

**MIT
Libraries**

| **DSpace@MIT**

MIT Open Access Articles

This is a supplemental file for an item in DSpace@MIT

Item title: Scalability strategies for
automated reaction mechanism generation

Link back to the item: <https://hdl.handle.net/1721.1/124333>



Massachusetts Institute of Technology

Scalability strategies for automated reaction mechanism generation

Agnes Jocher^a, Nick M. Vandewiele^a, Kehang Han^a, Mengjie Liu^a, Connie W. Gao^a, Ryan J. Gillis^a, and William H. Green^a

a. Department of Chemical Engineering, Massachusetts Institute of Technology, 77
Massachusetts Avenue, Cambridge, MA 02139, USA

ajocher@mit.edu, nickvandewiele@gmail.com, kehanghan@gmail.com, mjliu@mit.edu,
connie.w.gao@gmail.com, rgillis@mit.edu, whgreen@mit.edu

Abstract

Detailed modeling of complex chemical processes, like pollutant formation during combustion events, remains challenging and often intractable due to tedious and error-prone manual mechanism generation strategies. Automated mechanism generation methods seek to solve these problems but are held back by prohibitive computational costs associated with generating larger reaction mechanisms. Consequently, automated mechanism generation software such as the Reaction Mechanism Generator (RMG) must find novel ways to explore reaction spaces and thus understand the complex systems that have resisted other analysis techniques. In this contribution, we propose three scalability strategies — code optimization, algorithm heuristics, and parallel computing —

that are shown to considerably improve RMG's performance as measured by mechanism generation time for three representative simulations (oxidation, pyrolysis, and combustion). The improvements create new opportunities for the detailed modeling of diverse real-world processes.

1 Introduction

2 Reaction mechanisms are valuable engineering tools that provide insight into the
3 evolution and fate of chemically reactive systems. Coupled with or integrated into
4 multiscale-multiphysics models, reaction mechanisms can help elucidate physical
5 phenomena that are driven by chemical kinetics. Ultimately, this insight can be used to
6 design and optimize a wide range of processes and materials. Traditionally, reaction
7 mechanisms are generated manually, but in the past decade, the computer-assisted
8 development of reaction mechanisms has emerged as a promising tool leveraging the
9 combined power of theoretical insight and extensive experimental data (Burke, 2016).
10 Several software packages are available that harness the power of automation in the field
11 of chemical kinetics (Vandewiele et al., 2012; Broadbelt et al., 1996; Van de Vijver et al.,
12 2015; Rangarajan et al., 2012; Warth et al., 2000; Blurock, 1995). In the following work,
13 we focus on the Reaction Mechanism Generator (RMG) project, which provides a
14 framework for the automated, computer-assisted development of reaction mechanisms,
15 with the goal of producing high-fidelity, predictive, and reproducible kinetic models
16 (Green, 2007).

17 Several recent advancements in kinetic modeling (Vandewiele et al., 2015; Allen et al.,
18 2014; Gao et al., 2015; Prozument et al., 2014; Carr et al., 2015; Seyedzadeh and West,

1 2016; Class et al., 2016) and novel chemical kinetic applications (Jalan et al., 2010, 2013;
2 Suleimanov and Green, 2015) are posing new challenges to RMG. For instance, to
3 improve fidelity, single-component fuel surrogates are being replaced by more realistic
4 multicomponent formulations (Narayanaswamy et al., 2016). This enhances
5 understanding of the impact of existing and novel feedstock blends on the overall process,
6 but greatly increases the complexity of the modeling, challenging automated strategies.
7 Other advancements in kinetic modeling aim to reduce the uncertainties in model
8 parameters. For instance, the Quantum Mechanics Thermodynamic Property (QMTP)
9 module for estimating thermodynamic properties using on-the-fly quantum chemistry
10 calculations (Magoon and Green, 2013) complements the existing Benson group
11 additivity methods (Benson, 1968) that suffered from limited accuracy for some classes
12 of molecules, including polycyclic species. Quantum chemistry has been used to improve
13 the thermochemical parameters in a kinetic model, at the cost of significantly increasing
14 the CPU time required (Magoon and Green, 2013). Methods for calculating pressure-
15 dependent rate coefficients improve estimates for rate coefficients of pressure-dependent
16 reactions (Allen et al., 2012). In addition, capabilities such as sensitivity analysis (Gao et
17 al., 2016) and uncertainty quantification (Gao et al., 2019) are being added to RMG. All
18 the additions may facilitate the construction of high-fidelity kinetic models, but they also
19 increase computational demands. Consequently, without addressing the scalability of the
20 algorithms, many of the current and upcoming features will be available in theory, but
21 unaffordable in practice. Although considerable effort has been spent in developing robust
22 and scalable numerical schemes for solving the sets of equations associated with specific
23 chemical processes, such as combustion (Shi et al., 2011), there has been little focus on

1 the scalability of computer-assisted reaction mechanism generation. This shortcoming is
2 addressed in the present contribution. We begin by assessing the computational
3 complexity of RMG simulations and identifying the most time-consuming bottlenecks for
4 three simulations selected to represent the diversity of computational challenges faced
5 by automatic mechanism generation. Barriers that inhibit scalability and performance are
6 discussed and guide the proposed changes to code optimization, algorithm heuristics,
7 and parallel computing. These enhancements are then explained and analyzed in detail.
8 Finally, the strategies are evaluated with respect to walltime and memory consumption in
9 the three RMG test simulations.

10 2 Automated reaction mechanism generation

11 2.1 RMG - Overall structure and key components

12 This section outlines the overall structure and the key components that govern the
13 computational performance of RMG simulations. Details regarding the chemical
14 meaningfulness and accuracy of the employed methods have been published previously
15 (Gao et al., 2016; Hansen et al., 2013; Zhang et al., 2018; Gudiyella et al., 2018). RMG
16 uses the rate-based algorithm of Susnow et al. (1997) to iteratively grow (“enlarge”) the
17 reaction mechanism one or more species at a time. The starting point in this enlargement
18 scheme is a user-defined set of species that corresponds to the reactants of the modeled
19 chemical process and acts as the initial species of the reaction mechanism to be enlarged.
20 An iterative procedure generates all reactions between the species of the current
21 mechanism, called the “core”. These reactions are synchronized to avoid double
22 occurrences and the resulting product species that are not yet part of the core are added

1 to the “edge”. The rate-based enlargement algorithm dictates that an edge species is
2 moved to the core when the flux R_i to the edge species i exceeds a certain threshold
3 value, Eq. 1.

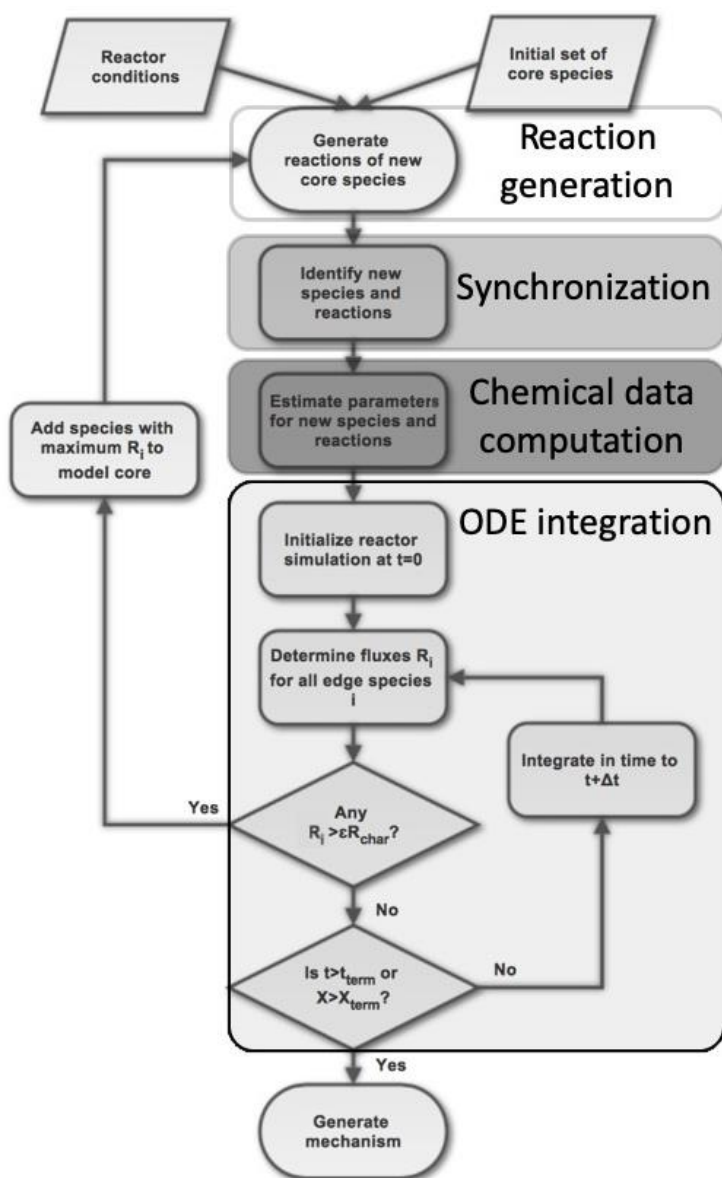
$$R_i = \frac{dC_i}{dt} > \epsilon \cdot R_{char} \quad \text{Eq. 1}$$

4 with ϵ being the user-specified error tolerance and R_{char} the characteristic flux of the
5 system, defined by Eq. 2

$$R_{char} = \sqrt{\sum_j R_j^2} \quad \text{Eq. 2}$$

6 for all core species j . An isothermal, isobaric batch reactor model is used to evaluate
7 species fluxes and species concentrations. A system of ordinary differential equations
8 (ODEs) containing only the reactions from the mechanism core is integrated in time until
9 the rate of formation of a species in the edge exceeds the threshold value. Whenever an
10 edge species is added to the core, the reactor is initialized back to time zero and species
11 fluxes and concentrations are re-evaluated through integration. This iterative enlargement
12 procedure is the main loop through which the mechanism is grown. The enlargement
13 procedure of the mechanism is complete when none of the edge species fluxes meets
14 the threshold value and the user-defined termination criteria, i.e., residence time for the
15 simulation ($t > t_{term}$) or a specific conversion of a given initial species ($X > X_{term}$) is reached.
16 The final model includes all the species and reactions in the core. To obtain all the
17 information that is required to decide which species should be added to the core, the
18 enlargement procedure passes through four key phases: reaction generation,

1 synchronization, chemical data computation, and finally, reaction system integration, see
 2 Fig. 1.



