

Learning the Language of Antibody Hypervariability Through Biological Property Prediction

by

Chiho Im

S.B. Computer Science and Engineering,
Massachusetts Institute of Technology, 2022

Submitted to the Department of
Electrical Engineering and Computer Science
in partial fulfillment of the requirements for the degree of
Master of Engineering in Electrical Engineering and Computer Science
at the

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

June 2023

© Chiho Im, MMXXIII. 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.

Authored by: Chiho Im
Electrical Engineering and Computer Science
May 12, 2023

Certified by: Bonnie Berger
Simons Professor of Mathematics
Thesis Supervisor

Accepted by: Katrina LaCurts
Chair, Master of Engineering Thesis Committee

Learning the Language of Antibody Hypervariability Through Biological Property Prediction

by

Chiho Im

Submitted to the Department of
Electrical Engineering and Computer Science
on May 12, 2023, in partial fulfillment of the
requirements for the degree of
Master of Engineering in Electrical Engineering and Computer Science

Abstract

Machine learning-based protein language models (PLMs) have proven to be successful in a variety of structure and function-prediction contexts. However, foundational PLMs (those trained on the corpus of all proteins) rely on evolutionary co-conservation of protein sub-sequences, but this distributional hypothesis does not hold for antibody hypervariable regions. Consequently, methods like AlphaFold 2 have relatively weak performance on antibody sequences. In this work, we propose AbMAP (**Antibody Mutagenesis-Augmented Processing**), a new transfer learning framework that fine-tunes foundational models specifically for antibody-sequence inputs by supervising on examples of antibody structure and binding specificity. We demonstrate how our feature representations can be applied to the accurate prediction of an antibody’s local and global 3D structures, mutational effects on antigen binding specificity, as well as identification of its paratope. The scalability of AbMAP newly enables large-scale analysis of human antibody repertoires. We find that the AbMAP representations of individual repertoires have remarkable overlap, more so than can be discerned by sequence analysis. Our findings provide robust evidence in support of the hypothesis that antibody repertoires across individuals converge towards similar structural and functional coverage. We anticipate AbMAP will accelerate efficient and effective design and modeling of antibodies and expedite antibody-based therapeutics discovery.

Thesis Supervisor: Bonnie Berger
Title: Simons Professor of Mathematics

Acknowledgments

It has been a great honor working in the Bonnie Berger Lab at MIT CSAIL. Members of the Berger Lab are all truly wonderful and incredibly talented, and this work would not have been possible without the generous help and support from everyone. First and foremost, I would like to express my sincerest gratitude to my thesis supervisor, Professor Bonnie Berger. Her vision and passion towards the advancement of computational biology and science overall have truly been an inspiration to me, shaping my future goals as a scientist.

I would especially like to thank our research scientist, Rohit Singh, for providing me with invaluable guidance and support throughout my two years in the Berger Lab. His amazing mentorship and insight allowed me to grow both as a scholar and a researcher. I would also like to thank members from Sanofi for both funding our project and giving helpful feedback and comments on our work. I also want to thank Sam Sledzieski for having insightful discussions with me on each other's work.

Finally, I am truly blessed to have such an amazing group of friends and family in my life. Wonjune, your unwavering support and encouragement have helped me navigate through the ups and downs of life. Their presence is a constant reminder of the importance of nurturing relationships and the positive impact it can have on my well-being. Mom, dad, and Egun - I love you all so much.

THIS PAGE INTENTIONALLY LEFT BLANK

Contents

1	Introduction	15
2	Methods	21
2.1	Datasets	21
2.1.1	SAbDab	21
2.1.2	LIBRA-Seq	22
2.1.3	CoV-AbDab	22
2.1.4	Thera-SAbDab	22
2.1.5	Human Antibody Repertoire	23
2.1.6	Datasets for Labeled Antibody Mutants	23
2.2	Embedding Generation Using Mutational Augmentation	23
2.3	Language Model Refinement with a Multi-task Architecture	24
2.4	Fixed-length Embeddings	25
2.5	Max-Entropy Regularization	26
2.6	Alignment of Raw Antibody Sequences	27
3	Biological Property Predictions	29
3.1	Effective fine-tuning of antibody embeddings	31
3.2	Antibody Structure Prediction	33
3.3	Mutational Variation Prediction	36
3.4	Paratope Prediction	38
4	Revealing Shared Landscapes Across Antibody Repertoires	41

4.1	LIBRA-Seq Cell Line Identification	41
4.2	The Landscape of Human Antibody Repertoires	43
4.3	Predicting SARS-CoV-2 Variant Neutralization Ability	46
5	Discussion and Conclusion	49
A	Tables	53
B	Figures	55

List of Figures

3-1	Overview of AbMAP embedding generation and its architecture. A) Given an input antibody sequence, our pipeline generates an embedding that can be applied for various downstream tasks including structure/property prediction as well as antibody repertoire analysis. B) AbMAP architecture comprises a projection module that applies contrastive augmentation and reduces the dimensionality of the input foundational PLM embedding to generate a variable length embedding, and a Transformer Encoder module that creates a {structure/function}-specific fixed-length embedding.	30
3-2	Average ground truth TM-scores antibody pairs in each of the 20 cosine similarity score bins. The embeddings used for cosine distance are A) raw PLM embeddings B) raw PLM embeddings on CDRs C) mutation-adjusted raw PLM embeddings on CDRs D) AbMAP fixed-length embeddings. Higher monotonicity in the plot implies that the model has successfully captured 3D structural information.	32
3-3	A) Chart of average Spearman’s rank scores for ddG scores prediction of different models with various train/test splits. B) Chart of average overlap of top- k_1, k_2 ddG scores for different k_1, k_2 values. (train: 0.02, test: 0.98 for this chart).	37

4-1	A) 2D-PCA of LIBRA-Seq dataset. Each antibody is colored by its source (donor) ID. B) KDE plot of 2D PCA for fixed-length embeddings of antibody repertoire sequences sampled from all 9 human subjects. C) Swarm plot of the ratio of sequences for each subject where the most similar sequence was from the same subject. The hit ratios for sequence alignment are more spread apart compared to that of AbMAP’s embedding distance. D) 2D PCA plot of 547 antibodies from Thera-SAbDab. The PCA and KDE plots seem to show a general clustering at the bottom right corner on the projection space.	42
4-2	PCA plot of SAbDab Training Set Abs (gray) and CoVAbDab Abs (rest) A) using our model’s fixed length embedding, B) using the top $k = 50$ features from embeddings in A) determined with Schema. . . .	46
B-1	Average ground truth TM-scores antibody pairs in each of the 20 cosine similarity score bins using AbMAP representations with raw PLM embeddings (CDRs only) as input; i.e. no contrastive augmentation. Monotonicity is noticeably lower when using raw PLM embeddings as input to AbMAP compared to using contrastively augmented PLM embeddings.	55
B-2	Comparison of AbMAP-B and sequence alignment at different sequence identity thresholds for template search in antibody structure prediction. Columns from left to right (TM-Score Chain H, TM-Score Chain L, RMSD Chain H, RMSD Chain L). Rows from top to bottom (Sequence identity: 0.9, 0.8, 0.7, 0.6, 0.5). For a pair of antibodies, AbMAP-B uses the cosine distance of fixed length sequence embeddings and sequence alignment uses the pairwise global sequence alignment score as a similarity metric.	56

B-3	(Left) Chart of average Spearman’s rank scores for ddG scores prediction of Bepler & Berger based models (AbMAP and raw) with various train/test splits. (Right) Chart of average overlap of top- k_1, k_2 ddG scores for different k_1, k_2 values. This experiment was conducted for different train/test splits (train: 0.02, test: 0.98 for this chart). While the raw embeddings are comparable to augmented embeddings at higher splits, their performance notably drops at lower training set sizes.	57
B-4	KDE plot: 2D PCA of AbMAP embeddings for antibody repertoire from each human subject.	58

THIS PAGE INTENTIONALLY LEFT BLANK

List of Tables

3.1	Comparison of AbMAP-B and other models for antibody structure prediction using both RMSD (C-alpha) and TM-Score as metrics. For AbMAP-B, ProtBert, and ESM-1b, the predicted template structure was selected from a set of antibodies whose sequence identity is below 0.7. Structure prediction was conducted on both individual CDR fragments and the whole Fv chain.	34
3.2	Comparison of our model and ProtBert’s representations for antibody paratope prediction by training a separate predictor using the representations as input. The performance of our model’s representations was also compared against Parapred.	39
4.1	Results for 5-fold cross validation on neutralization prediction task using logistic regression for different models.	47
A.1	Structure prediction scores for AbMAP-P and AbMAP-E. Similar to AbMAP-B, ProtBert, and ESM-1b in Fig. 3.1, the predicted template structure was selected from a set of antibodies whose sequence identity is below 0.7. Structure prediction was conducted on both individual CDR fragments and the whole chain.	53

A.2	Test set Spearman’s rank scores computed for AbMAP-B,E when using different datasets for train and test (i.e. we evaluated the prediction of functional similarity using embeddings trained on structure data, and vice versa). While embeddings trained solely on structure data perform quite well when predicting function, embeddings trained only on function data can infer far less about antibody structure.	54
A.3	Test set mean loss and Spearman’s rank scores computed for AbMAP-B when using different number of Transformer Encoder layers in its architecture. Increase in number of layers did not necessarily show improvement in model performance.	54

Chapter 1

Introduction

In modern therapeutics, antibodies have been some of the most promising drug candidates [25]. This therapeutic success has been due to the remarkable structural diversity of antibodies, allowing them to recognize an extremely wide variety of potential targets. This diversity originates from their hypervariable regions which are critical to the functional specificity of antibodies. Experimental design of an antibody against a target of interest has historically been done by approaches like immunization or with directed evolution techniques like phage display selection [42]. However, the generation and screening process is slow and expensive. It also does not systematically explore the possible structural space, potentially leading to candidates with suboptimal binding characteristics. Furthermore, downstream considerations (e.g., developability or function-specific engineering) can not be easily accommodated [44]. There is thus a need for computational methods that can design a new antibody from scratch for a given target [27] or more efficiently refine a small set of experimentally-determined candidates. General protein structure-prediction techniques (e.g., AlphaFold 2 [19]) can struggle to predict antibody structures since the latter’s hypervariable regions (also known as the complementarity determining regions, CDRs) display evolutionarily novel structure patterns. One direction towards this has been to model the 3D structure of the entire antibody, or just its CDRs [35, 17], but these have had limited accuracy. They are also slow and require many minutes per antibody (or CDR) structure, making it infeasible to perform large-scale computational exploration or

analyze an individual’s antibody repertoire, which may contain millions of sequences.

More recently, machine learning techniques used in natural language processing have been applied to generate high-dimensional protein representations [34, 2, 11, 29]. Protein language models (PLMs) capture structural features implicitly and also enable protein-property prediction. In the context of antibodies, one approach is to simply use PLMs trained on the corpus of all proteins (e.g., ESM-1b [34]). We refer to these as “foundational” PLMs, the machine learning term for broad, general-purpose models [3]. However, the CDRs of antibodies explicitly violate the distributional hypothesis underlying foundational PLMs: sequence variability in CDRs is *not* evolutionarily constrained. Indeed, the corresponding lack of high-quality multiple sequence alignments (MSAs) for antibodies is a key reason why AlphaFold 2 works less well on them than regular proteins [19]. Therefore, another set of approaches (e.g., AntiBERTa [21], IgLM [40]) has been proposed: these train the PLM just on antibody and B-cell receptor sequence repertoires. While these approaches better address the CDRs’ hypervariability, they have the disadvantage of not being trained on the diverse corpus of all protein sequences and thus can not access the substantial insights available from foundational PLMs [2, 34, 11]. Moreover, existing approaches like AntiBERTa expend valuable explanatory power on modeling also the non-CDRs of the antibody, which are not very diverse and substantially less crucial to antibody binding-specificity. Lastly, neither set of approaches takes advantage of the 3D structures available for over 6,700 antibodies in the PDB. Here, we argue that a more effective approach is to combine the strengths of the two approaches. We present a transfer learning approach that starts with a foundational PLM but adapts it for improved accuracy on the hypervariable regions by training on antibody-specific corpora. Such training lets us take advantage of available antibody structures as well as high-throughput single-cell assays of antibody binding specificity. Moreover, this approach can also easily be ported to new foundational PLMs as they are introduced (e.g., ESM-2 [24] instead of ESM-1b).

