G. Schneider, “Drug discovery with explainable artificial intelligence,” *Nature Machine Intelligence*, vol. 2, no. 10, pp. 573–584, 2020.
- [254] A. Cherkasov, E. N. Muratov, D. Fourches, A. Varnek, I. I. Baskin, M. Cronin, J. Dearden, P. Gramatica, Y. C. Martin, R. Todeschini, *et al.*, “QSAR modeling: where have you been? Where are you going to?” *Journal of Medicinal Chemistry*, vol. 57, no. 12, pp. 4977–5010, 2014.
- [255] E. N. Muratov, J. Bajorath, R. P. Sheridan, I. V. Tetko, D. Filimonov, V. Poroikov, T. I. Oprea, I. I. Baskin, A. Varnek, A. Roitberg, *et al.*, “QSAR without borders,” *Chemical Society Reviews*, vol. 49, no. 11, pp. 3525–3564, 2020.
- [256] A. Kendall and Y. Gal, “What Uncertainties Do We Need in Bayesian Deep Learning for Computer Vision?” In *Advances in Neural Information Processing Systems 30*, 2017.
- [257] N. S. Eyke, W. H. Green, and K. F. Jensen, “Iterative experimental design based on active machine learning reduces the experimental burden associated with reaction screening,” *Reaction Chemistry & Engineering*, vol. 5, no. 10, pp. 1963–1972, 2020.
- [258] A. G. Roy, J. Ren, S. Azizi, A. Loh, V. Natarajan, B. Mustafa, N. Pawlowski, J. Freyberg, Y. Liu, Z. Beaver, *et al.*, “Does Your Dermatology Classifier Know What It Doesn’t Know? Detecting the Long-Tail of Unseen Conditions,” *arXiv preprint arXiv:2104.03829*, 2021.
- [259] J. P. Janet, C. Duan, T. Yang, A. Nandy, and H. J. Kulik, “A quantitative uncertainty metric controls error in neural network-driven chemical discovery,” *Chemical Science*, vol. 10, no. 34, pp. 7913–7922, 2019.
- [260] K. Tran, W. Neiswanger, J. Yoon, Q. Zhang, E. Xing, and Z. W. Ulissi, “Methods for comparing uncertainty quantifications for material property predictions,” *Machine Learning: Science and Technology*, vol. 1, no. 2, p. 025 006, 2020.
- [261] L. V. Jospin, W. Buntine, F. Boussaid, H. Laga, and M. Bennamoun, “Hands-on Bayesian Neural Networks—a Tutorial for Deep Learning Users,” *arXiv preprint arXiv:2007.06823*, 2020.
- [262] B. Lakshminarayanan, A. Pritzel, and C. Blundell, “Simple and Scalable Predictive Uncertainty Estimation using Deep Ensembles,” in *Advances in Neural Information Processing Systems 30*, 2017, pp. 6402–6413.
- [263] Y. Gal and Z. Ghahramani, “Dropout as a Bayesian approximation: Representing model uncertainty in deep learning,” in *Proceedings of the 33rd International Conference on Machine Learning*, 2016.
- [264] L. Hirschfeld, K. Swanson, K. Yang, R. Barzilay, and C. W. Coley, “Uncertainty Quantification Using Neural Networks for Molecular Property Prediction,” *Journal of Chemical Information and Modeling*, vol. 60, no. 8, pp. 3770–3780, 2020.
- [265] J. S. Smith, B. Nebgen, N. Lubbers, O. Isayev, and A. E. Roitberg, “Less is more: Sampling chemical space with active learning,” *The Journal of Chemical Physics*, vol. 148, no. 24, p. 241 733, 2018.

- [266] J. Vandermause, S. B. Torrisi, S. Batzner, Y. Xie, L. Sun, A. M. Kolpak, and B. Kozinsky, "On-the-fly active learning of interpretable Bayesian force fields for atomistic rare events," *npj Computational Materials*, vol. 6, no. 1, pp. 1–11, 2020.