3 **Figure 1: The schematic outlines the workflow for automated reaction mechanism generation with one reactor**
 4 **for an isothermal, isobaric condition. The enlargement step that integrates a species from the edge into the**
 5 **core consists of four key components: reaction generation, synchronization, calculation of chemical data such**
 6 **as thermodynamic properties of species and finally solving the system of ODEs including the calculation of**
 7 **species concentrations and fluxes.**

1 2.1.1 Reaction generation

2 In the reaction generation phase, the species that was added in the previous enlargement
3 iteration is used as a reactant to generate new reactions and species. For unimolecular
4 reactions, this is done by iteratively matching the new reactant to the templates of a series
5 of built-in unimolecular reaction families. For bimolecular reactions, the new species is
6 paired with each species already present in the model core and then matched to the
7 bimolecular templates. When a set of compatible reactants and template are found,
8 products are created by applying the reaction family-specific recipe. This recipe precisely
9 describes the transformation from reactants to products with instructions to break and
10 form specific bonds, gain or lose electrons, etc. The reaction generation step results in
11 the creation of a large number of, potentially redundant, new reactions and species which
12 requires additional processing called synchronization.

13 2.1.2 Synchronization

14 During synchronization, RMG iterates over the list of newly generated reactions and
15 species and determines if any are redundant. If a newly generated species is identical to
16 a species in memory, RMG deletes the new structure and replaces it with a reference to
17 the existing structure in memory. A similar mechanism is in place for reactions. RMG
18 requires synchronization because tracking unique species and reactions allows
19 distinction between the model's core and edge. This distinction is required for the iterative
20 enlarging procedure and allows for the proper formulation of the system of ODEs to solve.
21 Furthermore, synchronization avoids the redundant computation of chemical data. The
22 uniqueness of species is established through isomorphism comparisons of the
23 corresponding chemical graphs that represent the species. The uniqueness of reactions

1 is determined by comparing the reactants and products of the reactions and hence also
2 relies on graph isomorphisms. If there are multiple transition states of comparable energy
3 connecting these species then the rates should be summed together to give the overall
4 flux from some set of reactants to some set of products.

5 2.1.3 Chemical data computation

6 After the synchronization phase has determined a list of new species and reactions, the
7 chemical data computation phase calculates thermodynamic properties of the species or
8 kinetic parameters of the reaction. Both thermodynamic property and kinetic parameter
9 calculations rely heavily on performing subgraph isomorphism queries on a large number
10 of substructures contained in RMG's databases. The uncertainty of the models due to the
11 uncertainty of the thermodynamic and kinetic parameters can be estimated using recently
12 added tools (Gao, 2016). Different models and observables have different degrees of
13 sensitivity to these uncertainties precluding any blanket statements about uncertainty
14 requirements.

15 2.1.4 Reaction system integration

16 In this final step, the state variables are re-computed, and a new species is selected and
17 added to the core. The system of ODEs that corresponds to the new state of the reaction
18 mechanism "core" is solved. The first edge species whose formation rate exceeds a
19 threshold, see Eq. 1 and Eq. 2, is chosen as the next species to be added to the core.
20 From here the cycle restarts, with the next reaction generation step based on the newly
21 added species.

1 2.2 Performance bottlenecks

2 To assess performance bottlenecks of the automated reaction mechanism generation
 3 tool, three simulations are analyzed in depth. They are selected to represent a diverse
 4 collection of cases that make use of a large part of RMG's features and settings, Tab. 1.

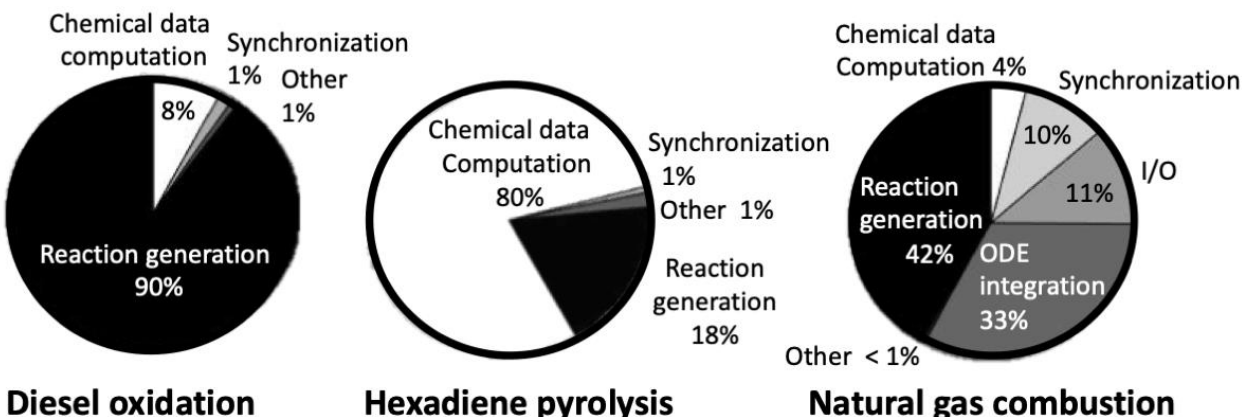
5 **Table 1: Characteristics of three representative RMG simulations.**

		Initial stage of diesel oxidation	Hexadiene pyrolysis	Rich natural gas combustion
Initial mixture		Six-component diesel surrogate, oxygen	Hexadiene, methane, hydrogen, nitrogen	Natural gas, air
Isothermal and isobaric operating conditions		500 K	1350 K	1400 K
		200 bar	1 bar	20 bar
Thermodynamic property calculation method		Benson group additivity	QMTP	Benson group additivity
Initial molar ratio		n-C ₁₁ : n-C ₁₃ : n-C ₁₆ : n-C ₁₉ : n-C ₂₁ : n-decylbenzene : O ₂ = 1.00 : 1.27 : 1.67 : 1.20 : 0.67 : 0.80 : 0.33	1,3-hexadiene : CH ₄ : H ₂ : N ₂ = 1.00 : 152.00 : 23.00 : 1288.00	CH ₄ : C ₂ H ₆ : C ₃ H ₈ : O ₂ : N ₂ = 1.00 : 0.09 : 0.02 : 0.71 : 0.33
Kinetics source		<ul style="list-style-type: none"> ▪ Training reactions ▪ RMG 'default' kinetics families ▪ Rate rules 	<ul style="list-style-type: none"> ▪ Training reactions ▪ RMG 'default' kinetics families ▪ Rate rules 	<ul style="list-style-type: none"> ▪ Training reactions ▪ RMG 'default' kinetics families ▪ Rate rules
Final model size (Base case)	Core	Species: 26 Reactions: 24	Species: 131 Reactions: 4689	Species: 376 Reactions: 17,739
	Edge	Species: 13,429 Reactions: 41,355	Species: 44,604 Reactions: 156,906	Species: 679,127 Reactions: 1,691,134

1 The first case is a partial oxidation of a diesel surrogate, which consists of five larger
2 linear alkanes and n-decylbenzene. The automatic mechanism generation of this system
3 leads to a large number of generated reactions even in the early stages of the simulation.
4 In the second case, a kinetic mechanism for the pyrolysis region of a 1,3-hexadiene doped
5 methane flame is generated. This model could be employed to study soot formation
6 (Sharma et al., 2007; Sharma et al., 2010, Vermeire et al., 2017). To address the need
7 for accurate thermodynamic properties of cyclic and polycyclic species in these studies,
8 the Quantum Mechanics Thermodynamic Property (QMTP) method is used rather than
9 Benson group additivity. Case three generates a mechanism for natural gas combustion
10 containing methane, ethane, and propane as the initial reactants. Due to the reactant
11 molecules being small in size, only a small number of reactions per iteration is generated
12 early in the simulation, contrasting case one. However, this simulation is run long enough
13 that the generated mechanism grows comparatively large in size. The main purpose of
14 the present work is to improve the scalability of our automated reaction generation tool.
15 To test the improvements within a reasonable time frame, the initial stage of diesel
16 oxidation and rich natural gas combustion simulations were evaluated after fourteen days
17 of wall clock time. For the hexadiene pyrolysis case, that is dominated by the computation
18 of thermodynamic properties with the QMTP method, results are evaluated for 130 core
19 species.

20 Figure 2 compares the four key phases of mechanism enlargement and their contribution
21 to the overall mechanism generation time for each of the three cases. These cases

1 will be called “base cases” in the remainder of the manuscript.



2 **Figure 2: Contribution of key enlargement phases to overall kinetic mechanism generation time.**

3 For the diesel oxidation case, which involves longer n-alkanes with multiple reacting sites,
4 the reaction generation phase is by far the largest performance bottleneck. The
5 synchronization step, the chemical data computation phase, and the ODE integration step
6 have minor or negligible contributions to the overall simulation time. In case two,
7 hexadiene pyrolysis, which employs the computationally expensive QMTP method for
8 estimating thermodynamic properties, the chemical data computation phase becomes the
9 limiting factor. Finally, as the RMG simulation advances to larger mechanism sizes in the
10 natural gas combustion case, a significant amount of time is spent in both the reaction
11 generation phase as well as the ODE integration phase. In this third case, the chemical
12 data computation, synchronization and input/output (I/O) events play a minor role. These
13 observations suggest that the computational load of RMG simulations cannot be isolated
14 in a single routine, but rather, is dynamic in nature and heavily dependent on the chemical
15 system under investigation. For example, while improvements to assignment, tree-
16 structure search, and reaction path degeneracy calculations could accelerate mechanism
17 generation in the diesel oxidation case, these strategies would have negligible impact in
18 the hexadiene pyrolysis case. As a result, multi-faceted strategies that focus on

1 computational performance across the spectrum of simulations are required to
2 significantly reduce the wall clock time spent generating the kinetic reaction mechanisms.

