

# Mutexa: A Computational Ecosystem for Intelligent Protein Engineering

Zhongyue J. Yang<sup>1-5,\*</sup>, Qianzhen Shao<sup>1</sup>, Yaoyukun Jiang<sup>1</sup>, Christopher Jurich<sup>1,3</sup>, Xinchun Ran<sup>1</sup>,  
Reecan J. Juarez<sup>1,6</sup>, Bailu Yan<sup>7</sup>, Sebastian L. Stull<sup>1</sup>, Anvita Gollu<sup>1</sup>, Ning Ding<sup>1</sup>

<sup>1</sup>*Department of Chemistry, Vanderbilt University, Nashville, Tennessee 37235, United States*

<sup>2</sup>*Center for Structural Biology, Vanderbilt University, Nashville, Tennessee 37235, United States*

<sup>3</sup>*Vanderbilt Institute of Chemical Biology, Vanderbilt University, Nashville, Tennessee 37235,*

*United States* <sup>4</sup>*Department of Chemical and Biomolecular Engineering, Vanderbilt University,*

*Nashville, Tennessee 37235, United States* <sup>5</sup>*Data Science Institute, Vanderbilt University,*

*Nashville, Tennessee 37235, United States* <sup>6</sup>*Chemical and Physical Biology Program, Vanderbilt*

*University, Nashville, Tennessee 37235, United States* <sup>7</sup>*Department of Biostatistics, Vanderbilt*

*University, Nashville, Tennessee, 37205, United States*

## **Corresponding Author**

\*Email: zhongyue.yang@vanderbilt.edu Phone: 615-343-9849

ABSTRACT. Protein engineering holds immense promise in shaping the future of biomedicine and biotechnology. This review focuses on our ongoing development of Mutexa, a computational ecosystem designed to enable "intelligent protein engineering". In this vision, researchers can seamlessly acquire sequences of protein variants with desired functions as biocatalysts, therapeutic

peptides, and diagnostic proteins by interacting with a computational machine, similar to how we use Amazon Alexa in these days. The technical foundation of Mutexa has been established through the development of database that integrates enzyme structures with their respective functions (e.g., IntEnzyDB), workflow software packages that enable high-throughput protein modeling (e.g., EnzyHTP and LassoHTP), and scoring functions that map the sequence-structure-function relationship of proteins (e.g., EnzyKR and DeepLasso). We will showcase the applications of these tools in benchmarking the convergence conditions of enzyme functional descriptors across mutants, investigating protein electrostatics and cavity distributions in SAM-dependent methyltransferases, and understanding the role of non-electrostatic dynamic effects in enzyme catalysis. Finally, we will conclude by addressing the future steps and challenges in our endeavor to develop new Mutexa applications that facilitate the selection of beneficial mutants in enzyme engineering.

## **KEYWORDS**

High-throughput computation, intelligent protein engineering, enzyme simulation, predictive modeling

## **1. Introduction**

Protein engineering refers to the process of optimizing protein sequences for enhanced physical (e.g., thermal stability, solubility, and complex stoichiometry), chemical (e.g., reactivity, substrate specificity, selectivity, and substrate scope), biological, and pharmaceutical functions. Typical strategies in protein engineering include directed evolution,<sup>1-4</sup> gene shuffling/recombination,<sup>5, 6</sup> site-directed mutagenesis,<sup>7, 8</sup> and protein truncation and fusion.<sup>9, 10</sup> Enabled by protein engineering, researchers can create enzymes to transform difficult<sup>11-14</sup> or even

new-to-nature reactions<sup>15, 16</sup>, develop peptides with targeted therapeutic effects,<sup>17, 18</sup> innovate diagnostic tools for early-stage cancer detection,<sup>19-21</sup> and advance our understanding of fundamental life processes.<sup>22, 23</sup>