- [267] J. Klicpera, S. Giri, J. T. Margraf, and S. Günnemann, "Fast and Uncertainty-Aware Directional Message Passing for Non-Equilibrium Molecules," *arXiv preprint arXiv:2011.14115*, 2020.
- [268] J. Klicpera, J. Groß, and S. Günnemann, "Directional Message Passing for Molecular Graphs," in *International Conference on Learning Representations*, 2020.
- [269] D. Nix and A. Weigend, "Estimating the Mean and Variance of the Target Probability Distribution," in *Proceedings of IEEE International Conference on Neural Networks*, 1994.
- [270] C. M. Bishop, "Mixture density networks," 1994.
- [271] I. Gilitschenski, R. Sahoo, W. Schwarting, A. Amini, S. Karaman, and D. Rus, "Deep Orientation Uncertainty Learning based on a Bingham Loss," in *International Conference on Learning Representations*, 2019.
- [272] A. Amini, A. P. Soleimany, W. Schwarting, S. N. Bhatia, and D. Rus, "Uncovering and Mitigating Algorithmic Bias through Learned Latent Structure," in *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, 2019, pp. 289–295.
- [273] P. Gurevich and H. Stuke, "Gradient conjugate priors and multi-layer neural networks," *Artificial Intelligence*, vol. 278, p. 103 184, 2020.
- [274] G. Scalia, C. A. Grambow, B. Pernici, Y.-P. Li, and W. H. Green, "Evaluating Scalable Uncertainty Estimation Methods for Deep Learning Based Molecular Property Prediction," *Journal of Chemical Information and Modeling*, vol. 60, no. 6, pp. 2697–2717, 2020.
- [275] M. Sensoy, L. Kaplan, and M. Kandemir, "Evidential Deep Learning to Quantify Classification Uncertainty," in *Advances in Neural Information Processing Systems 31*, 2018, pp. 3179–3189.
- [276] A. Amini, W. Schwarting, A. Soleimany, and D. Rus, "Deep Evidential Regression," in *Advances in Neural Information Processing Systems 33*, 2020, pp. 14 927–14 937.
- [277] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. The MIT Press, 2016.
- [278] Z. Wu, B. Ramsundar, E. N. Feinberg, J. Gomes, C. Geniesse, A. S. Pappu, K. Leswing, and V. Pande, "MoleculeNet: a benchmark for molecular machine learning," *Chemical Science*, vol. 9, no. 2, pp. 513–530, 2018.
- [279] C. Gorgulla, A. Boeszoermyenyi, Z.-F. Wang, P. D. Fischer, P. W. Coote, K. M. P. Das, Y. S. Malets, D. S. Radchenko, Y. S. Moroz, D. A. Scott, *et al.*, "An open-source drug discovery platform enables ultra-large virtual screens," *Nature*, vol. 580, no. 7805, pp. 663–668, 2020.
- [280] R. Ramakrishnan, P. O. Dral, M. Rupp, and O. A. Von Lilienfeld, "Quantum chemistry structures and properties of 134 kilo molecules," *Scientific Data*, vol. 1, pp. 1–7, 2014.

- [281] S. Jastrzębski, M. Szymczak, A. Pocha, S. Mordalski, J. Tabor, A. J. Bojarski, and S. Podlowska, “Emulating Docking Results Using a Deep Neural Network: A New Perspective for Virtual Screening,” *Journal of Chemical Information and Modeling*, vol. 60, no. 9, pp. 4246–4262, 2020.
- [282] D. E. Graff, E. I. Shakhnovich, and C. W. Coley, “Accelerating high-throughput virtual screening through molecular pool-based active learning,” *Chemical Science*, vol. 12, no. 22, pp. 7866–7881, 2021.
- [283] O. Trott and A. J. Olson, “AutoDock Vina: improving the speed and accuracy of docking with a new scoring function, efficient optimization, and multithreading,” *Journal of Computational Chemistry*, vol. 31, no. 2, pp. 455–461, 2010.