3 3 Scalability strategies

4 3.1 Challenges and tradeoffs

5 Scalability strategies for automated reaction mechanism generation face four main
6 challenges: continuity, programming language, information-rich data structures, and data
7 consistency.

8 3.1.1 Continuity

9 The RMG project has existed for almost two decades. Many incremental developments
10 have grown RMG into a codebase containing about 125 400 lines of Python code and 13
11 900 lines of Cython code, including a separate text-based database with about 277 700
12 lines of Python code for chemical data storage. A single codebase simultaneously
13 maintains a Django-powered website (<http://rmg.mit.edu>), serves as production code on
14 high performance clusters, and supports desktop versions on Linux and Mac OSX.
15 Furthermore, as a collaborative, open-source project with over 2000 commits added in
16 the last year by more than twenty developers in distributed development teams, the
17 capabilities of RMG are constantly expanding. However, each new update poses a risk
18 to the accuracy and performance of the software. While more than one thousand unit and
19 regression tests are embedded in the automated continuous integration system to catch
20 blatant bugs, subtler changes resulting in structural or numerical deviations in the final
21 converged mechanisms can arise from the synergistic effects of coupling multiple
22 features. These errors are more difficult to discover. Therefore, any code modification

1 must strike a balance between performance and software reuse and robustness, which
2 inevitably results in compromises to both.

3 3.1.2 Programming language

4 One such compromise is the choice of the programming language. RMG is implemented
5 in Python, after a 2012 rewrite from the original version in Java (Song, 2004). Python has
6 a rich ecosystem of scientific libraries that enable fast-paced development and ease of
7 portability. RMG currently contains over twenty external dependencies ranging from low-
8 level numerical libraries such as NumPy (van der Walt, 2011) and SciPy (Jones and
9 Oliphant) to graph visualization toolkits such as GraphViz (Gansner and North, 1999).
10 However, due to the interpreted, dynamically typed nature of Python, it bears significant
11 performance penalties compared to other popular choices of programming languages for
12 scientific software.

13 3.1.3 Information-rich data structures

14 The fundamental entities that are manipulated in RMG are “Species” and “Reaction” data
15 structures. They are highly nested, object-oriented data structures that systematically
16 categorize the layers of molecular information needed by the different parts of RMG.
17 Moreover, the species and reaction data structures do not merely store the primitive
18 variables required within a simulation, they also contain meta-variables for analysis,
19 visualization and post-processing of the generated reaction mechanisms. For example,
20 upon calculation of the thermodynamic properties of a species via the Benson group
21 additivity method (Benson, 1968), the matched groups and origin of the groups will be
22 stored as part of the Species data structure and can be verified *posteriori* to confirm their
23 validity. The downside of using such high-level, information-rich data structures is that

1 their manipulation is difficult to translate into the arithmetic operations of mathematical
2 data structures such as arrays of floating-point numbers. Instead, the types of operations
3 that species and reactions undergo are complex and require accessing substantial out-
4 of-cache data, which in turn leads to suboptimal performance on architectures designed
5 for compute-bound mathematical problems.

6 3.1.4 Data consistency

7 RMG extensively makes use of references to existing objects in memory. This strategy,
8 inherent to many object-oriented software paradigms, not only prevents anomalies during
9 the RMG simulation, but also reduces memory requirements and avoids redundant
10 computation and duplication of associated data. The extensive use of object references
11 is not only a result of the choice of programming model but is also inherent to the nature
12 of automated reaction mechanism generation. Typically, for mechanisms generated by
13 RMG consisting of about 100 species and 3000 reactions, every unique species will be
14 created 300 times during each reaction generation phase before being called in
15 synchronization. Further, each species appears in seventy-five distinct reactions on
16 average. The extensive use of object references in lieu of newly created objects in
17 memory undoubtedly leads to performance gains, it does, however, pose problems for
18 parallel computing.

19 3.2 Proposed scalability strategies

20 Accounting for the iterative nature of RMG's mechanism generation approach, the
21 identified performance bottlenecks, and the discussed challenges and tradeoffs for
22 scalability strategies, we have focused on performance improvements within the
23 individual phases of a mechanism enlargement loop. While post-processing features such

1 as periodic checkpointing and intermediate mechanism writing may significantly increase
2 the computational cost, they can be deactivated and are therefore not further considered
3 for performance optimization. Conventional strategies to reduce the computational load
4 of simulations, such as those that are employed by computational fluid dynamic software,
5 e.g., domain discretization and decomposition, don't apply to RMG. The strategies
6 proposed in this work can be roughly categorized into code optimization, algorithm
7 heuristics, and parallel computing.

8 3.2.1 Code optimization

9 RMG uses several code optimization techniques. To reduce the performance penalty
10 imposed by Python, large parts of RMG that contain critical code, such as the VF2 (Foggia
11 et al., 2001; Cordella et al., 2004) algorithm used for (sub-)graph isomorphism have been
12 converted to Cython (Behnel et al., 2011). The Cython project compiles Python-style code
13 to C through static typing, leading to an expected speed up of about an order of magnitude
14 for numerically intensive code. Another technique involves prescreening to determine that
15 molecule pairs are non-identical before calling the graph isomorphism algorithm.
16 Currently, the prescreening uses the molecular formula of the molecule as a fingerprint
17 and filters out about 80% of the molecule pairs. Finally, the integration of the ODE system
18 is performed by DASPK (Li and Petzold, 1999), a Fortran-based differential algebraic
19 system solver accessed through the PyDAS (PyDAS) Python interface.

20 3.2.2 Algorithm heuristics

21 To further target the reaction generation bottleneck, especially for simulations with large
22 molecules, a "reaction filtering" approach is implemented. Reaction filtering reduces the
23 number of reactions generated with each iteration trading off kinetic mechanism

1 completeness for speed-up by neglecting species with extremely low concentrations. Per
2 the laws of mass action, the rates of reactions whose reactants have very low
3 concentrations must be very low as well. Very slow reactions necessarily contribute little
4 to the overall rate of formation of high flux, and thus important, species. If a very slow
5 reaction contributes significantly to the overall rate of formation of a species, that species
6 must be formed in only low concentrations and thus can be excluded from the mechanism.
7 Despite the limited importance of slow reactions to the overall mechanism, the time spent
8 in the reaction generation phase is independent of the kinetic (un)importance of a
9 reaction. A heuristic approach is chosen that prevents species from undergoing reactions
10 if the highest achievable reaction rate between those species is below a threshold. In
11 each iteration, prior to the generation of reactions, the following inequality is verified:

$$k_{max} \cdot \prod_i c_i < \epsilon \cdot R_{char} , \quad \text{Eq. 3}$$

12 with k_{max} being an upper limit for the rate coefficient of a reaction between the reactants
13 i , c_i being the concentration of reactants evaluated at any given time in the reaction
14 system, ϵ being the model error tolerance and R_{char} being the characteristic flux defined
15 in [Eq. 1Eq. 1](#). The left-hand side of Eq. 3 can be interpreted as the highest reaction rate
16 possible for the examined reaction. The choice of $\epsilon \cdot R_{char}$ allows user input on the
17 tolerance through ϵ , while ensuring that the threshold follows the dynamics of the
18 simulated reaction system through R_{char} . Values for the upper limit of k_{max} are chosen
19 based on the reaction order. For unimolecular reactions, k_{max} is set to $(k_B \cdot T)/h$, where
20 k_B is the Boltzmann constant, T is the temperature and h is Planck's constant. This
21 maximum can be interpreted as the rate coefficient from the Eyring equation (Eyring,
22 1935) when the Gibbs free energy of activation is equal to zero. For bimolecular reactions,

1 k_{max} is set to $10^8 \text{ m}^3 \text{ mol}^{-1} \text{ s}^{-1}$, a typical high-pressure-limit rate coefficient for radical
2 recombination reactions where the radicals are polyatomic molecules. Since barrier-less
3 radical recombination reactions are among the fastest bimolecular reactions, it is
4 expected that this value for k_{max} should not exclude important reactions. To ensure the
5 chemical accuracy of a generated model the validity of the chosen k_{max} has to be carefully
6 studied. Here we just employ the reasonable default RMG values.

7 On a similar basis, another heuristic approach, called “species flux pruning” (Han et al.,
8 2017), was implemented to reduce the memory requirements of RMG simulations. In this
9 algorithm, normalized fluxes are used as a metric to identify negligible species during
10 model generation to prune them and the reactions in which they participate. A potential
11 risk of species flux pruning is the erroneous identification of important species as
12 unimportant and their subsequent deletion, leading to a loss of model accuracy. However,
13 careful selection of pruning parameters allows reduced memory requirements for building
14 accurate kinetic models, without loss of accuracy. For example, for a model with 200 to
15 300 species the memory requirement can be reduced by about a factor of four, allowing
16 creation of larger models that were unreachable using earlier versions of RMG due to
17 limited hardware resources.