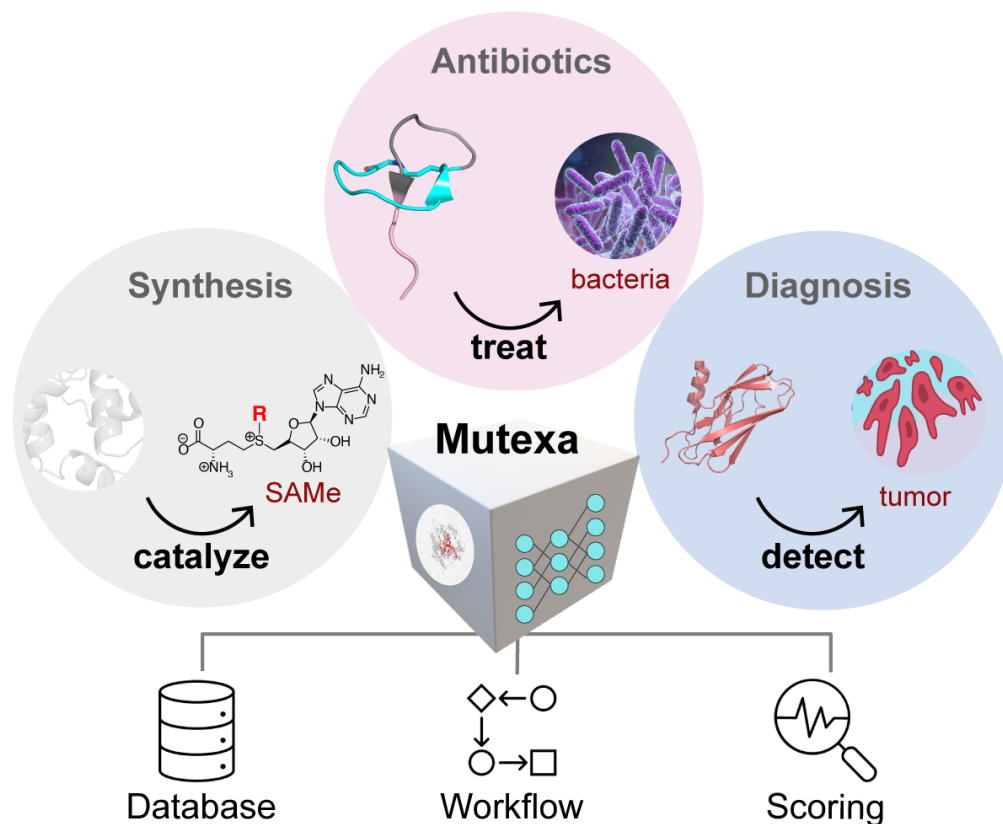
A “holy grail” challenge in protein engineering is the effective identification of desired protein variants within a mutation landscape.<sup>24, 25</sup> This difficulty results from the combinatorial explosion associated with sequence mutation. Sampling mutations across only a dozen amino acid sites creates an astronomical number of variants. Despite the advances of screening strategies for protein engineering, the success rate for identifying beneficial mutants is around 1% or lower.<sup>26-33</sup> *De novo* design of new functional proteins provides a promising alternative, but the hit rate to identify successful designs among all design candidates is similar to the chance of experimental discovery.<sup>34-37</sup> The time-consuming, labor-intensive, and expensive process of experimental screening is largely unavoidable.

To reduce the size of mutant libraries for functional screening, computational approaches have been augmented with protein engineering.<sup>25, 38-40</sup> These methods, such as bioinformatics,<sup>24, 41</sup> classical molecular simulations,<sup>42, 43</sup> quantum chemistry,<sup>44-47</sup> and data-driven modeling,<sup>22, 48-52</sup> span over a wide breadth of computational sub-fields. Each type of the modeling strategy has a specific strength. Bioinformatics reveals the evolutionary coupling and pattern behind the function-encoding sequence spots; classical molecular simulation elucidates the dynamics and conformational ensembles that constitute effective protein-protein/ligand interactions or enzyme catalysis; quantum chemistry informs the variation of electronic structure that underlies enzymatic reactions or covalent inhibition; and data-driven modeling predicts the formal, non-linear relationships between sequence, structure, and function. Each type of these computational methods may fall short in accuracy, efficiency, resolution, or reproducibility. The combination of these

computational approaches shows a great promise to establish an integrative strategy that we call “intelligent protein engineering”. Intelligent protein engineering aims to guide experimental discovery of desired protein mutants by effectively shrinking the sheer number of mutations that have to be screened. Intelligent protein engineering has the potential to save extensive amount of experimental efforts for identification of functional protein variants.

With a long-term goal to create a platform that enables intelligent protein engineering, our lab has been building a computational ecosystem called Mutexa. Mutexa is short for “Alexa for mutants”, and we believe that how people engineer proteins in the future should be similar to the way we use Amazon Alexa in these days – if researchers intend to obtain the sequences of protein variants with desired functions, they just need to ask for help from a computational machine. Mutexa integrates high-throughput computation, bioinformatics, quantum chemistry, multiscale simulation, and data-driven modeling to identify protein mutants that can enhance functions including enzyme catalysis, peptide therapeutics, and disease biomarker detection.<sup>23</sup> Over the past three years, my lab has been establishing the technical foundation of Mutexa by developing 1) a database that integrates enzyme structure and function data (IntEnzyDB<sup>53, 54</sup>), 2) software tools for high-throughput construction and modeling of enzymes (EnzyHTP<sup>55, 56</sup>) and lasso peptides (LassoHTP<sup>55</sup>), and 3) scoring functions to predict the impact of mutations on substrate-positioning dynamics,<sup>23, 57</sup> enzymatic kinetic resolution (EnzyKR<sup>58</sup>), and peptide antimicrobial activity (DeepLasso<sup>59</sup>). The database, workflow software, and scoring functions will be discussed in detail in Sections 2, 3, and 4, respectively. In addition, we will briefly introduce the applications of these tools to investigate the convergence in determining enzyme functional descriptors across enzyme mutants,<sup>60</sup> distribution of protein electrostatics and cavity for SAM-dependent methyltransferases,<sup>61</sup> and understanding of non-electrostatic dynamic effects in mediating enzyme

catalysis.<sup>62</sup> Finally, we will conclude by addressing the next steps and challenges in building new Mutexa applications for biocatalyst development.