We introduce AbMAP (Antibody Mutagenesis-Augmented Processing), a scalable transfer-learning framework that is applicable to any foundational PLM, and unlocks

greater accuracy in the prediction of an antibody’s structure and its key biochemical properties. Our broad conceptual advance is to address the weakness of foundational PLMs on antibody hypervariable regions by a supervised learning approach that is trained on antibody structure and binding-specificity profiles. Specifically, we introduce three key advances: a) maximally leverage available data by focusing the learning task only on antibody hypervariable regions; b) a contrastive augmentation approach to refine the baseline PLM’s hypervariable region embeddings so that they better capture antibody structure and function; and c) a multi-task supervised learning formulation that considers antibody protein structure as well binding specificity to supervise the representation. Since the function of an antibody is determined primarily by its hypervariable region, we focus on modeling this region (and its immediate framework neighborhood) rather than the full sequence. Furthermore, decades of research has culminated in the identification of reliable CDR signatures in antibody structures, making it computationally easy to identify such regions in an antibody sequence [28, 49]. A model like AntiBERTa [21] may expend significant portion of its capacity on attempting to re-discover such signatures during training. In contrast, we can better leverage the available data by focusing our model’s capacity on just the residues within antibody CDRs (we include two flanking residues on each side of a CDR to also capture some framework information).

Our contrastive augmentation approach is designed to hone in on the subspace of foundational PLM embeddings most relevant to an antibody. Consider the embedding of a CDR residue from a foundational PLM: it captures information about the residue and its overall context. However, this context was learned from the corpus of all proteins while the hypervariability in an antibody CDR implies a different distributional context. We therefore generate new sequences by *in silico* mutagenesis in the CDRs of the original sequence and obtain foundational PLM embeddings for these mutants. In photo-editing, the “contrast” operation increases the intensity of color differences. Analogously, to intensify the relative PLM contribution of the CDR residues, we subtract the mean of mutated embeddings from the original wild-type embedding. By subtracting away the non-CDR context, we accentuate the CDR-specific context and

the contribution of the original residues.

We apply AbMAP representations to a number of downstream tasks: identifying structural templates, predicting the binding energy changes ($\Delta\Delta G$) resulting from mutations, and identifying the antibody’s paratope. We formulate structure prediction with AbMAP as a *template-search* task where, for a query antibody, we search a template database of AbMAP antibody embeddings. Even without template refinement, our model is able to outperform many of the state-of-the-art structure prediction techniques, both general (e.g., AlphaFold 2 [19]) and antibody-specific (e.g., DeepAb [35]). It especially excels on the prediction of individual CDR structures, which is critical for accurate design. For $\Delta\Delta G$ prediction, we compare the performance of foundational PLM embeddings against their AbMAP-adjusted embeddings. Using AbMAP improves on prediction accuracy overall, achieving especially high precision in its top hits. On paratope prediction, AbMAP compares favorably with foundational PLMs as well ParaPred [23], despite having a smaller representation with fewer degrees of freedom. Finally, we also applied AbMAP to identify antibodies that can neutralize more than one SARS-CoV-2 variant; its predictions are substantially more accurate than those gleaned from ProtBert or ESM-1b directly.

AbMAP unlocks a deeper understanding of the diversity and similarity in antibody repertoires across individuals. An advantage of AbMAP’s CDR-focused representation is that isotype switching in the human body (where a heavy chain’s constant region is replaced while preserving its hypervariable region) makes such a representation more appropriate for characterizing antigen specificity in a repertoire. Analyzing Briney et al.’s data on the antibody repertoires of multiple individuals, we observe substantial diversity within each individual’s repertoire, as has been previously reported [4]. However, we find that repertoires across individuals are remarkably similar in the AbMAP embedding space, in contrast to marked sequence-level variations across individuals. This suggests that, despite sequence diversity, similar binding profiles are being activated in each individual. We also document that antibodies that have entered clinical trials appear in a specific region of the embedding space, one that also corresponds to a high density of native human antibodies. Thus, our

method can help evaluate candidate antibody sequences for druggability before expensive *in vitro* and pre-clinical trials are performed.

THIS PAGE INTENTIONALLY LEFT BLANK

Chapter 2

Methods

2.1 Datasets

2.1.1 SAbDab

SAbDab (Structural Antibody Database) [10] is a database of antibody structures where each structure is paired with various metadata including the heavy and light chains of antibodies in the PDB complex. Limiting ourselves to PDB entries where both the heavy and light chains were available, we obtained 3,785 pairs of heavy and light chain antibody sequences and PDB structures from SAbDab. We divided them into 3,000 pairs for train and 785 pairs for evaluation. From the 3,000 antibodies selected for training, we randomly sampled 100,000 pairs of antibodies and computed the pairwise similarity scores using their PDB structural data. For each of the sampled 100,000 SAbDab antibody pairs, we computed the TM-Scores of the pairs of antibodies (separately for heavy and light chains), and used them as the ground-truth label when supervised against the predicted similarity score from our model. We used TM-Score as the similarity metric as it is length-independent and is a standard metric used in well-known protein structure prediction tasks such as CASP [51].

2.1.2 LIBRA-Seq

To train the model concurrently on the functionality of antibodies in addition to their structural properties, we used LIBRA-seq, which maps 4,644 B-cell receptor sequences to their antigen specificity [38]. The antibody sequences were divided into 3,715 and 929 antibodies for train and evaluation, respectively. The 4,644 antibody sequences, each with a heavy and a light chain, were paired with three scores that indicate the binding specificity to the surface proteins from three different antigens: BG505, CZA97, and H1-A/New Caledonia/20/1999 [15, 36, 47, 33]. For each antibody, this set of three scores (each standardized to $\mu = 0, \sigma = 1$) serves as its binding specificity vector. For a pair of antibodies in the LIBRA-seq dataset, we scored their functional similarity as the dot-product of their vector representations. From the 3,715 antibodies for training, we randomly sampled 100,000 pairs of antibodies and computed these pairwise similarity scores. Similar to the structural similarity prediction mentioned above, these scores were then used as ground-truth labels for the supervised training of our model’s functional similarity prediction task.

2.1.3 CoV-AbDab

CoV-AbDab is a database of all published/patented antibody sequences (7,964 sequences in total as of July 19, 2022) that are capable of binding to coronaviruses including SARS-CoV-2, SARS-CoV-1, and MERS-CoV [30]. To evaluate our trained model, we selected from these the 2,077 sequences that bind to wild-type SARS-CoV-2. On these, we assessed the model’s ability to predict if the antibody could also neutralize at least one SARS-CoV-2 variant (alpha, beta, omicron, etc.), using only the antibody sequence information.

2.1.4 Thera-SAbDab

Thera-SAbDab is a set of antibodies (547 in total) collected from the PDB that contain immunotherapeutic variable domain sequences [31]. Each of these sequences were labeled with data such as the highest clinical trial passed (Phase-I, Phase-II, etc.).

We used our language model to compare the distribution of therapeutic antibodies against the human antibody repertoire.

2.1.5 Human Antibody Repertoire

The Briney et al. (2019) dataset contains over 3 billion B-Cell antibody sequence reads across 9 human subjects [4]. For our experiments, we sampled non-unique 100,000 antibody sequence reads from each subject, where we labelled each unique antibody with the read count within each sample population.

2.1.6 Datasets for Labeled Antibody Mutants

The dataset provided in Desautels et al. (2020) was used to validate our antibody language model through ddG score predictions [8]. They used a machine learning model to search the mutational combinatorial space of the m396 antibody, and calculated the binding propensity of these mutants against the SARS-CoV-2 spike protein’s receptor binding domain. They applied five standard software packages to score the physicochemical characteristics of the mutant’s binding, including STATIUM, FoldX and Rosetta [5, 37, 20]. To validate our model, we used these pre-computed scores to formulate regression-based property prediction tasks, with only the mutant sequence as input.

2.2 Embedding Generation Using Mutational Augmentation

Given a heavy or a light chain sequence of length n of an antibody, we created its *CDR-specific* embedding by performing the following set of operations. We first create an $\mathbb{R}^{n \times d}$ embedding using an foundational protein language model such as the Bepler & Berger [2], ESM-1b [34], or ProtBert [11], where $d = 2200$ (Bepler & Berger), 1280 (ESM-1b), or 1024 (ProtBert). We note that Bepler & Berger has 6165 dimensions across three layers but we used only the last layer’s weights. Then, using

ANARCI [9], we number the amino acid residues in the input sequence using the Chothia scheme [1]. Then, we identify parts of the embedding that correspond to the CDRs based on the numbering scheme. Next, we generate k (here, $k = 100$) new antibody sequences through *in silico* mutagenesis of the original sequence, where the mutation is performed (sampling from a uniform distribution over all amino acids) with certain probability (here, 0.5) for each residue in the identified CDRs. This procedure is repeated k times for the original sequence, generating k new sequences which are then embedded using a foundational PLM as in step 1. Then, each of k new embeddings is subtracted from the embedding of the original sequence, and these adjusted embeddings are averaged. Finally, we extract and concatenate parts of the averaged difference embedding that belong to the CDRs as determined previously by ANARCI, creating an $\mathbb{R}^{n' \times d}$ *CDR-specific* embedding. The length n' usually ranges between 20 to 30.

2.3 Language Model Refinement with a Multi-task Architecture

Once we generate a CDR-specific embedding for a given antibody sequence, we use it as input to our model which outputs a fixed-length feature that can be used for further downstream task-specific similarity prediction. We curated two datasets for pairs of antibodies, one labeled with their 3D structure similarity, and the other with functional similarity (with regards to antigen binding profiles) to train the model. The overall pipeline as well as the diagram of our antibody language model is shown in Fig. 3-1. Our model consists of an MLP projection module whose parameters are shared across the training of each task (structure and function). This projection module reduces the dimension of the input embedding ($\mathbb{R}^{n' \times d} \rightarrow \mathbb{R}^{n' \times d'}, d' < d$) so that the language model is forced to retain as much information about the antibody sequence as possible while trying to execute downstream prediction tasks. For example, when using the Bepler & Berger language model as the foundational PLM, the input CDR-

specific embedding is $\mathbb{R}^{n' \times 2200}$ and our model’s projection module outputs a $\mathbb{R}^{n' \times 256}$ embedding.

Then, the embedding with reduced dimensions, outputted by the projection module is fed into two separate PyTorch (v1.11.0) Transformer Encoder modules, one each for a downstream similarity score prediction task (structure and function) [45]. We add sinusoidal positional encodings to the input embeddings in order to inject information about the relative positions of amino acid residues in each sequence. The representations for each residue in the embedding outputted by this Transformer Encoder module now incorporate antibody-related structural and functional information by leveraging attention from other residues in the CDR.

In a multi-task training framework such as our setting, it is important that the calculated losses from each task is combined carefully so that the shared parameters in our model is correctly optimized. Rather than using a set ratio to weigh the two MSE losses, we assigned a new learnable parameter α for weighing the losses. The overall MSE loss L_{MSE} is calculated as a weighted sum of the losses calculated for structural similarity prediction ($L_{structure}$) and functional similarity prediction ($L_{function}$):

$$L_{MSE} = \alpha L_{structure} + \frac{1}{\alpha} L_{function} \quad (2.1)$$

where α is updated iteratively through gradient calculations.