18 3.2.3 Parallel computing

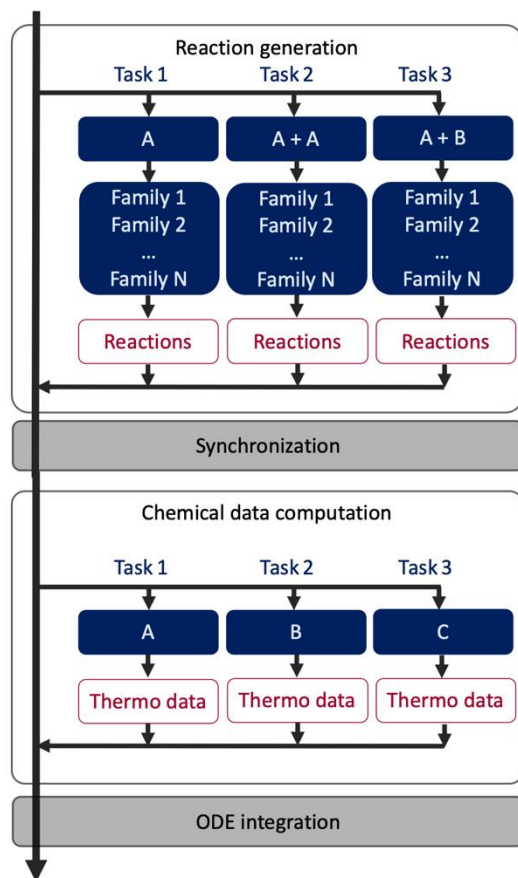
19 For several years processor clock speeds have been relatively stable and are not
20 expected to increase substantially in the near future (Millett and Fuller, 2011). However,
21 new chip designs allow for the addition of more cores to a single processor to increase
22 code performance. In addition, depending on the hardware manufacturer, each core can
23 allocate a number of threads. Taking advantage of the hardware architecture,

1 computationally intensive parts of RMG such as reaction generation and chemical data
2 computation can be executed in parallel. Python's multiprocessing package effectively
3 side-steps the *Global Interpreter Lock* (GIL) by using processes instead of threads. The
4 GIL mechanism is used by the CPython interpreter to assure that only one thread
5 executes Python bytecode at a time aiming to make the object model, e.g., Python
6 dictionaries, implicitly safe against concurrent access (Python). In this work, Python's
7 multiprocessing package's pool object is used to parallelize function execution across
8 multiple input values by distributing the input data across processes, so called data
9 parallelism (Python).

10 Due to the dominant overall contribution of the reaction generation and chemical data
11 calculation phases to the automated mechanism generation time, parallelizing these two
12 computationally intensive parts of RMG is expected to result in considerable performance
13 gains. The chemical data computing phase may become more demanding in the future
14 as more advanced quantum chemical methods replace the semi-empirical methods
15 currently supported (MOPAC2016). A schematic of the parallel workflow is presented in
16 Fig. 3.

17 Despite the potential for parallel computing in the aforementioned target areas, highly
18 non-uniform task workloads hamper the parallel efficiency as it becomes increasingly

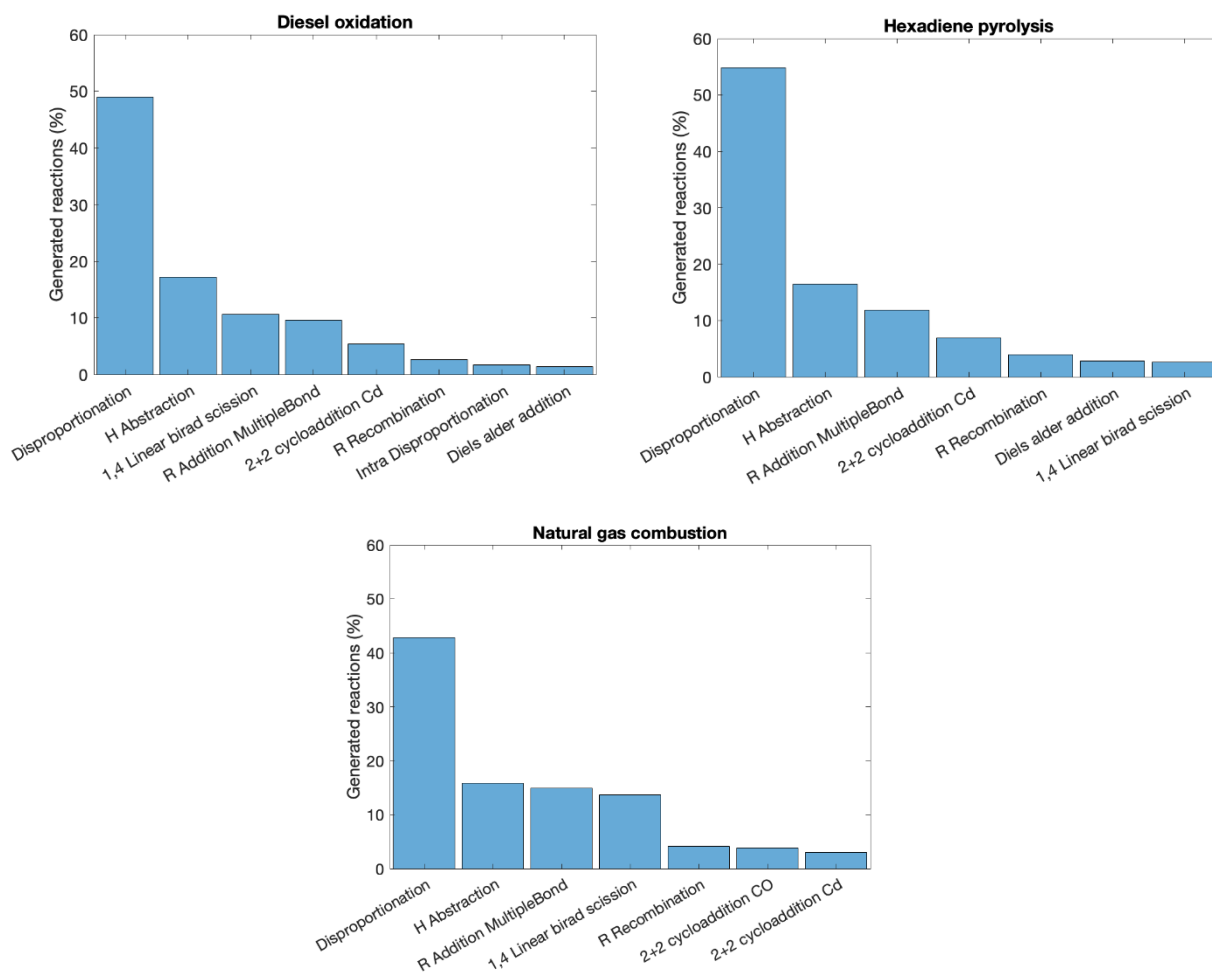
1



2 **Figure 3: Parallelization scheme for the reaction generation and chemical data computation phases as part of**
3 **one enlargement iteration in RMG. The scheme distributes the generation of reactions by creating one task per**
4 **molecule and molecule-molecule pair and subsequently creates reactions through application of the list of**
5 **reaction families. The generated reactions created per task are sent to the base process, bundled together and**
6 **further processed in the synchronization phase, where newly generated species are identified that have no**
7 **thermodynamic properties associated before creating one task per species in the chemical data computation**
8 **phase for submission to the QMTP part of RMG. After all tasks are completed the spawned processes are**
9 **joined together to continue the ODE integration.**

10 difficult for the task scheduler to achieve a uniform workload across the available
11 processes. Figure 4 presents the percentage of generated reactions per reaction family

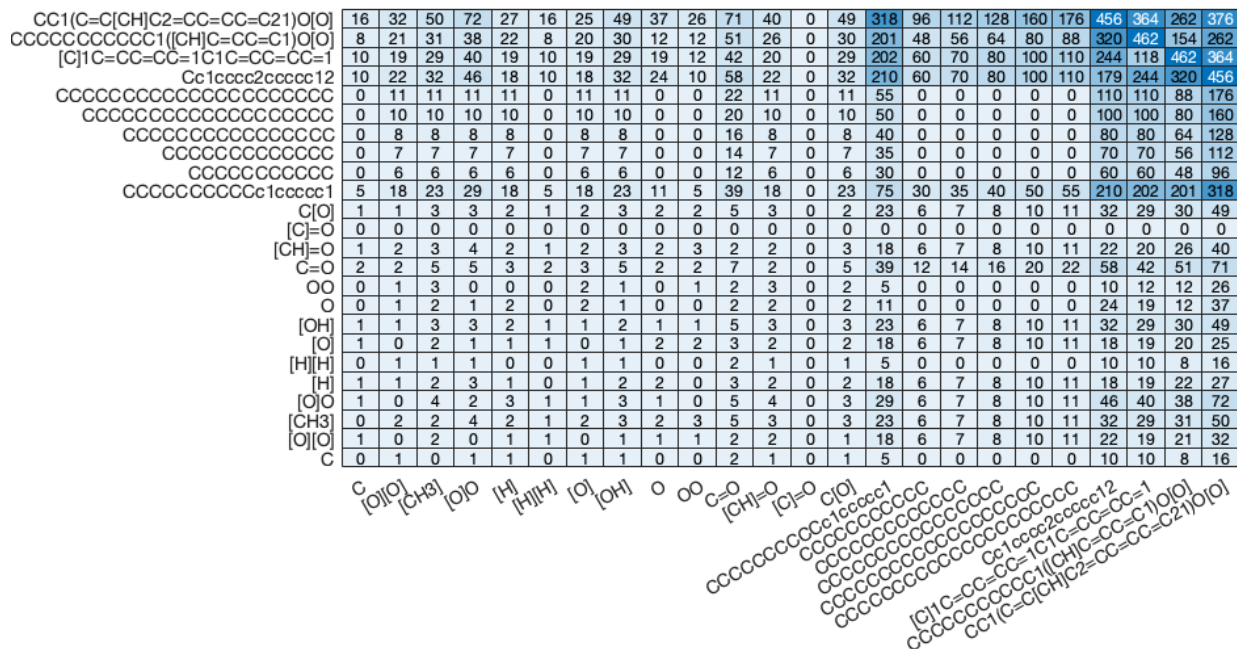
1 over the course of the simulation for the selected cases in Tab. 1.