**Scheme 1.** Overview of Mutexa, a computational ecosystem for protein engineering. Mutexa consists of three components, including a database that integrates structure and function information of proteins, a workflow software that allows automatic, high-throughput modeling for proteins, and a scoring function that describes sequence-structure-function relationship of proteins. Combining the three basic components, new applications for predictive modeling are being developed into Mutexa, including tools that enable enzyme engineering for non-native substrates or new-to-nature reactions, peptide engineering for antimicrobial uses, and binder protein engineering for disease biomarker recognition.

## 2. IntEnzyDB: an Integrated Structure-Function Enzymology Database

Building an integrated database that merges related enzyme sequence, structure, and function data in one place is essential for developing accurate physical methods and holistic data-driven models for enzyme engineering. However, data collection, cleaning, and joining present as

three major challenges. Data collection is often impeded by different design (e.g., relational, object-oriented, or hybrid), storage hierarchy, query mechanism, and API protocols of various types of existing databases. Data cleaning is tricky because existing data entries involve missing or inaccurate mutational spot labels and experimental conditions, as well as manual typos and rounding errors. Data joining between enzyme structure and function data is challenging due to inconsistent keys – enzyme kinetics databases typically store data entries by EC number and often lack PDB IDs, causing barriers for one-to-one mapping to structural databases.

To address these challenges, my lab developed an integrated structure-kinetics enzymology database, IntEnzyDB, for facile data-driven modeling and machine learning.<sup>53, 54</sup> The database merges related enzyme sequence, structure, and function data in one place to address the challenges associated with the collection, cleaning, and joining of enzymology data. In contrast to object-oriented databases that store enzyme records in separate data files,<sup>63</sup> IntEnzyDB employs a relational database architecture with a flattened data structure. This approach enhances scalability and enables the integration of additional enzyme function data, such as folding stability and solubility, into the database. Noticeably, Fleischmann et al. has employed relational architecture to build IntEnz, which is an integrated enzymology database for nomenclature and classification of enzyme-catalyzed reactions.<sup>64</sup>

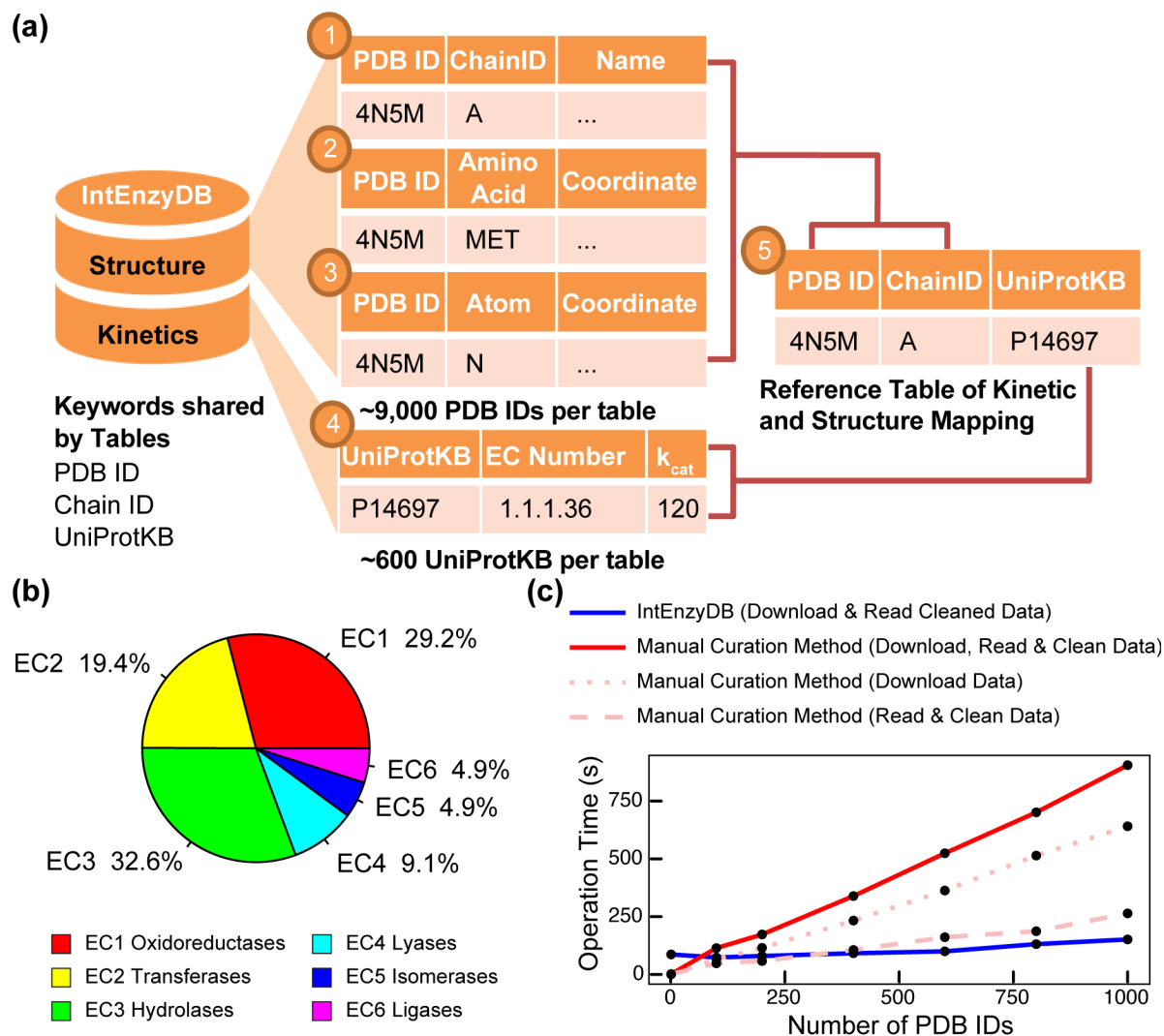
To store enzyme kinetics and structure information, IntEnzyDB implements five data tables (Figure 1a). We curated three tables to store cleaned enzyme structure data derived from RCSB PDB,<sup>63</sup> including a chain table that contains general protein structure information (e.g., nomenclature, organism, resolution, etc.), an amino acid table that contains amino acid attributes, properties, and physiochemical parameters, and an atomic structure table that contains the atom types and Cartesian coordinates. We curated one table for kinetics data derived from BRENDA,<sup>65</sup>