2.4 Fixed-length Embeddings

In addition to the per-residue embeddings (variable length) focusing on the CDRs of an antibody’s sequence, our model can also generate fixed-length embeddings by performing pooling operations on the variable length embedding outputted by the Transformer Encoder along the sequence dimension. Specifically, we use the Log-SumExp (LSE) operation, defined as:

$$LSE(\mathbf{x}) = \log \left[\sum_{j=1}^n \exp(x_j) \right] \quad (2.2)$$

where $\mathbf{x} = (x_1, \dots, x_n)$, to compute a smooth maximum over the sequence embedding. We also use mean pooling to generate another fixed length embedding of same length. The two are then concatenated into a single representation. Overall, the variable length feature in space $\mathbb{R}^{n' \times d'}$ is transformed into a fixed length feature in space \mathbb{R}^q , where $q = 2d'$ (here, $q = 512$). This fixed-length vector for each input antibody sequence is then used for similarity score prediction using the cosine distance metric.

2.5 Max-Entropy Regularization

In order to prevent the learned representations being overfit to the datasets used for supervision, we include a regularization loss adapted from Shannon entropy [39]:

$$H(X) = - \sum_{i=1}^n P(x_i) \log P(x_i) \quad (2.3)$$

The regularization is applied to the antibody’s L^2 -normalized fixed-length feature of length q , which is outputted by the task-specific Transformer Encoder module following the nonlinear projection module. The square of each entry in the fixed-length feature is treated as the probability of an arbitrary discrete random variable X , corresponding to $P(x_i)$ in the above equation. To induce a regularizing effect, we want the squared entries in the feature to form a uniform distribution. For each L^2 -normalized feature, the squared entries are non-negative and sum to 1, like a probability distribution. We therefore set the regularization loss in a max-entropy formulation as:

$$L_{reg}(u) = \sum_{i=1}^q u_i^2 \log u_i^2 \quad (2.4)$$

where $u = [u_1, u_2, \dots, u_q]$ is a task specific feature of length q .

Therefore, with a regularization parameter λ , the total loss computed during a single feed-forward step is:

$$L_{total} = L_{MSE} + \lambda L_{reg} \quad (2.5)$$

where λ was empirically determined as 0.0005, and was used throughout our training.

2.6 Alignment of Raw Antibody Sequences

We used the Needleman-Wunsch algorithm [26] to compute the per-residue similarity of two protein (antibody) sequences using the implementation of the algorithm in BioPython (v1.69). We set the match (identical character) and mismatch scores to 1 and 0, respectively, and did not assign gap penalties.

THIS PAGE INTENTIONALLY LEFT BLANK

Chapter 3

Biological Property Predictions

For a given antibody sequence, we start from an embedding generated by a foundational PLM and refine it with AbMAP. There are three main steps in our refinement: identification of CDR regions, augmenting the foundational PLM embeddings to focus on the CDRs, and an attention-based fine-tuning of the embedding to better capture antibody structure and function. We first apply ANARCI [9], a hidden Markov model approach, to demarcate the boundaries of CDR regions. ANARCI leverages well-known canonical patterns of antibody structure (e.g., a disulfide bond that spans CDR-H1 and CDR-H2 [28], a Tryptophan residue located immediately after CDR-L1 [49] etc.) to identify CDR regions with high confidence. We extend each ANARCI-reported segment by two residues on either side, allowing us to compensate for potential errors in ANARCI’s inference and also include the bounding framework residues. A similar design choice was made by Liberis et al. in Parapred [23]. Next, we apply a procedure that we term “contrastive augmentation”: we perform in-silico mutagenesis in the CDR regions by randomly replacing a CDR residue with another amino acid and generate foundational PLM embeddings for each mutant. We then compute the augmented embedding as the difference between the embeddings of the original sequence and the average over mutants. Our augmentation aims to subtract away the sub-space of the embedding that does *not* correspond to the CDR regions and, akin to masked language modeling, it highlights the contribution to embedding of a specific amino acid by contrasting it against potential replacements. We

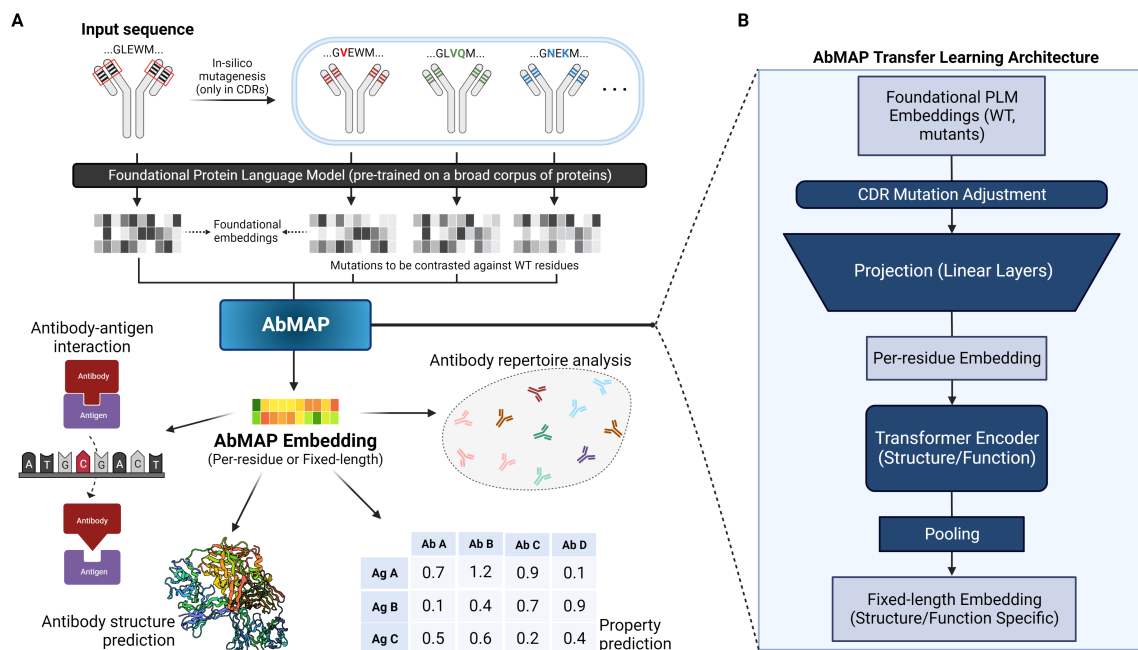


Figure 3-1: **Overview of AbMAP embedding generation and its architecture.** A) Given an input antibody sequence, our pipeline generates an embedding that can be applied for various downstream tasks including structure/function prediction as well as antibody repertoire analysis. B) AbMAP architecture comprises a projection module that applies contrastive augmentation and reduces the dimensionality of the input foundational PLM embedding to generate a variable length embedding, and a Transformer Encoder module that creates a {structure/function}-specific fixed-length embedding.

then optimize this augmented embedding with a Siamese neural network architecture with a single transformer layer that takes pairs of antibody sequences as inputs and seeks a final representation for each antibody where Euclidean distances capture structural/functional information (**Methods**). Our approach can be applied to any foundational PLM. Here, we have applied it on Bepler & Berger, ESM-1b [34], and ProtBert [11] foundational PLMs, producing AbMAP representations that we denote as AbMAP-(B, E, P) respectively.

3.1 Effective fine-tuning of antibody embeddings

We first assessed the effectiveness of our refinement and fine-tuning approach. After a random selection of 785 antibodies with available structures that were not in AbMAP’s training or validation set, we generated from them 10,000 random pairs and assessed how pairwise structural similarity correlated with representation similarity. We chose to evaluate structural similarity over the whole Fv, rather than just the CDR fragments, since we thought it to be a stricter test of the thesis that CDR-specific drivers are the primary determinants of overall structural variability across Fv structures. We evaluated representations from the baseline foundational PLM and each step of our refinement scheme: a) baseline (i.e., foundational) PLM representation for the whole protein, b) baseline PLM just for CDRs, c) with contrastive augmentation on CDRs and, d) also with the supervised transformer layer (residue-wise averaging for a-c). For each representation, we grouped the antibody pairs into 20 bins by cosine similarity and, in each bin, computed the distribution of TM-scores of structural similarity between the pairs.

We assessed the embedding–structure relationship on consistency (measured as the Spearman rank correlation between average TM-score and cosine similarity across bins), as well as discriminative power (measured as the TM-score difference between the first and last bin). We show the results of AbMAP-B, (i.e. our method applied to the Bepler & Berger embedding) on heavy chain antibodies in Fig. 3-2. While the baseline foundational PLM is quite powerful on its own, it has somewhat limited

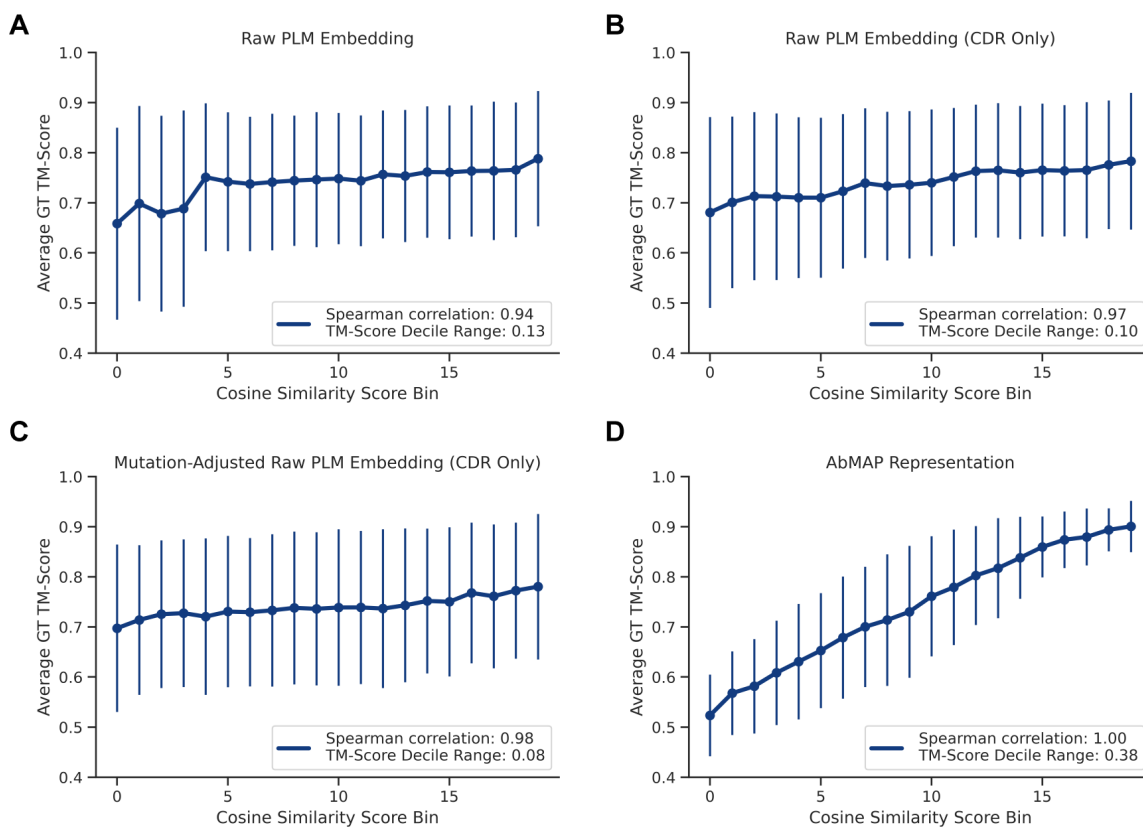


Figure 3-2: Average ground truth TM-scores antibody pairs in each of the 20 cosine similarity score bins. The embeddings used for cosine distance are A) raw PLM embeddings B) raw PLM embeddings on CDRs C) mutation-adjusted raw PLM embeddings on CDRs D) AbMAP fixed-length embeddings. Higher monotonicity in the plot implies that the model has successfully captured 3D structural information.

consistency, especially on pairs with lower representation similarity. The CDR-specific embedding and contrastive augmentation improve the consistency (with Spearman correlation increasing from 0.94 to 0.98), while the final supervised layer is crucial to achieving strong discriminative power, with the TM-score separation between the first and last bin increasing by 375%. Thus, applying a CDR-specific representation, augmenting it contrastively, and fine-tuning it in a supervised setting can accentuate antibody structural features more comprehensively and accurately than raw whole-protein PLM representations.