2
3 **Figure 4: Generated reactions per reaction family through the course of a simulation. Out of 66 available**
4 **reaction families in RMG, only a limited number of families is responsible for the large majority of the generated**
5 **reactions. Only reaction families that contribute more than 1% to the generated reactions are presented.**

6 For all three cases most reactions are generated using the disproportionation and H-
7 abstraction reaction family. However, the number of generated reactions is highly
8 dependent on the molecules that are paired together. For a molecule pairing that is
9 highly reactive, often with a larger number of atoms, submitting one task including all

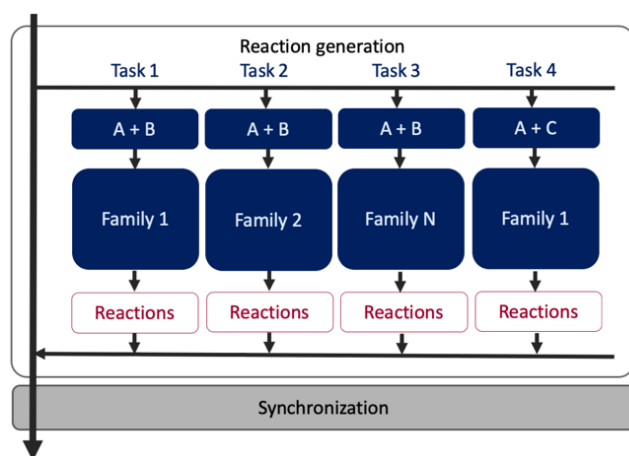
1 available reaction families in RMG will lead to rather large loads in comparison to the
 2 load associated with a more unreactive molecule pairing. Figure 5 is a heat map of the
 3 number of generated reactions throughout an RMG simulation between two species
 4 represented as SMILES. It can be observed that the number of reactions per molecule
 5 pair is highly dependent on the nature of the species and varies widely.



6 **Figure 5: Heat map showing the number of generated reactions between two species presented as SMILES**
 7 **during an RMG simulation. Darker colors imply a higher number of generated reactions.**

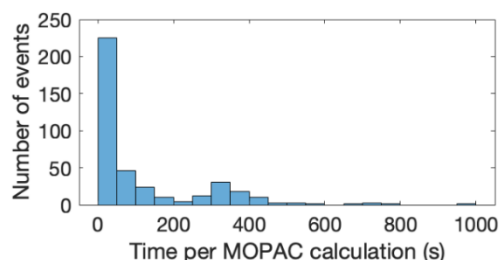
8 In response to this imbalance, parallelization across the reaction families for molecules
 9 larger than ten non-hydrogen atoms is implemented in RMG to create tasks with a
 10 computational load that is more uniformly distributed. A threshold of ten non-hydrogen
 11 atoms has been chosen based on scaling tests exhibiting a performance loss when
 12 employing family splitting with species that have a smaller number of atoms. The

1 schematic of the workflow is outlined in Fig. 6.



2
3 **Figure 6: Reaction generation schematic with family splitting to address unbalanced task loads. In Task 1,**
4 **bimolecular reactions using reactant A and B and RMG’s reaction family “Family 1” are created. In Task 2,**
5 **bimolecular reactions using reactant A and B and “Family 2” are created and in Task 3, bimolecular reactions**
6 **using several families “Family N” are created. The same task generation method is then applied for bimolecular**
7 **reactions using reactant A + C and “Family 1” like presented in Task 4.**

8 Load distribution is relatively well balanced for the parallel calculation of thermodynamic
9 properties. Figure 7 shows a histogram of the time to completion for the calculation of the



10 **Figure 7: Distribution of the individual times to successfully complete an QMTP calculation for 390 unique**
11 **hydrocarbons (C₂ – C₁₈) using MOPAC2016 at the PM7 level.**

12 thermodynamic properties of 390 hydrocarbons ranging from C₂ to C₁₈ using QMTP
13 methods. The majority of the calculations take less than 25 s to complete, although some
14 require more than 50 s. Due to the large number of events less than 25 s, load balancing

1 is expected to be less of an issue in this case. Furthermore, it is outside the scope of this
 2 work to modify the computation time of semi-empirical methods called from within RMG.

3 4 Benchmarking the scalability strategies

4 4.1 Parallel efficiency and load distribution

5 To evaluate the parallel performance, the elapsed wall clock time or real-time and the
 6 parallel efficiency for the reaction generation and thermodynamic property calculation
 7 phase are analyzed outside an RMG simulation in three standalone test cases, A to C,
 8 with results shown in Tab. 2.

9 **Table 2: Standalone testcases to analyze wall clock time and parallel efficiency for parallel reaction**
 10 **generation and chemical data computation. The reactants are obtained from cases presented in Tab. 1.**

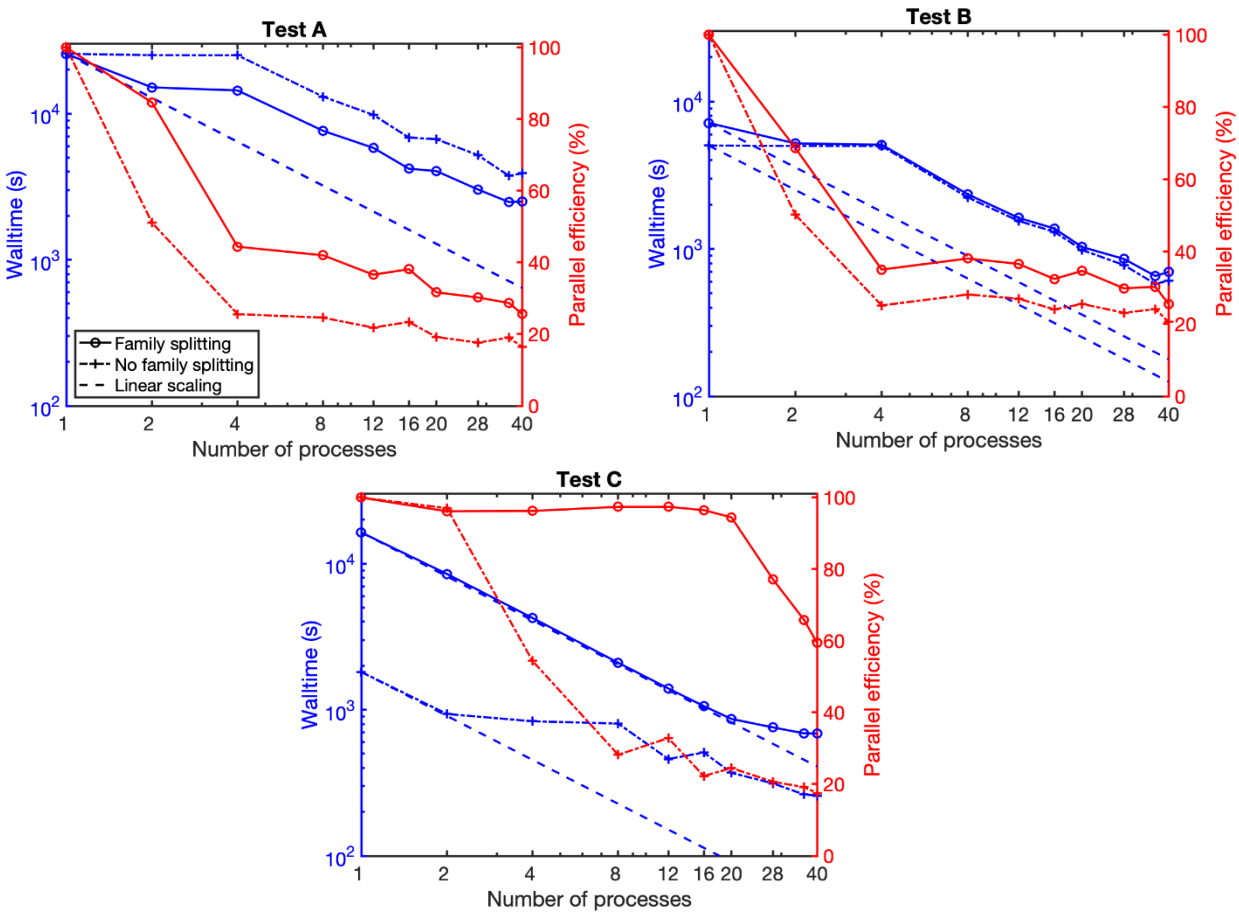
		A	B	C	Thermodynamic property calculation
Number of reactants		52	50	100	390
Reactants obtained from		Initial stage of diesel oxidation	Hexadiene pyrolysis	Rich natural gas combustion	Unique hydrocarbons ranging from C ₂ to C ₁₈
Generated reactions		25 987	51 953	37 675	-
Wall clock time reduced to	Family splitting	10%	10%	4%	5%
	No family splitting	15%	12%	14%	
Speedup factor 40 processes	Family splitting	10.3	10.2	24.0	18.0
	No family splitting	6.6	8.3	7.0	

11
 12 The average number of atoms per reactant decreases from test A to C. Strong scaling
 13 tests were chosen for the analysis meaning the scaling behavior of a fixed problem size

1 with an increasing number of processes was assessed. The parallel efficiency measures
 2 the efficiency of processes to execute a given parallel algorithm and is calculated with
 3 Eq. 4:

$$E = \frac{T_1}{m \cdot T_m}, \quad \text{Eq. 4}$$

4 where T_1 and T_m are wall clock times to complete the problem, using one up to m equipped

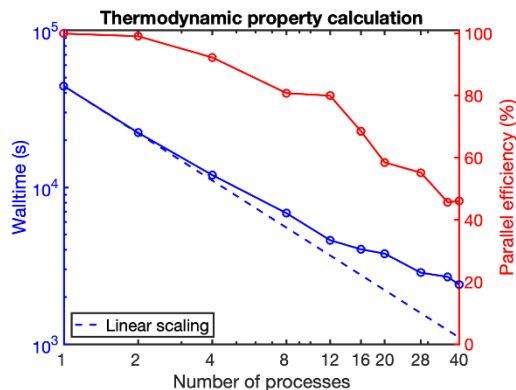


5 **Figure 8: Wall clock time and parallel efficiency for strong scaling tests using core species generated in the initial diesel**
 6 **oxidation, hexadiene pyrolysis, and natural gas combustion cases as reactants for reaction generation. The maximum**
 7 **achievable speed-up is presented as dashed lines for linear scaling.**