- [284] V. Kuleshov, N. Fenner, and S. Ermon, “Accurate Uncertainties for Deep Learning Using Calibrated Regression,” in *Proceedings of the 35th International Conference on Machine Learning*, 2018, pp. 2796–2804.
- [285] P. I. Frazier, “A tutorial on Bayesian optimization,” *arXiv preprint arXiv:1807.02811*, 2018.
- [286] J. M. Hernández-Lobato, J. Requeima, E. O. Pyzer-Knapp, and A. Aspuru-Guzik, “Parallel and Distributed Thompson Sampling for Large-scale Accelerated Exploration of Chemical Space,” in *Proceedings of the 34th International Conference on Machine Learning*, 2017, pp. 1470–1479.
- [287] S. M. Corsello, J. A. Bittker, Z. Liu, J. Gould, P. McCarren, J. E. Hirschman, S. E. Johnston, A. Vrcic, B. Wong, M. Khan, *et al.*, “The Drug Repurposing Hub: a next-generation drug library and information resource,” *Nature Medicine*, vol. 23, no. 4, pp. 405–408, 2017.
- [288] C. Blundell, J. Cornebise, K. Kavukcuoglu, and D. Wierstra, “Weight Uncertainty in Neural Network,” in *Proceedings of the 32nd International Conference on Machine Learning*, 2015.
- [289] G. Imbalzano, Y. Zhuang, V. Kapil, K. Rossi, E. A. Engel, F. Grasselli, and M. Ceriotti, “Uncertainty estimation for molecular dynamics and sampling,” *The Journal of Chemical Physics*, vol. 154, no. 7, p. 074102, 2021.
- [290] D. A. Rufa, H. E. B. Macdonald, J. Fass, M. Wieder, P. B. Grinaway, A. E. Roitberg, O. Isayev, and J. D. Chodera, “Towards chemical accuracy for alchemical free energy calculations with hybrid physics-based machine learning/molecular mechanics potentials,” *BioRxiv*, 2020.
- [291] K. Schütt, P.-J. Kindermans, H. E. S. Felix, S. Chmiela, A. Tkatchenko, and K.-R. Müller, “SchNet: A continuous-filter convolutional neural network for modeling quantum interactions,” in *Advances in Neural Information Processing Systems 30*, 2017, pp. 991–1001.
- [292] R. J. Townshend, M. Vögele, P. Suriana, A. Derry, A. Powers, Y. Laloudakis, S. Balachandar, B. Anderson, S. Eismann, R. Kondor, *et al.*, “ATOM3D: Tasks on Molecules in Three Dimensions,” in *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*, 2021.

- [293] S. Kearnes, “Pursuing a Prospective Perspective,” *Trends in Chemistry*, 2020.
- [294] K. P. Murphy, “Conjugate Bayesian Analysis of the Gaussian distribution,” *def*, vol. 1, no. $2\sigma^2$, p. 16, 2007.
- [295] K. T. Schütt, P. Kessel, M. Gastegger, K. A. Nicoli, A. Tkatchenko, and K.-R. Müller, “SchNetPack: A Deep Learning Toolbox For Atomistic Systems,” *Journal of Chemical Theory and Computation*, vol. 15, no. 1, pp. 448–455, 2019.
- [296] M. Wenlock and N. Tomkinson, “Experimental in Vitro DMPK and Physicochemical Data on a Set of Publicly Disclosed Compounds.,” 2015.
- [297] H. Zhu, T. M. Martin, L. Ye, A. Sedykh, D. M. Young, and A. Tropsha, “Quantitative structure- activity relationship modeling of rat acute toxicity by oral exposure,” *Chemical Research in Toxicology*, vol. 22, no. 12, pp. 1913–1921, 2009.
- [298] L. Ruddigkeit, R. Van Deursen, L. C. Blum, and J.-L. Reymond, “Enumeration of 166 billion organic small molecules in the chemical universe database GDB-17,” *Journal of Chemical Information and Modeling*, vol. 52, no. 11, pp. 2864–2875, 2012.