Sabio-RK,<sup>66</sup> ProtaBank,<sup>67</sup> and Design2Data.<sup>68</sup> The table contains information of enzyme kinetic assays such as EC number, substrate, mutation, temperature, turnover number, Michaelis constant, and so on. Finally, we curated one reference table to achieve one-to-one mapping between enzyme kinetics and PDB based on foreign keys PDB ID, Chain ID, and UniProtKB. Using IntEnzyDB, we created a data table comprising 4243  $k_{cat}/K_M$  values for enzymes with single amino acid substitutions. The dataset includes 691 wild-type (WT) enzymes, 2592 enzyme mutants, and 943 substrates. Of the stored  $k_{cat}/K_M$  values, 29.2% pertain to oxidoreductases (EC 1), 19.4% to transferases (EC 2), 32.6% to hydrolases (EC 3), 9.1% to ligases (EC 4), 4.9% to isomerases (EC 5), and 4.9% to lyases (EC 6) (Figure 1b).

To assess the efficiency of retrieving enzyme structure data using IntEnzyDB, we compared it against a manual curation strategy (Figure 1c). Unlike the manual approach, IntEnzyDB allows the user to filter and download pre-cleaned and tabulated structural data directly using SQL queries. Our results indicate that for processing 200 enzymes, IntEnzyDB is approximately two times faster than the manual curation approach (80 s vs 173 s), and for 1000 enzymes, it is around six times faster (151 s vs 905 s). The results indicate that the operating time using IntEnzyDB is nearly independent of the data size, which is particularly beneficial when handling a large amount of structural data (i.e., thousands or more). The flattened data structure of IntEnzyDB likely accounts for its high data processing efficiency. By loading all data entries at once, IntEnzyDB outperforms the traditional approach, where data tables and files are accessed serially in CPU. While processing smaller amounts of data (e.g., for one enzyme structure), IntEnzyDB may take longer (86 s vs 1.9 s) than the manual approach. However, IntEnzyDB can save a substantial amount of time when handling large amounts of structural data by avoiding repeatedly opening and reading files.

Although only ~10 mins are saved when operating on the 1000 structures in the benchmark (Figure 1), time savings are expected to proportionally increase with the number of data entries. We expect that more quality data of enzyme structure and function will be collected and stored in coming years. As such, IntEnzyDB provides an efficient solution for extracting enzyme structural features for statistical analysis or machine learning. The quality structure and function data stored in IntEnzyDB also provide benchmark sets for systematic assessment and development of new molecular modeling methods used in enzyme engineering. As the next steps for developing IntEnzyDB, we will further expand the mapped structure-kinetics data table by using predicted structures and active site annotation. Text mining strategies will be implemented to enable more comprehensive data validation and expansion. We will incorporate new types of enzymology data to IntEnzyDB, including stability, solubility, expressibility, and even molecular modeling data derived from high-throughput simulations.<sup>56</sup> The incorporation of a diverse range of quality data from molecular level to macroscopic scale has the potential to enhance the learning efficiency, predictive accuracy, and generalizability of the models.