3.2 Antibody Structure Prediction

We approached structure prediction with AbMAP as a template-matching task: we searched through a database of antibody templates to find the example that we expect to be closest in structure to a query antibody. The template can later be refined (e.g., the Evoformer module [19]), a direction we hope to explore in future work (**Discussion**). For benchmarking purposes, the AbMAP template quality is itself an indication of the method’s ability to recapitulate structure. We constructed the template database from the set of SAbDab structures during AbMAP’s training; re-using these examples allowed us to evaluate on the remaining structures in SABDab. We first applied CD-HIT [22, 14] to remove any template entries with greater than 70% sequence identity of test entries (our results are robust to this threshold as shown in Fig. B-2). Using fixed-length AbMAP representations, we obtained the k (here, $k = 10$) templates closest to the query embedding by Euclidean distance (**Methods**). The medoid of these k representations (in the Euclidean space) was reported as the matching template. Choosing the medoid instead of the very closest template offers some robustness to variability in embedding quality across queries and templates. As the underlying foundational PLMs improve, we expect this k could be lowered. We applied this process to generate templates from each foundational PLM as well as its AbMAP variants. For the former, the fixed-length embedding was obtained by taking the mean over the residues of the sequence’s PLM embedding. In addition to

Model	Metric	Chain	CDR 1	CDR 2	CDR 3	Whole Fv
AbMAP-B	TM-Score	H	0.62 \pm 0.007	0.67 \pm 0.006	0.54 \pm 0.007	0.86 \pm 0.006
ProtBert	TM-Score	H	0.33 \pm 0.003	0.31 \pm 0.003	0.28 \pm 0.003	0.69 \pm 0.005
ESM-1b	TM-Score	H	0.49 \pm 0.007	0.51 \pm 0.006	0.44 \pm 0.007	0.79 \pm 0.004
DeepAb	TM-Score	H	0.51 \pm 0.008	0.55 \pm 0.008	0.24 \pm 0.004	0.54 \pm 0.005
OmegaFold	TM-Score	H	0.61 \pm 0.007	0.67 \pm 0.006	0.34 \pm 0.005	0.79 \pm 0.005
AlphaFold2	TM-Score	H	0.28 \pm 0.003	0.30 \pm 0.003	0.30 \pm 0.004	0.65 \pm 0.004
AbMAP-B	TM-Score	L	0.62 \pm 0.007	0.76 \pm 0.007	0.65 \pm 0.007	0.89 \pm 0.004
ProtBert	TM-Score	L	0.34 \pm 0.003	0.60 \pm 0.007	0.52 \pm 0.010	0.80 \pm 0.005
ESM-1b	TM-Score	L	0.41 \pm 0.005	0.61 \pm 0.007	0.39 \pm 0.005	0.82 \pm 0.004
DeepAb	TM-Score	L	0.40 \pm 0.006	0.66 \pm 0.009	0.38 \pm 0.006	0.52 \pm 0.005
OmegaFold	TM-Score	L	0.63 \pm 0.006	0.69 \pm 0.006	0.58 \pm 0.006	0.83 \pm 0.005
AlphaFold2	TM-Score	L	0.27 \pm 0.003	0.36 \pm 0.007	0.31 \pm 0.003	0.58 \pm 0.004
AbMAP-B	RMSD	H	0.43 \pm 0.016	0.38 \pm 0.013	0.43 \pm 0.025	2.11 \pm 0.082
ProtBert	RMSD	H	1.40 \pm 0.030	1.21 \pm 0.028	0.64 \pm 0.031	3.07 \pm 0.078
ESM-1b	RMSD	H	0.54 \pm 0.017	0.76 \pm 0.022	0.50 \pm 0.023	2.69 \pm 0.063
DeepAb	RMSD	H	0.56 \pm 0.016	0.55 \pm 0.015	0.81 \pm 0.031	0.72 \pm 0.028
OmegaFold	RMSD	H	0.35 \pm 0.013	0.37 \pm 0.013	0.75 \pm 0.035	2.39 \pm 0.067
AlphaFold2	RMSD	H	0.70 \pm 0.032	0.37 \pm 0.030	1.24 \pm 0.063	4.40 \pm 0.077
AbMAP-B	RMSD	L	0.41 \pm 0.015	0.18 \pm 0.006	0.44 \pm 0.022	1.42 \pm 0.044
ProtBert	RMSD	L	1.39 \pm 0.041	0.15 \pm 0.008	0.64 \pm 0.044	0.76 \pm 0.012
ESM-1b	RMSD	L	0.68 \pm 0.022	0.22 \pm 0.008	0.65 \pm 0.028	2.38 \pm 0.062
DeepAb	RMSD	L	0.72 \pm 0.021	0.32 \pm 0.011	0.84 \pm 0.030	1.16 \pm 0.084
OmegaFold	RMSD	L	0.40 \pm 0.016	0.17 \pm 0.007	0.46 \pm 0.017	2.31 \pm 0.076
AlphaFold2	RMSD	L	1.13 \pm 0.046	0.24 \pm 0.024	0.64 \pm 0.033	4.79 \pm 0.084

Table 3.1: Comparison of AbMAP-B and other models for antibody structure prediction using both RMSD (C-alpha) and TM-Score as metrics. For AbMAP-B, ProtBert, and ESM-1b, the predicted template structure was selected from a set of antibodies whose sequence identity is below 0.7. Structure prediction was conducted on both individual CDR fragments and the whole Fv chain.

the foundational PLMs, we compared the performance of AbMAP against DeepAb, OmegaFold, and AlphaFold, some of the state-of-the-art deep learning-based methods for antibody structure prediction [35, 48, 19]. To quantify the similarity between the predicted and ground truth structures, we computed the TM-scores and RMSD (Root Mean Square Deviation) between the predicted and ground-truth Fv structures. While we consider both metrics, we believe TM-score is more appropriate since it is robust to variations in protein size, unlike RMSD. Since an antibody’s CDRs play crucial roles in its function, it is important that structure prediction models achieve high local accuracy on CDR structures. Accordingly, we also evaluated the methods on their predictions of specific CDRs. For individual CDR structure prediction, we prepared our embeddings separately, supervising with similarity scores for each CDR (H1-3, L1-3) instead of the whole structure.

Overall, as shown in Fig. 3.1, AbMAP (we show its -B variant here, with others reported in Table A.1) is able to achieve high accuracy in structure prediction, despite no further refinement of the reported template. Compared to their respective foundational PLMs, each of the corresponding AbMAP variants performed substantially better. Overall, AbMAP-B performed better than other variants, possibly because the underlying foundation model is trained on both sequence and structure. Notably, AbMAP also improved over dedicated structure-prediction methods broadly. In particular, AlphaFold 2 performed substantially worse than others. This is consistent with reports of it underperforming on targets (like antibodies) where high-quality multiple sequence alignments are not available [50]. AbMAP was also competitive with the language model-based OmegaFold, outperforming it on the TM-score metric and being roughly equivalent with it on the RMSD metric. On both metrics, the relative performance of AbMAP was especially strong on the crucial CDR-H3 region, suggesting that our generated embeddings may contain rich information about the CDRs that contribute the most to the antibody’s activity and specificity.

While the structures predicted by AbMAP can certainly be used for downstream tasks, we recommend directly using the PLM itself for downstream property-prediction tasks like paratope or ddG prediction. Historically, explicit elucidation of the atomic

coordinates of the antibody structure has been viewed as a prerequisite for such downstream tasks since they are informed by the structure’s physicochemical properties. However, we believe that the implicitness of representations encoded in a PLM like AbMAP allows a task-specific neural network greater power in marginalizing over unknowns and uncertainties in the structure (e.g., conformational flexibility); this implicit richness is lost when resorting to a single, fixed 3-D structure. Moreover, AbMAP offers the choice between a fixed-length embedding for property-prediction tasks and a per-residue variable-length embedding for tasks like *in-silico* mutagenesis; either embedding may be used as desired by the user.

3.3 Mutational Variation Prediction

A key application of computational antibody modeling is low-N antibody design and optimization: the task is to computationally extrapolate the effect of combinatorial mutations starting from a small training set of antibodies to a broad set of antibody candidates, using the results to guide the next round of assays. PLM-based *in silico* mutagenesis can play a key role in speeding up the design and development of antibody-based therapeutics. We assessed the generalization performance of AbMAP in estimating the binding efficacy of m396 mutants to SARS-CoV-2. The original wild-type variant of m396 targets the receptor binding domain (RBD) of the SARS-CoV-1 spike protein. During the pandemic, Desautels et al. sought to adapt this antibody to target the RBD of the SARS-CoV-2 spike protein. They generated 90,000 mutant *in silico* and estimated each mutant’s binding efficacy by computing its ddG scores from five energy functions (FoldX Whole, FoldX Interface Only, Statium, Rosetta Flex, Rosetta Total Energy) [37, 20, 5]. The effort required substantial high-performance computing resources from the Lawrence Livermore National Laboratory, and needed over 200,000 hours of CPU time [8]. We sought to predict the ddG scores for these set of mutants after training on as little as 0.5% of the examples. We evaluated two prediction architectures: i) using AbMAP’s variable-length embedding as input to a transformer layer followed by a two-layer feed-forward network (averaged

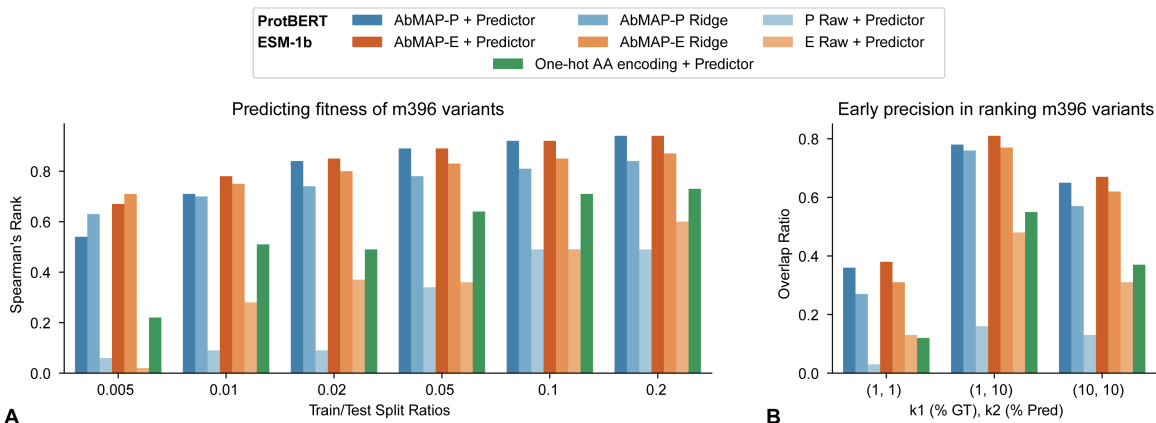


Figure 3-3: A) Chart of average Spearman's rank scores for ddG scores prediction of different models with various train/test splits. B) Chart of average overlap of top- k_1, k_2 ddG scores for different k_1, k_2 values. (train: 0.02, test: 0.98 for this chart).

over residues), and ii) using AbMAP's fixed-length embedding as input to a ridge regression. The target variable in both cases was Desautels et al.'s binding efficacy score for the mutant. The first architecture offers greater explanatory capacity while the second architecture can be applied even with very limited training data. We trained these architectures using embeddings from AbMAP-B/E/P as well as their corresponding foundational PLMs. We also used a simple baseline embedding that uses a one-hot encoding of the amino acids. The intuition behind this baseline is that, given sufficiently many examples, the one-hot encoding can be leveraged effectively by the downstream predictor but its usefulness would diminish when fewer training examples are available. We trained the different models for 100 epochs each, and varied the train/test split ratio to examine how the performance degrades as fewer training examples are provided.