8 processes, respectively. The scaling tests were carried out on a MacBook Pro (Early
 9 2015, 2.9 GHz Intel Core i5, 8 GB 1867 MHz DDR Memory) and our in-house cluster

1 (Silicon Mechanics, Rackform R4422.v6, two Intel® X®(R) CPU E5-2630 v4 @ 2.20GHz
2 each equipped with ten cores, 128 GB 2400 MHz DDR4 Memory). Hyperthreading allows
3 for the execution of up to forty processes per node. Both environments showed similar
4 scaling, while the MacBook was only tested up to four processes. The cluster results are
5 presented here. Scaling tests were performed with and without splitting the reaction
6 families over several tasks to improve load balancing, see Fig. 8. In test C family splitting
7 allows for almost linear scaling and a parallel efficiency of over 90% up to twenty
8 processes. The parallel efficiency then decreases to 60% for forty processes, probably
9 due to the use of hyperthreading. Without family splitting, the parallel efficiency in test C
10 decreases to about 50% with just four processes and to less than 20% for forty processes.
11 However, due to relatively small reactant sizes and low reactivities in case C, reaction
12 generation is a rather fast process and using family splitting considerably increases the
13 number of tasks to be submitted to the task scheduler. This larger number of tasks causes
14 the test C wall clock time with family splitting to be significantly higher than the wall clock
15 time without family splitting when using one to forty processes. For test B, with medium
16 sized reactants, the wall clock times with or without family splitting are very similar. In test
17 A, with relatively large and reactive species, wall clock times are reduced through family
18 splitting. These observations lead to the heuristic of employing family splitting only for
19 reactants larger than ten non-hydrogen atoms. Finally, using forty processes instead of a
20 single process reduces the wall clock time to at least 15% of its original value. These
21 improvements correspond to speed-up factors ranging between 6.6 and twenty-four. The
22 parallel efficiency was improved with family splitting, but in tests A and B still fell rapidly
23 with only a few processes in use.

1 Similar to the analysis for reaction generation, Fig. 9 presents the elapsed wall clock time
2 and parallel efficiency for the calculation of thermodynamic properties in an isolated test,
3 i.e. not embedded in an RMG simulation.

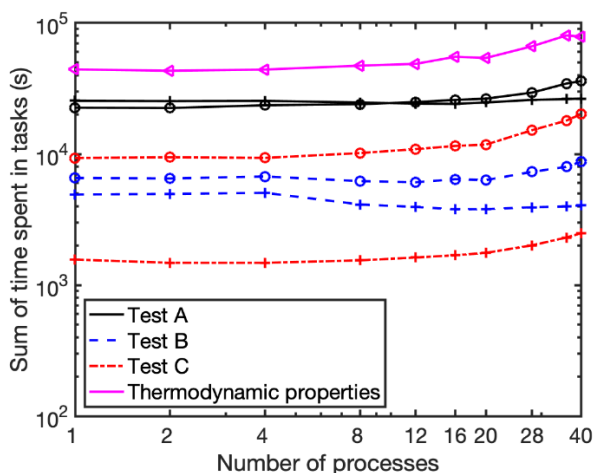


4 **Figure 9: Strong scaling test for the calculation of thermodynamic properties of 390 unique hydrocarbons**
5 **ranging from C₂ to C₁₈ with MOPAC2016 at the PM7 level.**

6 Parallelization allowed for a wall clock time reduction to about 5% of its original value,
7 corresponding to a speed-up factor of eighteen when using forty processes. Only small
8 deviations from linear scaling as well as a smooth decline in parallel efficiency can be
9 observed up to twelve processes. For twelve to forty processes the parallel efficiency
10 decrease is relatively steep, ending below 50%. As in the reaction generation tests, the
11 decrease in parallel efficiency implies that further increasing the number of processes will
12 result in only marginal reductions of the overall wall clock time. To improve our
13 understanding of the decreasing parallel efficiency, two more possible bottlenecks
14 besides load balancing are investigated: contention for shared resources and
15 communication time.

1 4.1.1 Shared resources

2 Shared resources can refer to shared variables within sections of the code designated as
3 critical. Critical sections are defined to ensure shared data is accessed in a serial way,
4 thus avoiding race conditions. As a result, only one process may execute a critical section,
5 or access a shared variable within a critical section simultaneously. Contention for shared
6 resources during the parallel generation of reactions and calculation of thermodynamic
7 properties is assessed by calculating the sum of the time spent in individual tasks as a
8 function of the number of processes. Figure 10 shows that the total time spent in the tasks
9 remains generally constant up to twenty processes, after which it gently increases for
10 most of the tests.



11 **Figure 10: Sum of the times spent in tasks generating reactions or calculating thermodynamic properties as a**
12 **function of the number of processes. Circles represent tests including the splitting of reaction families across**
13 **multiple tasks while crosses represent tests without family splitting.**

14 Based on the results shown in Figure 10, it becomes clear that writing to shared variables
15 within critical sections is not an important factor when analyzing the parallel efficiency
16 decrease. This is likely mainly due to the internal forking algorithm used in Python's
17 multiprocessing package that generates copies of the base process when spawning

1 processes and the fact that for reaction generation in RMG, processes only read the static
2 RMG databases but never write data to them, which would require a locking mechanism
3 to avoid race conditions. Writing to shared variables only occurs when the generated
4 reactions from a process are added to the collection of reactions on the base process or
5 the calculated thermodynamic properties are returned to the molecule. Both actions are
6 handled internally by the `multiprocessing` package.

7 4.1.2 Communication overhead

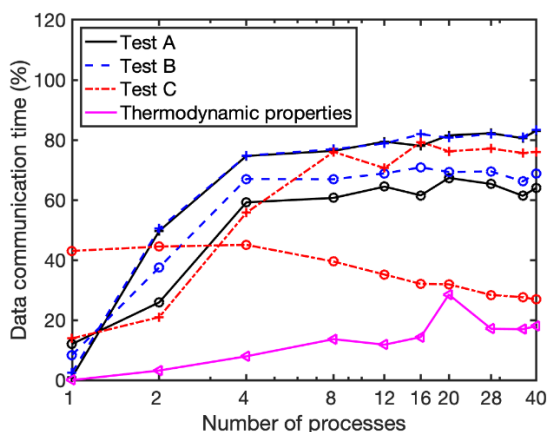
8 In the case of a physical shared resource such as system memory, the off-chip bandwidth
9 puts a limit on the maximum data transfer rate that can be achieved between the system
10 memory and the CPU chip. As discussed before, Python's GIL prevents the simultaneous
11 execution of multiple threads. Therefore, the parallel execution of tasks is achieved by
12 running them through processes with separate memory spaces and, as a result, data
13 communication occurs through message passing. This step requires the serialization of
14 the object hierarchy into a binary representation, the transfer of that byte stream, and
15 finally the inversion operation, or deserialization of the stream into the object hierarchy
16 containing the original data. The current parallel reaction generation scheme requires two
17 communication operations: sending the resonance structure pair from the base process
18 to the parallel processes and sending a list of generated reactions back from the parallel
19 processes to the base process. Since the transferred objects consist of highly layered

1 and complex objects, data communication contributes significantly to the overall wall
2 clock time.

3 The time spent on data communication T_{DC} is calculated with

$$T_{DC} = T_m - \frac{\sum T_{Task}}{m}, \quad \text{Eq. 5}$$

4 where T_{Task} is the time spent in one task, treated in parallel with m other tasks. Therefore,
5 the second term on the right-hand side represents the average sum of time spent in a
6 task. As the time spent on data communication varies for different processes due to the



7 **Figure 11: Time spent on data communication as a contribution to the overall wall clock time and as a function**
8 **of the number of processes. Circles represent tests including the splitting of reaction families across multiple**
9 **tasks while crosses represent tests without family splitting.**

10 distribution of task sizes, T_{DC} represents an average across the processes as opposed to
11 a process-specific metric. The percentage of time spent on data communication as a
12 portion of the overall wall clock time is presented as a function of the number of processes
13 in Fig. 11. The test cases are identical to those used previously and defined in Tab. 2.
14 Generally, for reaction generation, a larger fraction of time is spent on data
15 communication as the number of processes increases, eventually taking a majority of wall
16 clock time. The exception to this is test C with family splitting. Due to the submission of

1 a large number of tasks with relatively short execution times, the data communication
2 time as a percentage of the overall wall clock time is larger for a smaller number of
3 processes. This trend also does not hold for thermodynamic property calculation. In that
4 case, the data communication time in percent is considerably lower as compared to the
5 reaction generation cases.

6 Overall, the results suggest that the drop in parallel efficiency for reaction generation is
7 mainly due to time spent on data communication, which is observed to take up to 80% of
8 the overall wall clock time.

9 4.2 RMG simulations

10 After having assessed the scalability characteristics of the individual components, we now
11 assess their impact on wall clock time and memory consumption for the RMG simulations
12 summarized in Tab. 1. First, we compared reaction filtering, pruning, or multiprocessing
13 in isolation to the base cases. Next, we employed the heuristic algorithm improvements
14 — reaction filtering and pruning — together. Finally, we used all the improvement
15 strategies together. Table 3 presents the obtained wall times for each case.