**Figure 1.** The architecture, kinetics data statistics, and performance benchmark for IntEnzyDB. (a) The database architecture is based on five tables, including three tables for enzyme structure information (chain-level, amino acid-level, and atom-level), one table for enzyme kinetics, and a reference table that includes foreign keys from the structure and kinetics tables. The mapping of the tables is established using the PDB ID, Chain ID, and UniProtKB keys. (b) The distribution of kinetics data for six enzyme commission classes. (c) The comparison of operation time between IntEnzyDB and manual curation methods. The operation time for downloading, reading, and cleaning data is measured for processing 1, 100, 200, 400, 600, 800, and 1000 PDB IDs, with data downloading and reading/cleaning indicated by dotted and dashed lines, respectively. The total operation time for the manual curation method is shown by the red solid line. All operation times are reported in seconds.

### 3. Software Tools that Enable High-throughput Molecular Simulations of Proteins

#### 3.1 EnzyHTP: a High-Throughput Computational Platform for Enzyme Modeling

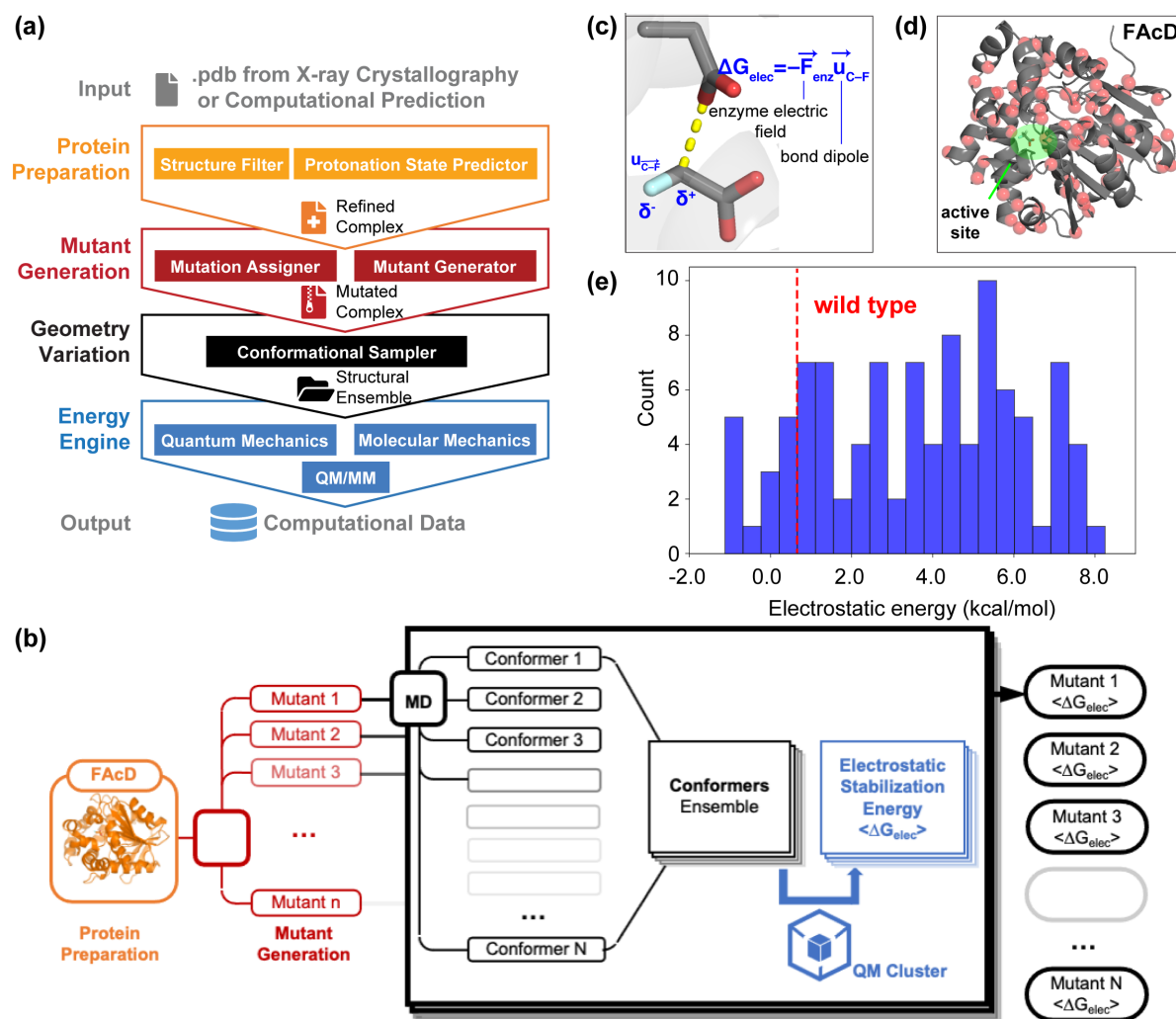
Different types of computational theories and methods, including quantum mechanics (QM), molecular mechanics (MM), and multiscale QM/MM modeling, have been extensively employed in protein engineering to guide the selection of function-enhancing enzyme mutants for late-stage functionalization<sup>69</sup>, polymer upcycling<sup>70, 71</sup>, degradation of environmental pollutants<sup>72, 73</sup>, and treatment of food allergies<sup>74-46</sup>. To maximize the potential of molecular simulations in biocatalyst development,<sup>75-78</sup> it is essential to perform enzyme modeling in an automatic and high-throughput fashion. To address this challenge, my lab developed a computational platform, EnzyHTP, to automate the entire life cycle of enzyme modeling in a high-throughput manner. EnzyHTP has four levels of operation arranged in a top-down hierarchy (Figure 2a). The four levels are protein preparation, mutant generation, geometry variation, and energy engine. Each level was implemented as an independent Python module. The protein preparation module emphasizes constructing computational models for enzyme structures obtained from X-ray crystallography experimental data or computational predictions. The mutant generation module is responsible for generating novel enzyme variants based on a common enzyme sequence and scaffold by altering an existing amino acid's sidechain type and conformation. The geometry variation module samples enzyme conformation and substrate reaction coordinates using external molecular dynamics or Monte Carlo software packages. The energy engine makes use of QM, MM, or multiscale QM/MM calculations using quantum chemistry toolboxes. In particular, the QM treatment of enzyme active sites and reacting species is critical to elucidating the catalytic mechanisms of enzymes and predicting the impact of mutations on enzyme catalysis.