In our evaluations, we assessed both the overall accuracy of AbMAP-based predictions as well as its ability to recapitulate the top ground-truth hits. For the overall analysis, we computed the Spearman rank correlation between predicted and ground-truth scores, averaging these correlations over the five energy-function categories. As shown in Fig. 3-3a, with just 20% of the examples (i.e., train/test split of 0.2), AbMAP-E and AbMAP-P's representations both achieve Spearman rank correlation of 0.94, indicating that AbMAP can effectively generalize from a limited training set,

thus reducing experimental and computational expenses. In contrast, the raw PLMs perform substantially worse. Indeed, for the largest training set size (0.2 split), the lower-dimensional one-hot encoding performs better than the foundational PLMs. The performance comparison for AbMAP-B and the foundational Bepler & Berger PLM are shown in Fig. B-3. Furthermore, the performance of AbMAP embeddings is more robust to smaller training set sizes than the baseline PLMs. The relative outperformance of AbMAP becomes more substantial as the number of training examples decreases. With just 0.5% of the examples (i.e., train/test split of 0.005), AbMAP is able to achieve high accuracy (Spearman rank correlations of 0.71 and 0.63 with the AbMAP-E and -P models, respectively). Notably, the ridge regression-based formulation starts to outperform the more complex model as the number of training examples decrease, as would be expected. High accuracy in such few-shot settings is crucial since they enable a small set of experimentally-assayed binding specificity/strength measurements to be extrapolated more broadly.

We also examined the early precision of our model, i.e., its ability to correctly prioritize mutants with high ground-truth scores. We checked how many sequences in the top $k_1\%$ ($k_1 = 1, 10$) of ground-truth scores overlapped with those in the top $k_2\%$ ($k_2 = 1, 10$) of the model’s predicted scores. Even at high cutoffs, when comparing the top 1% of predicted and ground truth scores by magnitude, there is 38% and 36% overlap when using AbMAP-E and AbMAP-P embeddings, respectively (Fig. 3-3b). Altogether, AbMAP-P/E are able to robustly predict mutational performance, both broadly and for the top hits. Crucially, they are able to operate much more effectively with limited training data, compared to the foundational PLM embeddings or one-hot encodings.

3.4 Paratope Prediction

An important antibody sub-structure is the paratope— the region that recognizes and binds to an antigen. In particular, each residue on the antibody backbone can be assigned a binary label indicating if it belongs to the antibody’s paratope. We applied

Embedding	Dimension	Predictor	Accuracy	AUPRC	AUROC
AbMAP-B	256	1-layer transformer	0.798	0.653	0.830
ProtBert	1024	1-layer transformer	0.797	0.646	0.842
One-hot + Physicochemical	28	Parapred	0.621	0.639	0.838

Table 3.2: Comparison of our model and ProtBert’s representations for antibody paratope prediction by training a separate predictor using the representations as input. The performance of our model’s representations was also compared against Parapred.

AbMAP-B to the paratope prediction task, comparing it against a dedicated machine learning method, Parapred [23], as well the ProtBert foundational PLM. Acquiring all SABDab heavy chain entries with at least one CDR in contact with an antigen, we labeled a residue as part of the paratope if it was within 5Å of the antigen. To avoid data snooping, we re-used AbMAP’s training set for this task-specific training (1195 entries), and evaluated on the remaining test set of entries (312 entries).

For paratope prediction, we specified a simple architecture that uses the per-residue, variable-length representation of AbMAP: a single transformer layer followed by two linear layers. Predictions from ProtBert were made using the same architecture. Notably, the ProtBert model has more parameters because the ProtBert embedding dimension (1024) is four times as large as AbMAP’s (256). The Parapred model takes individual CDR fragment strings as input, and we applied it separately on the three heavy chain CDRs. We calculated per-residue performance and report the overall statistics in Table 3.2. AbMAP-B achieves the highest overall accuracy for per-residue paratope prediction. The performance of Parapred reported here may be a slight overestimate since we use the trained model made available by Leem et al.’s PyTorch implementation [18]; their training may have utilized some of the examples that we use here in the test set. While ProtBert has similar accuracy to AbMAP, it uses many more model parameters due to the larger embedding dimensionality.

THIS PAGE INTENTIONALLY LEFT BLANK

Chapter 4

Revealing Shared Landscapes Across Antibody Repertoires

The scalability of our approach, along with its fidelity in capturing structure and function, enables systematic analyses of large antibody repertoires. Existing approaches are ill-suited to such analyses: while structure-prediction methods can not scale to large repertoires, purely sequence similarity-based analyses [4] will not be sensitive to structural/functional similarities between antibodies with different sequences. While PLMs directly address this concern, language models that learn the full antibody’s representation (e.g., AntiBERTa [21]) may misemphasize the framework diversity at the expense of CDR diversity, with the latter being the key determinants of antibody specificity.

4.1 LIBRA-Seq Cell Line Identification

As a preliminary evaluation, we assessed AbMAP on Setliff et al.’s LIBRA-seq study that profiled B-cell receptor (BCR) binding specificity against a panel of HIV and influenza-related antigens. We wondered if AbMAP is able to recover the cell of origin of the BCR. While we had used a subset of this dataset in AbMAP’s multi-task training, we note that no cell-of-origin information was provided during training. Furthermore, we evaluated on a subset of BCRs held out during training.

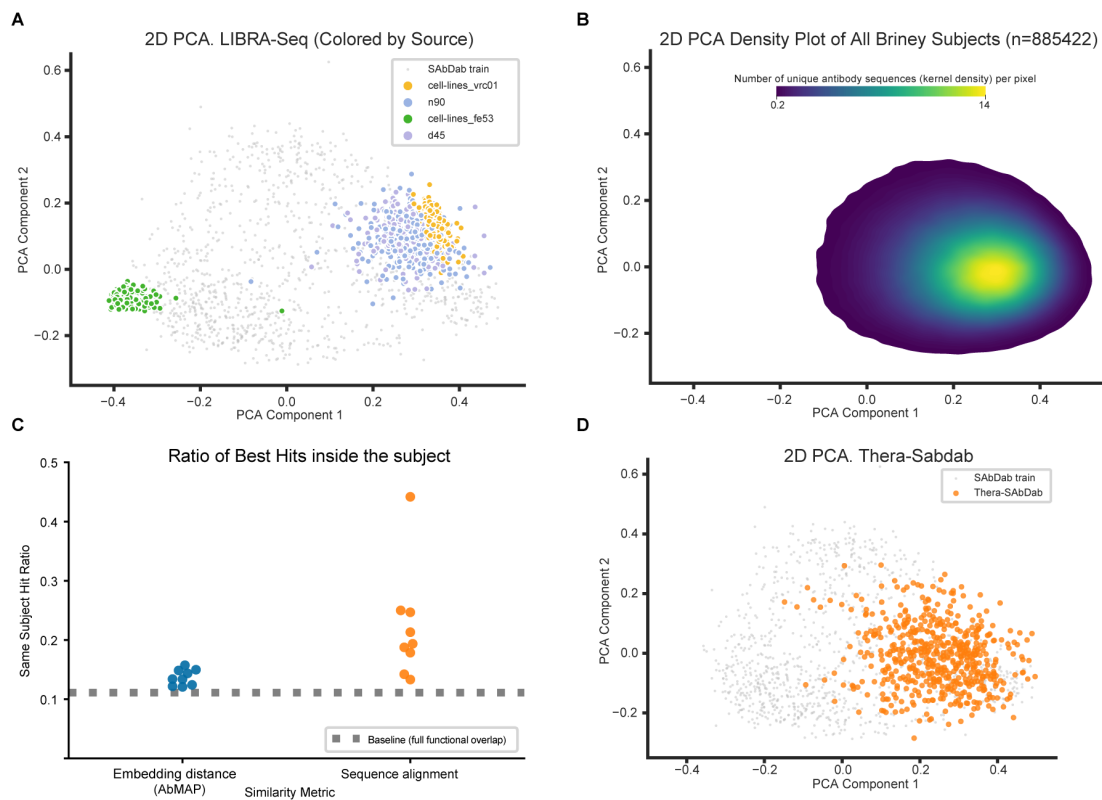


Figure 4-1: A) 2D-PCA of LIBRA-Seq dataset. Each antibody is colored by its source (donor) ID. B) KDE plot of 2D PCA for fixed-length embeddings of antibody repertoire sequences sampled from all 9 human subjects. C) Swarm plot of the ratio of sequences for each subject where the most similar sequence was from the same subject. The hit ratios for sequence alignment are more spread apart compared to that of AbMAP’s embedding distance. D) 2D PCA plot of 547 antibodies from Thera-SAbDab. The PCA and KDE plots seem to show a general clustering at the bottom right corner on the projection space.

For consistent visualization across analyses, we first standardized a 2-dimensional representation of the overall antibody space. Using the sketching algorithm Hopper [7], we selected 1,000 diverse antibodies from the SABDab training set. The farthest-first sampling approach of Hopper is motivated by the vertex k-center problem, and allows us to identify a set of antibodies that cover the entire space. We then computed a 2D reduction by mapping the fixed-length AbMAP-B representations of these antibodies to the top two principal components. In all analyses that follow, AbMAP fixed-length embeddings are reduced to 2D space by this specific mapping, enabling us to visualize different datasets (e.g., in Figures 4-1a-c) on the same axes. In most plots, we mark the 1,000 “anchor” antibodies as gray dots in the background, to help contextualize the visualization.

We computed AbMAP-B embeddings for 887 LIBRA-seq BCR sequences that were not part of the training set. As shown in Fig. 4-1a, our model representation was able to differentiate whether BCRs originated from engineered B-cell lines or real human donors. Furthermore, our model was able to discern the two different types of BCRs within the engineered cell lines: VRC01, a CD4 binding-site-directed HIV-1 bNAb (broadly neutralizing antibody), and Fe53, a bNAb recognizing the stem of group-1 influenza hemagglutinins (Fig. 4-1a). We note that there is some clonal diversity within each cell line, with cells in the line sharing the same lineage but having diversified through antigen exposure and somatic mutation [16].

4.2 The Landscape of Human Antibody Repertoires

We analyzed the human antibody datasets from the Briney et al.’s study of BCR repertoires across multiple individuals. As part of this analysis, we ask two key questions:

- Is the set of BCRs uniformly distributed over the embedding space, or does the distribution display “hotspots” of clustering? In case the latter is true, we also wondered if antibody drug candidates that have successfully passed pre-clinical evaluations were more likely to cluster in these hotspots.

- While an extraordinary sequence diversity of BCRs has been reported across individual repertoires [4], multiple analyses have suggested that these diverse repertoires converge to similar structure/function [6, 12, 32]. We wondered if such similarities across individuals would be more evident in our representation space than in the raw sequence space.

Recently, machine learning and experimental approaches have suggested a convergence of structure and function across human BCR repertoires [6, 12, 32]. Additionally, Friedensohn et al. reported a deep learning approach that shows extensive convergent selection in antibody repertoires of mice for a range of protein antigens and immunization conditions [13]. AbMAP-based embeddings enable us to go beyond previous approaches in providing a large-scale, principled approach to systematically quantifying the structural/functional convergence across individual repertoires. While previous work has primarily been focused on convergence in sequence fragments, our embedding-based approach encompasses convergence also in “paratope structural signatures” [32]. We therefore hypothesized that incorporating AbMAP representations in repertoire analysis could more accurately reveal the extent of structural/functional convergence across individuals.