16 Figure 12.a presents the wall clock time and memory consumption for the diesel oxidation
17 case. Pruning or multiprocessing alone did not have a significant impact on the number
18 of core species in the model after a wall clock time of fourteen days. However, reaction
19 filtering allowed for a larger number of core species. A combination of pruning, reaction
20 filtering and multiprocessing allowed for about twice as many core species to be
21 generated in the same wall clock time as compared to the base case. Steep increases in

22

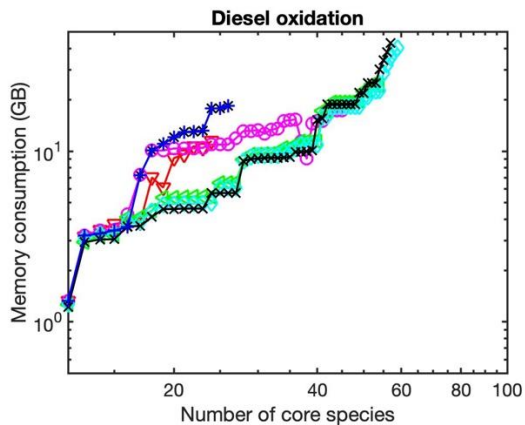
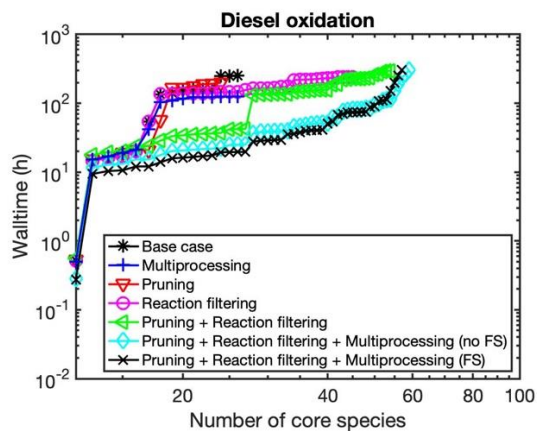
1 Table 3: Wall times for model generation with the suggested changes in RMG. Evaluated at 24, 125, and 218 core species
 2 generated for diesel oxidation, hexadiene pyrolysis, and natural gas combustion, respectively. FS indicates
 3 splitting reaction families over several tasks for reactants containing more than ten non-hydrogen atoms. No
 4 FS indicates that all reaction family matches for one reactant or reactant pair are evaluated in a single task.

Wall times (h)	Base case	Multiprocessing	Pruning	Reaction filtering	Pruning + Reaction filtering	Pruning + Reaction filtering + Multiprocessing	
						No FS	FS
Diesel oxidation	249.6	118.9	199.0	144.9	37.1	22.1	18.7
Hexadiene pyrolysis	95.2	47.4	81.9	16.0	15.8	-	7.4
Natural gas combustion	11.1	9.9	130	1.9	4.5	-	1.9

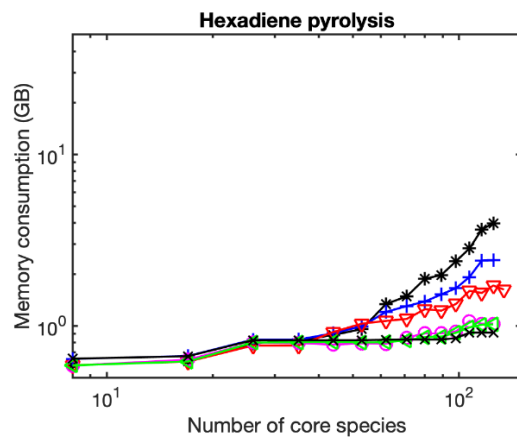
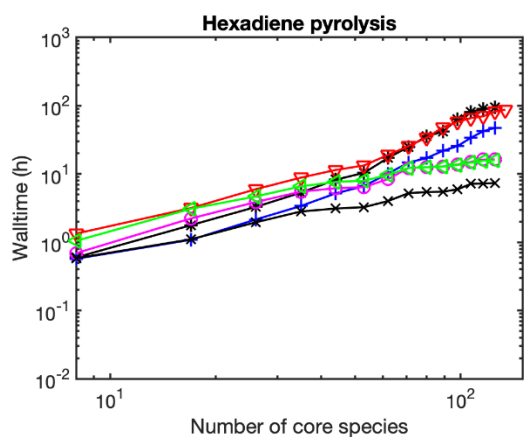
5

6 wall clock time for the addition of a single species were observed even when splitting the
 7 reaction families over several processes. We suspect that this is because a large
 8 computational load is associated with a single reactant pair and reaction family. Such a
 9 task is not divisible under the current family splitting strategy. A further reduction in wall
 10 clock time, especially for the first steep increase, could be obtained by generating more
 11 and smaller tasks for the work done within the reaction generation for one family. This is
 12 also expected to decrease the data communication time and therefore improve the
 13 parallel efficiency for this case. Speed-up factors for the three cases and the suggested
 14 code improvements are reported in Tab. 4. The base cases are also presented in Fig. 2.

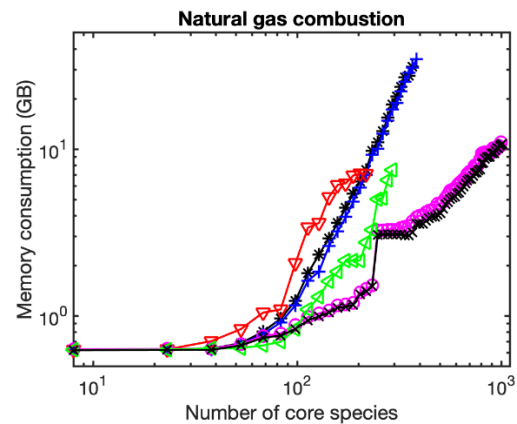
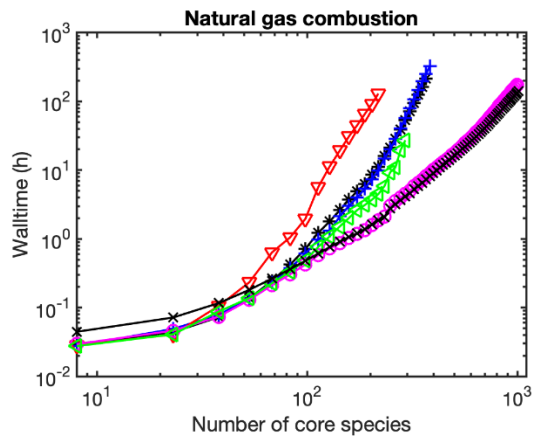
15 Monitoring the memory consumption is important because the forking algorithm limits the
 16 number of available parallel processes by the memory consumption of the base process



1



2



3

4 **Figure 12: Wall clock time and memory consumption as a function of the number of species added to the core**
 5 **employing the discussed scalability strategies either alone or combined for the cases outlined in Tab. 1.**

6

1 and the maximum available memory on the machine. Therefore, when using the parallel
 2 computing strategy, the memory consumption mimics the base case, limiting the number
 3 of allowed processes to below the available maximum of forty. Due to these restrictions
 4 for very memory intense cases, parallel computing might not be possible, causing the
 5 simulation to be carried out in serial. However, the combination of pruning and reaction
 6 filtering allows for a smaller memory footprint compared to the base case. Therefore,
 7 adding the multiprocessing option to these two heuristic strategies allows for better
 8 leveraging the parallel computation strategy. Similar observations can be made for the
 9 hexadiene pyrolysis and natural gas combustion cases presented in Fig. 12.b and c,
 10 respectively.

11 The speed-up between using algorithm heuristics alone and combining them with parallel
 12 computing is still about a factor of two. This seems to be a comparatively small

13

14 **Table 4: Speed-up factor compared with the base cases for the suggested changes in RMG. Evaluated at 24,**
 15 **125, and 218 core species generated for diesel oxidation, hexadiene pyrolysis, and natural gas combustion,**
 16 **respectively. FS indicates splitting reaction families over several tasks for reactants containing more than**
 17 **ten non-hydrogen atoms. No FS indicates that all reaction family matches for one reactant or reactant pair**
 18 **are evaluated in a single task.**

Speed-up	Base case	Multiprocessing	Pruning	Reaction filtering	Pruning + Reaction filtering	Pruning + Reaction filtering + Multiprocessing	
						No FS	FS
Diesel oxidation	1	2.1	1.3	1.7	6.7	11.3	13.3
Hexadiene pyrolysis	1	2.0	1.2	6.0	6.0	-	12.9
Natural gas combustion	1	1.1	0.1	6.0	2.5	-	6.0

19

1 improvement, however, as Amdahl's law (Amdahl, 1967) dictates, the maximum
2 achievable speed-up in an enlargement iteration is governed by the time spent in the
3 serial parts of the code.

4 5 Conclusions and perspectives

5 The complexity of modeling processes driven by chemical kinetics has been a long-
6 standing challenge. With the concerted efforts outlined in this paper, we have identified
7 critical bottlenecks and provided scalability solutions to make computer-aided reaction
8 mechanism construction a tool with practical usefulness in the generation of predictive
9 detailed reaction mechanisms. It was shown that bottlenecks for the overall performance
10 of RMG simulations are dynamic in nature and depend on the state of the simulation and
11 the methods used for the calculation of chemical data. A three-pronged approach
12 consisting of code optimization, algorithm heuristics and parallel computing was
13 employed to accelerate the essential phases of the rate-based enlargement procedure.
14 Tests that measured the individual scalability characteristics have shown speed-up
15 factors with respect to the use of a single process ranging between 6.6 and twenty-four
16 for reaction generation and of about eighteen for thermodynamic property calculation.
17 However, the parallel computing benchmarking also revealed issues that are of
18 fundamental importance to automated mechanism generation tools. The goal of
19 augmenting important entities with metadata competes with scalability and operation
20 speed. Therefore, we believe that future efforts should be directed at devising novel data
21 structures and approaches that fully take advantage of the vast computational resources
22 available today, while creating information-rich models that enable analysis and scrutiny.

1 To conclude, this work demonstrated how the new version 2.4.0 of the Reaction
2 Mechanism Generator opens up new opportunities for the construction of more
3 comprehensive and accurate mechanisms of chemical processes. It also creates
4 avenues for the modeling of real-world processes that were previously too complex to
5 model. RMG is a free-to-use open-source code under the MIT X/11 license. The
6 scalability improvements described in this paper are implemented in the latest version
7 found at <https://github.com/ReactionMechanismGenerator/RMG-Py>. Expecting that the trend
8 towards computer-assisted reaction mechanism generation continues, the identified
9 hurdles and proposed solutions in this work might be applicable in a broader context
10 outside the realm of RMG and may impact other adjacent computational approaches.
11 More specifically, even though each of the available automatic mechanism generators
12 are quite unique in their algorithms they are all expected to benefit from the described
13 strategies when adapted accordingly. Parallelization of thermodynamic property
14 calculation and reaction generation would be helpful in any software where the
15 thermodynamic data or reactions are algorithmically generated.