To demonstrate the high-throughput capability of EnzyHTP, we employed the software to investigate the impact of single mutations on the interior enzyme electrostatics for 100 fluoroacetate dehalogenase (FAcD) mutants (Figure 2b). The model enzyme, *Rhodopseudomonas*

*palustris* FAcD, hydrolyzes the C–F bond of fluoroacetate (FAc) via an  $S_N2$  mechanism (Figure 2c).<sup>79-83</sup> The cleavage of the C–F bond contributes to the rate-determining step, and enzyme electric field accelerates the reaction by stabilizing the dipole moment along the breaking C–F bond.<sup>84</sup> The electrostatic effect is quantified using electrostatic stabilization energy (i.e.,  $\Delta G_{\text{elec}}$ ), which is computed by the dot product between the electric field and the C–F bond dipole (Figure 2c). Using EnzyHTP, we created a Python workflow to compute  $\Delta G_{\text{elec}}$  values for 100 FAcD variants with random single amino acid substitution. The workflow first generates 100 variants using the mutant generation module based on a curated FAcD crystal structure (Figure 2b). The mutation spots are distributed over the entire FAcD enzyme scaffold (Figure 2d), with a spatial proximity to the active site ranging from 7 Å to 32 Å. The workflow performs an MD simulation for each variant and samples 100 conformers from a 1 ns MD production run. The structure involves a restrained pre-reaction complex in which the residue Asp<sup>110</sup> is aligned with the substrate C–F bond for a potential  $S_N2$  attack. A short propagation time is used for the MD simulations to ensure that the sampled enzyme conformers resemble the crystal structure. Third, the workflow computes the ensemble average of  $\Delta G_{\text{elec}}$  values (denoted by  $\langle \Delta G_{\text{elec}} \rangle$ ) using 100 conformational snapshots extracted from a 1 ns MD trajectory. The bond dipole is computed using a single-point QM calculation (HF/6-31G(d)) that consists of the substrate and Asp<sup>110</sup>, followed by the wavefunction-based localized molecular orbital (LMO) analysis using Multiwfn. The electronic field strength of a mutant is computed based on the RESP charges of enzyme atoms using Coulomb’s law. Solvent molecules and counterions are excluded. Using the workflow, we completed the computation of  $\langle \Delta G_{\text{elec}} \rangle$  values for 100 FAcD variants in 7 hours with 10 GPUs (NVIDIA V100 SMX2) and 160 CPUs (Xeon Gold 6248). In contrast, performing this process manually for 100 enzyme variants would

take several weeks due to tedious processes of mutant structure curation and file preparation, in addition to the computational runtime.