We acquired 100,000 randomly-chosen BCR sequences from 9 individuals each, filtering out identical sequences from an individual. Because of assay limitations in the original study, many of the sequences were truncated on one or both ends. This would be a challenge for language models that need to consider the full antibody; however, our focus on just the CDRs enabled us to recover embeddings for the vast majority ($n = 885,422$) of these BCRs. We applied AbMAP-B on these sequences and visualized their 2D reduction as described above. We found the distribution to be highly clustered, with a kernel density estimator of the distribution being unimodal (Fig. 4-1b). Notably, the human cell lines from the LIBRA-seq data overlap very well with the Briney dataset. Interestingly, while the VRC01 cell line falls within the high-occupancy region, the Fe53 cell line is well out of it. While both cell lines are engineered human B cell lines (Ramos) [46], they express distinct BCRs that neutralize different antigens (VRC01 antibody neutralizes HIV-1 while Fe53 antibody

neutralizes influenza).

We next assessed if the antibody repertoires across individuals were similar. In the original study, Briney et al. had observed substantial cross-individual diversity. While this makes sense given the vast space of possible antibody sequences, it is also somewhat puzzling—ultimately, most individuals need antibody-based protection from similar antigens (e.g., the flu, environmental stressors etc.). Indeed, when we visualized each individual’s repertoire in our embedding space, we found that the distributions looked remarkably similar (Fig. B-4), suggesting that each repertoire had similar structural/function coverage. For a more systematic assessment of cross-individual overlaps, we sampled 5,000 sequences from each subject (i.e., 45,000 sequences in total) and performed all-vs.-all pairwise comparisons, using either the raw sequence¹ or the AbMAP-B representation (**Methods**). From these pairwise comparisons, we obtained the nearest neighbor of each antibody across all individuals and computed the frequency with which this neighbor hails from the same individual. If the per-individual repertoires are independent and identically distributed (i.i.d.) samples from the same underlying distribution, the fraction of cases where the nearest neighbor is from the same individual should be $\frac{1}{9}$ (=11.1%). When using AbMAP-B embeddings to compute similarities, we found that this to nearly be the case (Fig. 4-1c): per-individual fractions averaged 0.14 ± 0.013 , with very little variation across individuals. In contrast, when using the sequence similarity metric, the average per-individual fraction was substantially higher (0.22 ± 0.087), with much greater variability across individuals. Thus, while per-individual repertoires seem to differ substantially when only sequence similarity is concerned, these repertoires are revealed to be much more similar in their structure/function coverage when represented by AbMAP.

We wondered if the distribution pattern of human antibody representations is therapeutically useful. Towards that, we mapped 547 antibodies from Thera-SABDab, a dataset of immunotherapeutic antibodies that have entered clinical trials [31]. Our hypothesis was that even if an antibody drug candidate is effective in vitro, it may not

¹The global alignment for sequence pairs was computed with no match or gap penalties

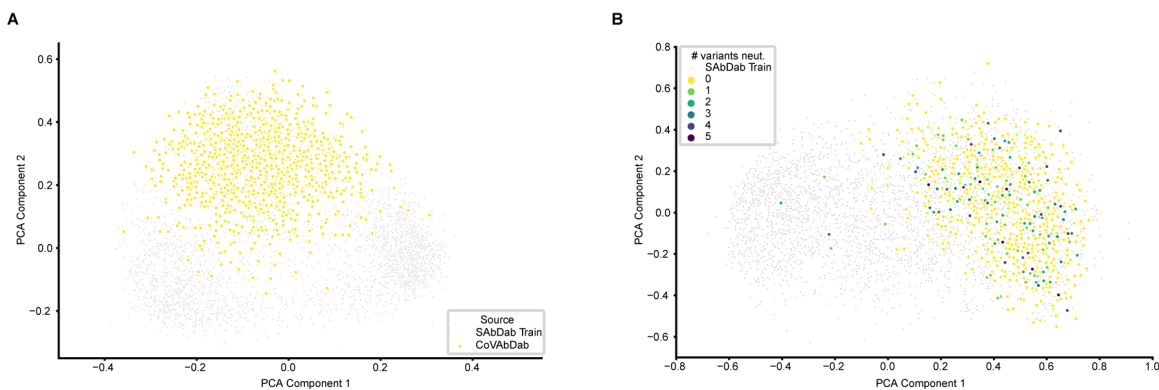


Figure 4-2: PCA plot of SAbDab Training Set Abs (gray) and CoVAbDab Abs (rest) A) using our model’s fixed length embedding, B) using the top $k = 50$ features from embeddings in A) determined with Schema.

have drug-like properties (e.g., low toxicity) unless it falls solidly within the realm of native-like antibodies. As Fig. 4-1d indicates, this does seem to be the case – the set of antibodies in Thera-SABDab are located in the high-occupancy region of the embedding space. Thus, a general drug-discovery recommendation we make is to sample candidate antibody drugs primarily from this region of the embedding space.

4.3 Predicting SARS-CoV-2 Variant Neutralization Ability

We also examined AbMAP’s capability on the prediction of antibody efficacy on neutralizing SARS-CoV-2 variants. This analysis also demonstrates how AbMAP could be applied in a low-N setting to optimize an existing panel of antibodies to address incremental mutations in the target. From CoV-AbDab, a dataset of coronavirus-binding antibodies, we obtained the set of 2,077 antibodies that were reported to neutralize the wild-type strain of SARS-CoV-2; we then computed fixed-length AbMAP-B embeddings for these. Setting aside 522 (approx 25%) randomly chosen antibodies, we plotted the rest (=1,555) on the 2D visualization described above (Fig. 4-2a). The set of wild-type neutralizing antibodies span a fairly large part of the overall space, suggesting that a substantial diversity of antibody structures is capable of neutral-

Embedding	Dimension	Mean Accuracy	Balanced Accuracy	F1 Score	AUROC
AbMAP-B	50	0.752	0.631	0.227	0.638
ProtBert	1024	0.767	0.565	0.155	0.602
ESM-1b	1280	0.775	0.550	0.169	0.555

Table 4.1: Results for 5-fold cross validation on neutralization prediction task using logistic regression for different models.

izing SARS-CoV-2 (we note that these antibodies vary in their reported viral target protein).

Seeking to further hone in on the variant neutralization capabilities of these antibodies, we applied Schema [41] to extract a linear projection of AbMAP-B embeddings that accentuates the difference between CoV-AbDab antibodies and a baseline set of 3,000 antibodies (the SAbDab subset used to train AbMAP). Schema computes a low-distortion linear projection of an embedding such that distances in the projected space better capture a user-specified co-variate; here, membership in CoV-AbDab (Fig. 4-2b). By applying Schema, we are able to enrich the SARS-CoV-2 relevant signal in the embedding so that a simple model, working with relatively limited training data, can capture the desired biological intuition. On the held-out set of 522 CoV-AbDab antibodies, we annotated each antibody with a binary label indicating if it neutralizes at least one SARS-CoV-2 variant (Alpha, Beta, Delta, Gamma, Omicron). We then applied logistic regression to predict this label from the Schema-projected representation. We evaluated the results as per 5-fold cross validation, also comparing fixed-length embeddings derived from ProtBert and ESM-1b (Table 4.1). AbMAP-based representation, even though it is of much lower dimensionality (50 compared to ESM-1b’s 1,280) has substantially higher balanced-accuracy and F1 score, indicating that our transfer learning approach is more effectively able to hone in on the subspace relevant to SARS-CoV-2 neutralization.

THIS PAGE INTENTIONALLY LEFT BLANK

Chapter 5

Discussion and Conclusion

We presented AbMAP, a transfer learning framework to adapt any foundational PLM (i.e., one trained on the broad corpus of protein sequences) to antibodies, whose hypervariable regions violate the evolutionary conservation that PLMs typically rely on. Our framework refines the raw PLM embedding by *in silico* mutation-based contrastive augmentation. To further shape the augmented embedding, we supervise its non-linear projection to a lower dimensional space such that Euclidean distances in the projected space better capture antibody structural (from PDB) and functional (from Setliff et al.’s LIBRA-seq [38]) similarity. In an ablation study, we found that training on just structural data recapitulated function but not vice-versa (Table A.2). Correspondingly, our multi-task formulation learned a higher weight for the structural task. To maximally leverage the limited structural data on antibodies, AbMAP focuses its capacity on the complementarity-determining regions (CDRs) and their flanking residues, these being the key drivers of antibody specificity. We present AbMAP-E, AbMAP-P, and AbMAP-B, adaptations of the **ESM-1b**, **ProtBert**, and the **Bepler & Berger** PLMs, respectively. The choice of foundational PLM may depend on the type of downstream prediction task using the resulting AbMAP embedding (e.g. AbMAP-B for structure prediction, and AbMAP-E/P for function/property prediction). Our transfer learning framework can be adapted to any new PLM, and we make both the training data and code available for doing so.

AbMAP represents a conceptual advance to language model design for antibodies.

Currently, two broad approaches have been espoused. One, embodied in techniques like OmegaFold [48], is to essentially focus on improving the general PLM, with the expectation that gains will accrue to antibodies as well. The other, represented by methods such as AntiBERTa [21], AbLang [43], and IgLM [40], is to essentially treat antibodies as an entirely new language-modeling task and train solely on corpora of antibody sequences. AbMAP represents a new middle path. We propose a transfer learning approach that can adapt to any foundational PLM and hence benefit from innovations in the underlying PLMs. Since antibody hypervariable regions are not evolutionarily conserved, foundational PLMs will likely remain weaker at modeling antibodies than regular proteins. On the other hand, the framework regions (which comprise $\sim 90\%$ of the sequence) have limited diversity across antibodies. Language models focused specifically on antibodies thus do not benefit from the full diversity of protein sequences that foundational PLMs access. Our approach starts from informative foundational PLMs and then adapts them to be more accurate on antibody structure and function. We additionally make the choice to focus AbMAP’s explanatory capacity on the CDRs and their flanking residues. This is an inductive bias, allowing us to limit the model complexity of the transfer learning layer (one layer of transformer) and enabling robustness as well as accuracy.

We expect AbMAP to be applicable in a variety of cases, e.g., property prediction and antigen binding. We now appraise some of AbMAP’s strengths in the context of structure prediction, as the problem is relatively well-studied and there exist methods focused solely on antibody structure prediction. Our implementation of structure prediction with AbMAP is as a template-finding task, and we leave the task of template refinement to future work. Even unrefined, however, AbMAP performs remarkably well, outperforming AlphaFold 2 and being competitive with OmegaFold in most cases. It outperforms the latter on the functionally-crucial CDR-H3 region. We wondered if our template-finding approach was biasing the results in AbMAP’s favor, and re-evaluated it at progressively lower levels of sequence homology (Fig.)B-2; we also compared AbMAP- $\{B,E,P\}$ with their respective foundational PLM baselines. AbMAP’s performance did not decline meaningfully at lower sequence homologies

and it substantially outperformed its foundational PLM baselines. Notably, the latter did not compare as favorably with OmegaFold or AlphaFold2, suggesting that it is our transfer learning innovations, rather than the template-based approach, that offers the gains.

AbMAP’s design is compatible with a variety of foundational PLM architectures. Here the three foundational models we use differ substantially in their approach: the Bepler & Berger model was training in a multi-task setting and uses protein structure information, the ProtBert model adapts the existing BERT architecture, and the ESM-1b model creates a transformer-based architecture from scratch. Across a variety of applications, the corresponding AbMAP variants outperform their foundational baselines. We believe this is due to our focus on a) CDRs, b) contrastive augmentation that accentuates the residues in the CDRs, and c) a modular architecture with the fine-tuning layer cleanly separated from the foundational PLM. Of these, we believe the last to play the most important role (Fig. 3-2). In our fine-tuning layer, we apply a single transformer layer to adapt contrastively-augmented representations into final AbMAP embeddings; we found that additional layers did not meaningfully increase performance in terms of test set loss and Spearman’s rank (Table A.3). The contrastive augmentation step, though it improves performance and lead to more stable results, also requires multiple invocations of the baseline PLM (one for each mutant). In situations where speed is crucial, this step can be removed for greater efficiency. Towards this, we offer two sets of pre-trained models: one with, and the other without, contrastive augmentation. While it is faster to train AbMAP without contrastive augmentation, the resulting representation space is less consistent with ground truth structural properties (Fig. B-1).