- [299] P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski, P. Peterson, W. Weckesser, J. Bright, *et al.*, “SciPy 1.0: fundamental algorithms for scientific computing in Python,” *Nature Methods*, vol. 17, no. 3, pp. 261–272, 2020.
- [300] G. Lample, A. Conneau, L. Denoyer, and M. Ranzato, “Unsupervised Machine Translation Using Monolingual Corpora Only,” in *International Conference on Learning Representations*, 2018.
- [301] M. A. Schorn, S. Verhoeven, L. Ridder, F. Huber, D. D. Acharya, A. A. Aksenov, G. Aleti, J. A. Moghaddam, A. T. Aron, S. Aziz, *et al.*, “A community resource for paired genomic and metabolomic data mining,” *Nature Chemical Biology*, vol. 17, no. 4, pp. 363–368, 2021.
- [302] E. Karunaratne, D. W. Hill, K. Dührkop, S. Böcker, and D. F. Grant, “Combining Experimental with Computational Infrared and Mass Spectra for High-Throughput Nontargeted Chemical Structure Identification,” *Analytical Chemistry*, vol. 95, no. 32, pp. 11 901–11 907, 2023.
- [303] D. Elser, F. Huber, and E. Gaquerel, “Mass2SMILES: deep learning based fast prediction of structures and functional groups directly from high-resolution MS/MS spectra.,” *BioRxiv*, pp. 2023–07, 2023.
- [304] R. Sutton, “The bitter lesson,” *Incomplete Ideas (blog)*, 2019.
- [305] K. McCloskey, E. A. Sigel, S. Kearnes, L. Xue, X. Tian, D. Moccia, D. Gikunju, S. Bazzaz, B. Chan, M. A. Clark, *et al.*, “Machine learning on DNA-encoded libraries: a new paradigm for hit finding,” *Journal of Medicinal Chemistry*, vol. 63, no. 16, pp. 8857–8866, 2020.
- [306] L. H. Yuen and R. M. Franzini, “Achievements, Challenges, and Opportunities in DNA-Encoded Library Research: An Academic Point of View,” *ChemBioChem*, vol. 18, no. 9, pp. 829–836, 2017.

- [307] D. Krentzel, S. L. Shorte, and C. Zimmer, “Deep learning in image-based phenotypic drug discovery,” *Trends in Cell Biology*, vol. 33, no. 7, pp. 538–554, 2023.
- [308] A. Pratapa, M. Doron, and J. C. Caicedo, “Image-based cell phenotyping with deep learning,” *Current Opinion in Chemical Biology*, vol. 65, pp. 9–17, 2021.
- [309] J. N. Spradlin, E. Zhang, and D. K. Nomura, “Reimagining druggability using chemo-proteomic platforms,” *Accounts of Chemical Research*, vol. 54, no. 7, pp. 1801–1813, 2021.
- [310] K. G. Hicks, A. A. Cluntun, H. L. Schubert, S. R. Hackett, J. A. Berg, P. G. Leonard, M. A. Ajalla Aleixo, Y. Zhou, A. J. Bott, S. R. Salvatore, *et al.*, “Protein-metabolite interactomics of carbohydrate metabolism reveal regulation of lactate dehydrogenase,” *Science*, vol. 379, no. 6636, pp. 996–1003, 2023.
- [311] R. N. Muchiri and R. B. van Breemen, “Affinity selection–mass spectrometry for the discovery of pharmacologically active compounds from combinatorial libraries and natural products,” *Journal of Mass Spectrometry*, vol. 56, no. 5, e4647, 2021.
- [312] R. Reher, A. T. Aron, P. Fajtová, P. Stincone, B. Wagner, A. I. Pérez-Lorente, C. Liu, I. Y. B. Shalom, W. Bittremieux, M. Wang, *et al.*, “Native metabolomics identifies the rivulariapeptolide family of protease inhibitors,” *Nature communications*, vol. 13, no. 1, p. 4619, 2022.