Figure 2e displays the distribution of  $\langle \Delta G_{\text{elec}} \rangle$  values for 100 FAcD variants. The computed  $\langle \Delta G_{\text{elec}} \rangle$  values exhibit a range of -1.1 kcal/mol to 8.2 kcal/mol. Comparing to the reference  $\langle \Delta G_{\text{elec}} \rangle$  value of the WT FAcD (i.e., 0.5 kcal/mol), a small proportion of mutations (~10%) cause a reduction in the  $\langle \Delta G_{\text{elec}} \rangle$  value, indicating the formation of a more favorable electrostatic environment that can between stabilize the developing C–F dipole in the FAcD mutant compared to the WT FAcD. However, the majority of mutations (~90%) have the opposite effect, which are likely to reduce or even abolish the catalytic effect. Despite an enhanced enzyme electric field strength for breaking the C–F bond, the 10% mutations are not necessarily the actual beneficial mutations due to the impact of mutation on other untested aspects, such as stability, solubility, expressibility, and so on. Our work on developing EnzyHTP software sets the basis for *in silico* high-throughput enzyme screening that identifies beneficial enzyme variants, which can accelerate the development cycle of new biocatalysts that catalyze non-native substrates or new-to-nature reactions. EnzyHTP will facilitate the comprehension of enzyme catalytic mechanisms across numerous enzymes within a protein family. EnzyHTP can also help generate computational data for our database IntEnzyDB that guides future statistical understanding and machine learning. Inspired by the code base and architecture of EnzyHTP, we are developing more high-throughput software suites to address specific challenges of automatic molecular modeling in protein engineering. For one, we developed a tool for automatic construction and modeling of lasso peptides. This will be discussed in Section 3.2.



**Figure 2.** The design framework and application of EnzyHTP. (a) The workflow of high-throughput enzyme modeling. The workflow comprises four levels of operation, namely protein preparation, mutant generation, geometry variation, and energy engine. The input to this framework is the enzyme structure, and the output is computational modeling data. (b) Application of EnzyHTP to compute the electrostatic stabilization energy values (i.e.,  $\langle \Delta G_{\text{elec}} \rangle$ ) for 100 fluoroacetate dehalogenase (FAcD) mutants. For each mutant, the workflow automatically conducts 1 ns molecular dynamics simulations, 100 single point quantum mechanics calculations, dipole moment analysis, and output an averaged  $\langle \Delta G_{\text{elec}} \rangle$  value. (c) Definition of electrostatic stabilization energy, which is computed as the dot product between the enzyme interior electric field and the dipole moment of the breaking C-F bond. (d) Spatial distribution of 100 single mutation spots on FAcD. (e) The distribution of  $\Delta G_{\text{elec}}$  values for 100 FAcD mutants, where the red dashed line indicates the  $\Delta G_{\text{elec}}$  value for the WT FAcD.

### 3.2 LassoHTP: a High-Throughput Tool for Lasso Peptide Structure Construction and Modeling

Lasso peptides are a class of ribosomally synthesized and post-translationally modified natural products. They were first discovered in 1991,<sup>85</sup> and have been increasingly reported as candidates for new antibiotics,<sup>86-89</sup> enzyme inhibitors,<sup>87, 90</sup> and receptor antagonists,<sup>85</sup> (e.g., microcin J25<sup>90, 91</sup>). Lasso peptides involve a 1-rotaxane topology<sup>92, 93</sup> with a macrolactam ring held in position by sterically bulky residues above and below the ring. The ring in the lasso peptide is formed by an isopeptide bond between the N-terminal  $\alpha$ -amino group and the carboxylate group of an aspartate or glutamate. Bioinformatic analyses estimate that the lasso peptides with a known structure and function occupy ~10% of all possible lasso peptides that exist in nature. To accelerate the discovery of functional lasso peptides, computational tools that allow the prediction of structures and functions of uncharacterized lasso peptides will help to prioritize pharmaceutically valuable lasso peptides for experimental assessments. However, due to the distinct topology of lasso peptides, computational tools that were designed for structural prediction of globular proteins (e.g., AlphaFold2<sup>94</sup>) or cyclic peptides<sup>95</sup> fail to inform the structure of lasso peptides with high fidelity.

To address this challenge, my lab developed LassoHTP as a tool for high-throughput lasso peptide structure prediction and conformational sampling. LassoHTP converts a user-defined lasso peptide sequence (with annotation of ring, loop, and tail) into a three-dimensional structure and a conformational ensemble using three software modules, including a scaffold constructor, a mutant generator, and an MD simulator (Figure 3a). The scaffold constructor is responsible for generating a poly-alanine lasso peptide scaffold based on a structural library and tail extender function, while the mutant generator module mutates this scaffold to produce a lasso peptide structure that

corresponds to the user-defined sequence or sequences resulting from mutagenesis (Figure 3b). Finally, the MD simulator uses the AMBER software package<sup>96</sup> to parameterize the resulting lasso peptide structure and conduct MD simulations to output a conformational ensemble. The modular architecture of LassoHTP ensures its flexibility and versatility, similar to that of EnzyHTP.<sup>56</sup> Each module can be independently operated for building, modifying, or modeling a lasso peptide, and the three modules can be sequentially executed as part of an automatic workflow to convert user-defined lasso peptide sequences into conformational ensembles.