The design of AbMAP represents an explicit focus on the hypervariable region of the antibody at the expense of the framework region. We believe our design choice enables more accurate analysis of human repertoires. Due to the phenomenon on immunoglobulin class (i.e. isotype) switching, constant regions may be replaced while preserving hypervariable regions; this preserves antigen specificity while enabling different effector molecules to bind for downstream effects. Thus, our CDR-focused

approach more representatively captures the diversity of antigen-binding profiles in an individual repertoire. Nonetheless, in some cases of therapeutic antibody design, the framework region may play an important role. To address this, we followed Liberis et al. [23] and expanded the standard Chothia delineation of a CDR to also include two flanking residues on each end. Furthermore, as more antibody datasets become available (e.g., those measuring binding or functional specificity in diverse contexts), the additional training data may enable us to leverage more complex models that effectively capture both the framework and hypervariable regions.

The exploration of human immune repertoires, as well as the design and development of large-molecule therapeutics, require a deep understanding of antibody structure and function, and an ability to efficiently manipulate it in silico. In parallel, stunning advances in general-purpose modeling of proteins do not easily translate to antibodies because of their unique hypervariability. The transfer learning approach of AbMAP represents a general technique to adapt foundational PLMs for specific protein sets of interest – rather than training a dedicated language model for the subset, we argue it is more effective to leverage the fast-moving advances in foundational PLMs and fine-tune for the subset. We believe the scalable and accurate modeling of antibodies enabled by AbMAP will unlock a better understanding of the behavior of antibodies and empower the discovery of novel therapeutic biologics.

Appendix A

Tables

Model	Metric	Chain	CDR 1	CDR 2	CDR 3	Whole
AbMAP-P	TM-Score	H	0.59 ± 0.007	0.61 ± 0.006	0.47 ± 0.006	0.81 ± 0.005
AbMAP-E	TM-Score	H	0.61 ± 0.007	0.64 ± 0.006	0.51 ± 0.007	0.85 ± 0.004
AbMAP-P	TM-Score	L	0.58 ± 0.007	0.71 ± 0.007	0.57 ± 0.007	0.84 ± 0.005
AbMAP-E	TM-Score	L	0.60 ± 0.007	0.74 ± 0.007	0.61 ± 0.007	0.88 ± 0.004
AbMAP-P	RMSD	H	0.44 ± 0.017	0.45 ± 0.014	0.54 ± 0.031	2.16 ± 0.068
AbMAP-E	RMSD	H	0.40 ± 0.013	0.40 ± 0.013	0.50 ± 0.037	1.71 ± 0.057
AbMAP-P	RMSD	L	0.47 ± 0.018	0.18 ± 0.007	0.50 ± 0.020	2.05 ± 0.061
AbMAP-E	RMSD	L	0.41 ± 0.014	0.18 ± 0.006	0.51 ± 0.028	1.67 ± 0.051

Table A.1: Structure prediction scores for AbMAP-P and AbMAP-E. Similar to AbMAP-B, ProtBert, and ESM-1b in Fig. 3.1, the predicted template structure was selected from a set of antibodies whose sequence identity is below 0.7. Structure prediction was conducted on both individual CDR fragments and the whole chain.

Model	Train Data	Test Data	Spearman's ρ
AbMAP-B	SAbDab	LIBRA-Seq	0.57
AbMAP-B	LIBRA-Seq	SAbDab	0.028
AbMAP-E	SAbDab	LIBRA-Seq	0.60
AbMAP-E	LIBRA-Seq	SAbDab	0.049

Table A.2: Test set Spearman's rank scores computed for AbMAP-B,E when using different datasets for train and test (i.e. we evaluated the prediction of functional similarity using embeddings trained on structure data, and vice versa). While embeddings trained solely on structure data perform quite well when predicting function, embeddings trained only on function data can infer far less about antibody structure.

num. transformer layers	1	2	3
Mean Loss	0.272	0.281	0.279
Spearman's ρ	0.826	0.805	0.817

Table A.3: Test set mean loss and Spearman's rank scores computed for AbMAP-B when using different number of Transformer Encoder layers in its architecture. Increase in number of layers did not necessarily show improvement in model performance.

Appendix B

Figures

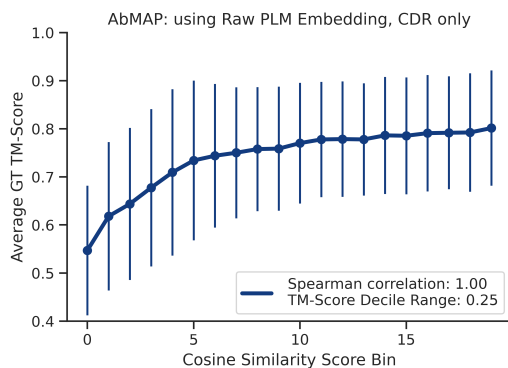


Figure B-1: Average ground truth TM-scores antibody pairs in each of the 20 cosine similarity score bins using AbMAP representations with raw PLM embeddings (CDRs only) as input; i.e. no contrastive augmentation. Monotonicity is noticeably lower when using raw PLM embeddings as input to AbMAP compared to using contrastively augmented PLM embeddings.



Figure B-2: Comparison of AbMAP-B and sequence alignment at different sequence identity thresholds for template search in antibody structure prediction. Columns from left to right (TM-Score Chain H, TM-Score Chain L, RMSD Chain H, RMSD Chain L). Rows from top to bottom (Sequence identity: 0.9, 0.8, 0.7, 0.6, 0.5). For a pair of antibodies, AbMAP-B uses the cosine distance of fixed length sequence embeddings and sequence alignment uses the pairwise global sequence alignment score as a similarity metric.

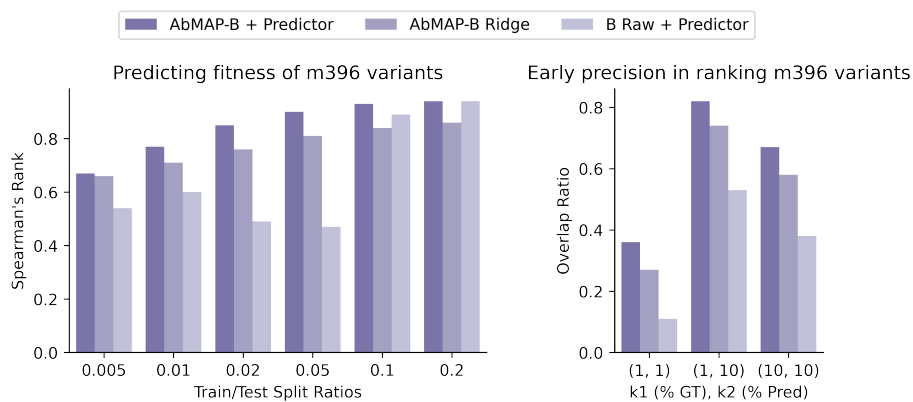


Figure B-3: (Left) Chart of average Spearman's rank scores for ddG scores prediction of Bepler & Berger based models (AbMAP and raw) with various train/test splits. (Right) Chart of average overlap of top- k_1, k_2 ddG scores for different k_1, k_2 values. This experiment was conducted for different train/test splits (train: 0.02, test: 0.98 for this chart). While the raw embeddings are comparable to augmented embeddings at higher splits, their performance notably drops at lower training set sizes.

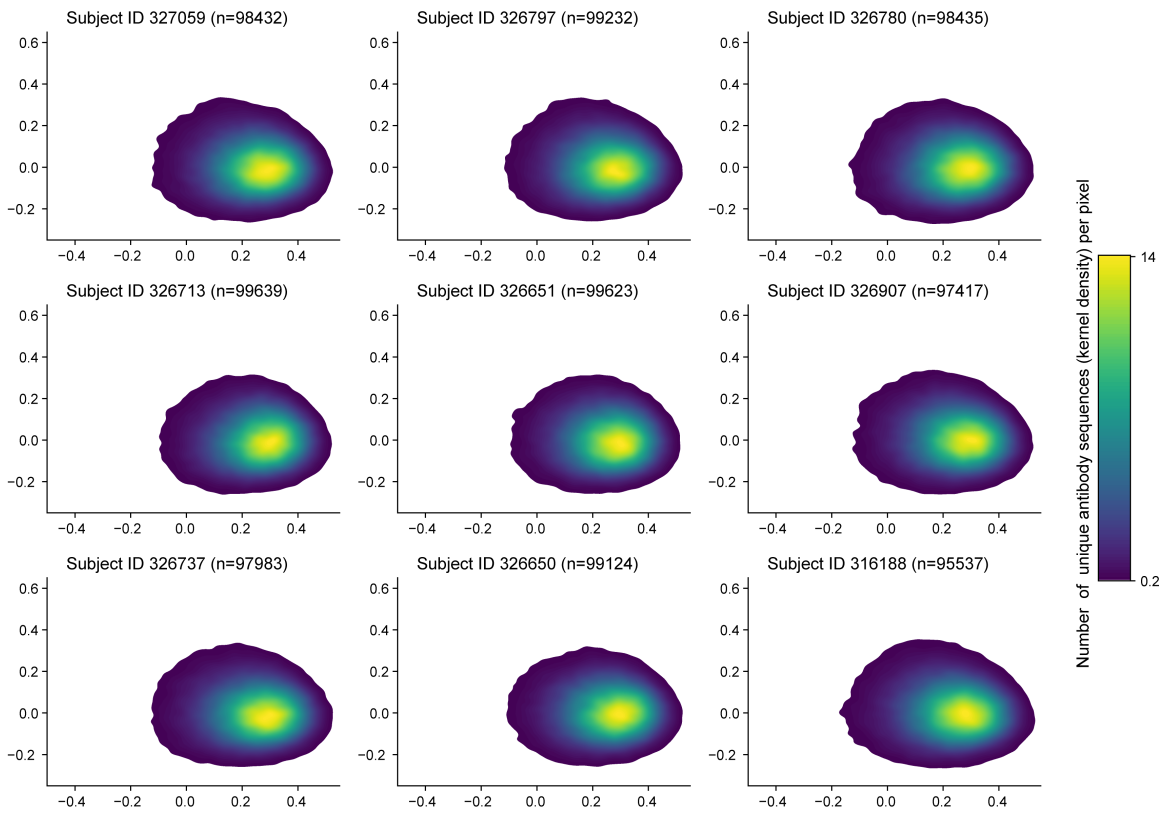


Figure B-4: KDE plot: 2D PCA of AbMAP embeddings for antibody repertoire from each human subject.

Bibliography

- [1] Bissan Al-Lazikani, Arthur M Lesk, and Cyrus Chothia. Standard conformations for the canonical structures of immunoglobulins. *Journal of molecular biology*, 273(4):927–948, 1997.
- [2] Tristan Bepler and Bonnie Berger. Learning protein sequence embeddings using information from structure. In *International Conference on Learning Representations*, 2019.
- [3] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- [4] Bryan Briney, Anne Inderbitzin, Collin Joyce, and Dennis R Burton. Commonality despite exceptional diversity in the baseline human antibody repertoire. *Nature*, 566(7744):393–397, 2019.
- [5] Joe DeBartolo, Mikko Taipale, and Amy E Keating. Genome-wide prediction and validation of peptides that bind human prosurvival bcl-2 proteins. *PLoS computational biology*, 10(6):e1003693, 2014.
- [6] Brandon J DeKosky, Oana I Lungu, Daechan Park, Erik L Johnson, Wissam Charab, Constantine Chrysostomou, Daisuke Kuroda, Andrew D Ellington, Gregory C Ippolito, Jeffrey J Gray, et al. Large-scale sequence and structural comparisons of human naive and antigen-experienced antibody repertoires. *Proceedings of the National Academy of Sciences*, 113(19):E2636–E2645, 2016.
- [7] Benjamin DeMeo and Bonnie Berger. Hopper: a mathematically optimal algorithm for sketching biological data. *Bioinformatics*, 36(Supplement_1):i236–i241, 2020.
- [8] Thomas Desautels, Adam Zemla, Edmond Lau, Magdalena Franco, and Daniel Faissol. Rapid in silico design of antibodies targeting sars-cov-2 using machine learning and supercomputing. *BioRxiv*, 2020.
- [9] James Dunbar and Charlotte M Deane. Anarci: antigen receptor numbering and receptor classification. *Bioinformatics*, 32(2):298–300, 2016.