16 6 Acknowledgements

17 The authors like to thank RMG developers, specifically Richard H. West, Joshua W. Allen
18 and Murat Keceli for helpful discussions. Nick M. Vandewiele acknowledges financial
19 support for a postdoctoral fellowship at MIT from the Belgian American Educational
20 Foundation (BAEF). Agnes Jocher acknowledges financial support from the DFG
21 Research Fellowship JO 1526/1-1. This research was supported by subcontract 7F-
22 30180 to MIT from UC Chicago Argonne LLC. It is a component of the Exascale

1 Computing Project (ECP), Project Number 17-SC-20-SC, a collaborative effort of two
2 DOE organizations, the Office of Science and the National Nuclear Security
3 Administration, responsible for the planning and preparation of a capable exascale
4 ecosystem including software, applications, hardware, advanced system engineering,
5 and early test bed platforms to support the nation's exascale computing imperative

6 Declarations of interest: The authors declare they have no conflicts of interest.

7 References

- 8 Allen, J. W., et al., 2012. Automatic estimation of pressure-dependent rate coefficients. *Phys. Chem. Chem. Phys.* 14,
9 1131–55.
- 10 Allen, J. W., et al., 2014. A coordinated investigation of the combustion chemistry of diisopropyl ketone, a prototype for
11 biofuels produced by endophytic fungi. *Combust. Flame* 161, 711–724.
- 12 Amdahl, G. M., 1967. Validity of the single processor approach to achieving large scale computing capabilities. *Proc.*
13 *AFIPS Spring Jt. Comput. Conf.* 30.
- 14 Behnel, S. et al., 2011. Cython: The Best of Both Worlds. *Comput. Sci. Eng.* 13, 31–39.
- 15 Benson, S. W., 1968. *Thermochemical Kinetics*. 1st, John Wiley & Sons Ltd.
- 16 Blurock, E. S., 1995. Reaction - System for modeling chemical reactions. *J. Chem. Inf. Comput. Sci.* 35.
- 17 Broadbelt, L. J., et al., 1996. Computer generated reaction modelling: Decomposition and encoding algorithms for
18 determining species uniqueness. *Comput. Chem. Eng.* 20, 113–129.
- 19 Burke, M. P., 2016. Harnessing the Combined Power of Theoretical and Experimental Data through Multiscale
20 Informatics. *Int. J. Chem. Kinet.* 48, 212–235.
- 21 Carr, A. G., et al., 2015. Supercritical Water Treatment of Crude Oil and Hexylbenzene: An Experimental and
22 Mechanistic Study on Alkylbenzene Decomposition. *Energy & Fuels* 29, 5290–5302.
- 23 Class, C. A., et al., 2016. Automatic mechanism generation for pyrolysis of di-tert-butyl sulfide. *Phys. Chem. Chem.*
24 *Phys.* 18, 21651–21658.
- 25 Cordella, L. P. P., et al., 2004. A (sub)graph isomorphism algorithm for matching large graphs. *IEEE Trans. Pattern*
26 *Anal. Mach. Intell.* 26, 1367–1372.
- 27 Eyring, H., 1935. The Activated Complex in Chemical Reactions. *J. Chem. Phys.* 3, 107.
- 28 Foggia, P., et al., 2001. A Performance Comparison of Five Algorithms for Graph Isomorphism. *Proceedings of the 3rd*
29 *IAPR TC-15 Workshop on Graph-based Representations in Pattern Recognition.* 188–199.
- 30 Gansner, E. R., North, S. C., 1999. An open graph visualization system and its applications to software engineering.
31 *Pr. Exper* 0, 1–5.
- 32 Gao, C. W., et al., 2015. JP-10 combustion studied with shock tube experiments and modeled with automatic reaction
33 mechanism generation. *Combust. Flame* 162, 3115–3129.

- 1 Gao, C. W., et al., 2016. Reaction Mechanism Generator: Automatic construction of chemical kinetic mechanisms.
2 Comput. Phys. Commun. 203, 212–225.
- 3 Gao, C. W., 2016, Ph.D. dissertation, MIT. <https://dspace.mit.edu/handle/1721.1/104205?show=full>
- 4 Green, W. H., 2007. Predictive Kinetics: A New Approach for the 21st Century. Adv. Chem. Eng. 32, 1–50.
- 5 Gudiyella, S., et al., 2018. Modeling Study of High Temperature Pyrolysis of Natural Gas. Ind. Eng. Chem. Res. 57
6 (22), 7404-7420.
- 7 Han, K., et al., 2017. On-the-Fly Pruning for Rate-Based Reaction Mechanism Generation. Comput. Chem. Eng. 100.
8 1-8.
- 9 Hansen, N., et al., 2013. The Predictive Capability of an Automatically Generated Combustion Chemistry Mechanism;
10 Chemical Structures of Premixed iso-Butanol Flames. 160. 2343-2351.
- 11 Jalan, A., et al., 2010. Predicting solvation energies for kinetic modeling. Annu. Reports Sect. 'C' Phys. Chem. 106,
12 211.
- 13 Jalan, A., et al., 2013. An Extensible Framework for Capturing Solvent Effects in Computer Generated Kinetic Models.
14 J. Phys. Chem. B 117, 2955–2970.
- 15 Jones, E., Oliphant, T. SciPy Open Source Scientific Tools for Python. www.scipy.org. (accessed 18 February 2019).
- 16 Li, S., Petzold, L. R. 1999. Design of New DASPK for Sensitivity Analysis.
- 17 Magoon, G. R., Green, W. H., 2013. Design and implementation of a next-generation software interface for on-the-fly
18 quantum and force field calculations in automated reaction mechanism generation. Comput. Chem. Eng. 52, 35–45.
- 19 Millett, L. I., Fuller, S. H., 2011. Computing Performance: Game Over or Next Level?, in Computer, vol. 44, no. , pp.
20 31-38. doi:10.1109/MC.2011.15
- 21 MOPAC2016, James J. P. Stewart, Stewart Computational Chemistry, Colorado Springs, CO,
22 USA, [HTTP://OpenMOPAC.net](http://OpenMOPAC.net) (accessed 18 February 2019).
- 23 Narayanaswamy, K., et al., 2016. A component library framework for deriving kinetic mechanisms for multi-component
24 fuel surrogates: Application for jet fuel surrogates. Combust. Flame. 165, 288-309.
- 25 Prozument, K., et al., 2014. A Signature of Roaming Dynamics in the Thermal Decomposition of Ethyl Nitrite: Chirped-
26 Pulse Rotational Spectroscopy and Kinetic Modeling. J. Phys. Chem. Lett. 5, 3641–8.
- 27 PyDAS. <https://github.com/jwallen/PyDAS> (accessed 18 February 2019).
- 28 Python 2.7.15 documentation. <https://docs.python.org/2/> (accessed 01 February 2019).
- 29 Rangarajan, S., et al., 2012. Language-oriented rule-based reaction network generation and analysis: Description of
30 RING. Comput. Chem. Eng. 45, 114–123.
- 31 Seyedzadeh Khanshan, F., West, R. H., 2016. Developing detailed kinetic models of syngas production from bio-oil
32 gasification using Reaction Mechanism Generator (RMG). Fuel 163, 25–33.
- 33 Sharma, S., et al., 2007. Automated reaction mechanism generator and applications to model hexadiene doped
34 methane flames. AIChE Annual Meeting, Conference Proceedings.
- 35 Sharma, S., et al., 2010. Modeling of 1-hexadiene, 2,4-hexadiene and 1,4-hexadiene doped methane flames: Flame
36 modeling, Benzene and Styrene formation. Combust. Flame 157, 1331-1345.
- 37 Shi, Y., et al., 2011. Redesigning combustion modeling algorithms for the Graphics Processing Unit (GPU): Chemical
38 kinetic rate evaluation and ordinary differential equation integration. Combust. Flame 158, 836–847.
- 39 Song, J., 2004. Building Robust Chemical Reaction Mechanisms: Next Generation of Automatic Model Construction
40 Software. MIT.
- 41 Suleimanov, Y. V., Green, W. H., 2015. Automated Discovery of Elementary Chemical Reaction Steps Using Freezing

- 1 String and Berny Optimization Methods. *J. Chem. Theory Comput.* 11, 4248–4259.
- 2 Susnow, R. G., et al., 1997. Rate-based construction of kinetic models for complex systems. *J. Phys. Chem. A* 101,
3 3731–3740.
- 4 Van de Vijver, R., et al., 2015. Automatic Mechanism and Kinetic Model Generation for Gas- and Solution-Phase
5 Processes: A Perspective on Best Practices, Recent Advances, and Future Challenges. *Int. J. Chem. Kinet.* 47, 199–
6 231.
- 7 van der Walt, S., et al., 2011. The NumPy Array: A Structure for Efficient Numerical Computation. *Comput. Sci. Eng.*
8 13, 22–30.
- 9 Vandewiele, N. M., et al., 2012. Genesys: Kinetic model construction using chemo-informatics. *Chem. Eng. J.* 207,
10 526–538.
- 11 Vandewiele, N. M., et al., 2015. Kinetic Modeling of Jet Propellant-10 Pyrolysis. *Energy Fuels* 29, 1, 413-427.
- 12 Vermeire, F. H., et al., 2017. Experimental and kinetic modeling study of the pyrolysis and oxidation of 1,5-hexadiene:
13 The reactivity of allylic radicals and their role in the formation of aromatics. *Fuel* 208, 779-790.
- 14 Warth, V., et al., 2000. Computer based generation of reaction mechanisms for gas-phase oxidation. *Comput. Chem.*
15 24, 541–560.
- 16 Zhang, P., et al., 2018. Modeling Study of the anti-knock tendency of substituted phenols as additives: An application
17 of the Reaction Mechanism Generator (RMG). *Phys. Chem. Chem. Phys.* 20, 10637–10649.