- [10] James Dunbar, Konrad Krawczyk, Jinwoo Leem, Terry Baker, Angelika Fuchs, Guy Georges, Jiye Shi, and Charlotte M Deane. Sabdab: the structural antibody database. *Nucleic acids research*, 42(D1):D1140–D1146, 2014.
- [11] Ahmed Elnaggar, Michael Heinzinger, Christian Dallago, Ghalia Rihawi, Yu Wang, Llion Jones, Tom Gibbs, Tamas Feher, Christoph Angerer, Martin Steinegger, et al. Prottrans: towards cracking the language of life’s code through self-supervised deep learning and high performance computing. *arXiv preprint arXiv:2007.06225*, 2020.
- [12] Katja Fink. Can we improve vaccine efficacy by targeting t and b cell repertoire convergence? *Frontiers in Immunology*, 10:110, 2019.
- [13] Simon Friedensohn, Daniel Neumeier, Tarik A Khan, Lucia Csepregi, Cristina Parola, Arthur R Gorter de Vries, Lena Erlach, Derek M Mason, and Sai T Reddy. Convergent selection in antibody repertoires is revealed by deep learning. *BioRxiv*, 2020.
- [14] Limin Fu, Beifang Niu, Zhengwei Zhu, Sitao Wu, and Weizhong Li. Cd-hit: accelerated for clustering the next-generation sequencing data. *Bioinformatics*, 28(23):3150–3152, 2012.
- [15] Ivelin S Georgiev, M Gordon Joyce, Yongping Yang, Mallika Sastry, Baoshan Zhang, Ulrich Baxa, Rita E Chen, Aliaksandr Druz, Christopher R Lees, Sandeep Narpala, et al. Single-chain soluble bg505. sosip gp140 trimers as structural and antigenic mimics of mature closed hiv-1 env. *Journal of virology*, 89(10):5318–5329, 2015.
- [16] Ning Jiang, Jiankui He, Joshua A Weinstein, Lolita Penland, Sanae Sasaki, Xiao-Song He, Cornelia L Dekker, Nai-Ying Zheng, Min Huang, Meghan Sullivan, et al. Lineage structure of the human antibody repertoire in response to influenza vaccination. *Science translational medicine*, 5(171):171ra19–171ra19, 2013.
- [17] Wengong Jin, Jeremy Wohlwend, Regina Barzilay, and Tommi Jaakkola. Iterative refinement graph neural network for antibody sequence-structure co-design. *arXiv preprint arXiv:2110.04624*, 2021.
- [18] Jinwoo Leem. Parapred - pytorch. <https://github.com/alchemab/parapred-pytorch>, 2021. Accessed: 2022-11-28.
- [19] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Židek, Anna Potapenko, et al. Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021.
- [20] Andrew Leaver-Fay, Michael Tyka, Steven M Lewis, Oliver F Lange, James Thompson, Ron Jacak, Kristian W Kaufman, P Douglas Renfrew, Colin A Smith, Will Sheffler, et al. Rosetta3: an object-oriented software suite for the

- simulation and design of macromolecules. In *Methods in enzymology*, volume 487, pages 545–574. Elsevier, 2011.
- [21] Jinwoo Leem, Laura S Mitchell, James HR Farmery, Justin Barton, and Jacob D Galson. Deciphering the language of antibodies using self-supervised learning. *Patterns*, page 100513, 2022.
- [22] Weizhong Li and Adam Godzik. Cd-hit: a fast program for clustering and comparing large sets of protein or nucleotide sequences. *Bioinformatics*, 22(13):1658–1659, 2006.
- [23] Edgar Liberis, Petar Veličković, Pietro Sormanni, Michele Vendruscolo, and Pietro Liò. Parapred: antibody paratope prediction using convolutional and recurrent neural networks. *Bioinformatics*, 34(17):2944–2950, 2018.
- [24] Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin, Robert Verkuil, Ori Kabeli, Yaniv Shmueli, et al. Evolutionary-scale prediction of atomic level protein structure with a language model. *bioRxiv*, 2022.
- [25] Ruei-Min Lu, Yu-Chyi Hwang, I-Ju Liu, Chi-Chiu Lee, Han-Zen Tsai, Hsin-Jung Li, and Han-Chung Wu. Development of therapeutic antibodies for the treatment of diseases. *Journal of biomedical science*, 27(1):1–30, 2020.
- [26] Saul B Needleman and Christian D Wunsch. A general method applicable to the search for similarities in the amino acid sequence of two proteins. *Journal of molecular biology*, 48(3):443–453, 1970.
- [27] Richard A Norman, Francesco Ambrosetti, Alexandre MJJ Bonvin, Lucy J Colwell, Sebastian Kelm, Sandeep Kumar, and Konrad Krawczyk. Computational approaches to therapeutic antibody design: established methods and emerging trends. *Briefings in bioinformatics*, 21(5):1549–1567, 2020.
- [28] Ponraj Prabakaran and Partha S Chowdhury. Landscape of non-canonical cysteines in human vh repertoire revealed by immunogenetic analysis. *Cell reports*, 31(13):107831, 2020.
- [29] Roshan Rao, Nicholas Bhattacharya, Neil Thomas, Yan Duan, Peter Chen, John Canny, Pieter Abbeel, and Yun Song. Evaluating protein transfer learning with tape. *Advances in neural information processing systems*, 32, 2019.
- [30] Matthew I. J. Raybould, Aleksandr Kovaltsuk, Claire Marks, and Charlotte M. Deane. CoV-AbDab: the Coronavirus Antibody Database. *Bioinformatics*, 37(5):734–735, 2021.
- [31] Matthew IJ Raybould, Claire Marks, Alan P Lewis, Jiye Shi, Alexander Bujotzek, Bruck Taddese, and Charlotte M Deane. Thera-sabdab: the therapeutic structural antibody database. *Nucleic acids research*, 48(D1):D383–D388, 2020.

- [32] Matthew IJ Raybould, Anthony R Rees, and Charlotte M Deane. Current strategies for detecting functional convergence across b-cell receptor repertoires. In *MAbs*, volume 13, page 1996732. Taylor & Francis, 2021.
- [33] Rajesh P Ringe, Gabriel Ozorowski, Anila Yasmeen, Albert Cupo, Victor M Cruz Portillo, Pavel Pugach, Michael Golabek, Kimmo Rantalainen, Lauren G Holden, Christopher A Cottrell, et al. Improving the expression and purification of soluble, recombinant native-like hiv-1 envelope glycoprotein trimers by targeted sequence changes. *Journal of virology*, 91(12):e00264–17, 2017.
- [34] Alexander Rives, Joshua Meier, Tom Sercu, Siddharth Goyal, Zeming Lin, Jason Liu, Demi Guo, Myle Ott, C Lawrence Zitnick, Jerry Ma, et al. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proceedings of the National Academy of Sciences*, 118(15), 2021.
- [35] Jeffrey A Ruffolo, Jeremias Sulam, and Jeffrey J Gray. Antibody structure prediction using interpretable deep learning. *Patterns*, 3(2):100406, 2022.
- [36] Rogier W Sanders, Ronald Derking, Albert Cupo, Jean-Philippe Julien, Anila Yasmeen, Natalia de Val, Helen J Kim, Claudia Blattner, Alba Torrents de la Peña, Jacob Korzun, et al. A next-generation cleaved, soluble hiv-1 env trimer, bg505 sosip. 664 gp140, expresses multiple epitopes for broadly neutralizing but not non-neutralizing antibodies. *PLoS pathogens*, 9(9):e1003618, 2013.
- [37] Joost Schymkowitz, Jesper Borg, Francois Stricher, Robby Nys, Frederic Rousseau, and Luis Serrano. The foldx web server: an online force field. *Nucleic acids research*, 33(suppl_2):W382–W388, 2005.
- [38] Ian Setliff, Andrea R Shiakolas, Kelsey A Pilewski, Aryn A Murji, Rutendo E Mapengo, Katarzyna Janowska, Simone Richardson, Charissa Oosthuysen, Nagarajan Raju, Larance Ronsard, et al. High-throughput mapping of b cell receptor sequences to antigen specificity. *Cell*, 179(7):1636–1646, 2019.
- [39] Claude Elwood Shannon. A mathematical theory of communication. *The Bell system technical journal*, 27(3):379–423, 1948.
- [40] Richard W Shuai, Jeffrey A Ruffolo, and Jeffrey J Gray. Generative language modeling for antibody design. *bioRxiv*, 2021.
- [41] Rohit Singh, Brian L Hie, Ashwin Narayan, and Bonnie Berger. Schema: metric learning enables interpretable synthesis of heterogeneous single-cell modalities. *Genome biology*, 22(1):1–24, 2021.
- [42] George P Smith. Filamentous fusion phage: novel expression vectors that display cloned antigens on the virion surface. *Science*, 228(4705):1315–1317, 1985.
- [43] Iain H. Moal Tobias H. Olsen and Charlotte M. Deane. Ablang: An antibody language model for completing antibody sequences. *bioRxiv*, 2022.

- [44] Bernhard Valldorf, Steffen C Hinz, Giulio Russo, Lukas Pekar, Laura Mohr, Janina Klemm, Achim Doerner, Simon Krah, Michael Hust, and Stefan Zielonka. Antibody display technologies: selecting the cream of the crop. *Biological Chemistry*, 403(5-6):455–477, 2022.
- [45] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [46] Grant C Weaver, Rina F Villar, Masaru Kanekiyo, Gary J Nabel, John R Mascola, and Daniel Lingwood. In vitro reconstitution of b cell receptor–antigen interactions to evaluate potential vaccine candidates. *Nature protocols*, 11(2):193–213, 2016.
- [47] James RR Whittle, Adam K Wheatley, Lan Wu, Daniel Lingwood, Masaru Kanekiyo, Steven S Ma, Sandeep R Narpala, Hadi M Yassine, Gregory M Frank, Jonathan W Yewdell, et al. Flow cytometry reveals that h5n1 vaccination elicits cross-reactive stem-directed antibodies from multiple ig heavy-chain lineages. *Journal of virology*, 88(8):4047–4057, 2014.
- [48] Ruidong Wu, Fan Ding, Rui Wang, Rui Shen, Xiwen Zhang, Shitong Luo, Chenpeng Su, Zuofan Wu, Qi Xie, Bonnie Berger, et al. High-resolution de novo structure prediction from primary sequence. *BioRxiv*, 2022.
- [49] Tai Te Wu and Elvin A Kabat. An analysis of the sequences of the variable regions of bence jones proteins and myeloma light chains and their implications for antibody complementarity. *The Journal of experimental medicine*, 132(2):211–250, 1970.
- [50] Rui Yin, Brandon Y Feng, Amitabh Varshney, and Brian G Pierce. Benchmarking alphafold for protein complex modeling reveals accuracy determinants. *Protein Science*, 31(8):e4379, 2022.
- [51] Yang Zhang and Jeffrey Skolnick. Scoring function for automated assessment of protein structure template quality. *Proteins: Structure, Function, and Bioinformatics*, 57(4):702–710, 2004.