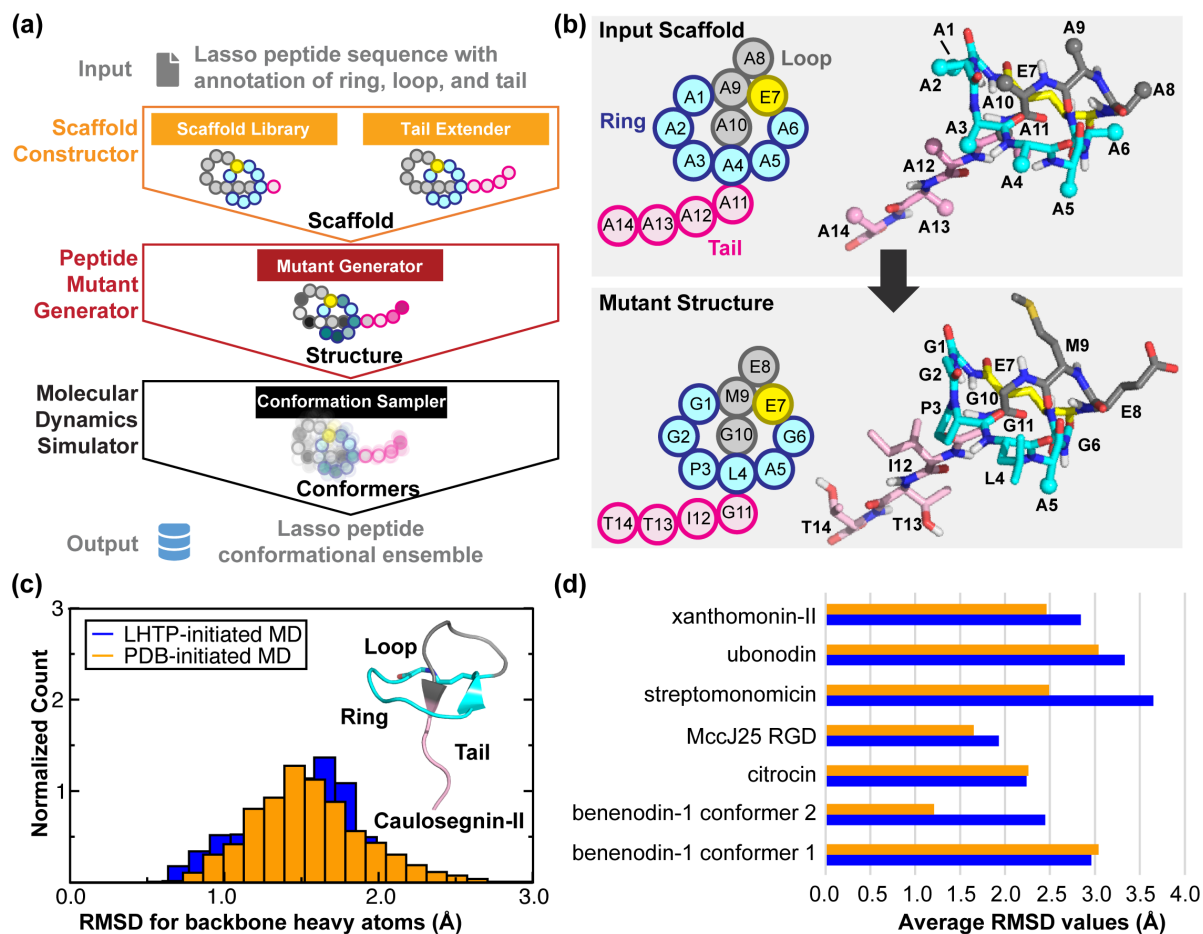
To test LassoHTP, we employed the software to predict conformational ensembles for different types of lasso peptides (called LHTP-initiated MD) and then benchmarked their consistency against the MD ensembles initiated from the corresponding crystal- or NMR-structures (called PDB-initiated MD, Figure 3c and 3d). The first test case is the WT caulosegnin-II<sup>97</sup> (PDB ID: 5D9E), as a crystal structure (resolution: 0.86 Å) exists for this peptide. For both LHTP-initiated and PDB-initiated MD ensembles, trajectories were simulated using identical force field parameters and the ensembles were constructed by evenly taking 1000 snapshots from a 100 ns MD trajectory. The RMSD value calculated from the LHTP-initiated ensemble (1.48 Å) closely align with that from the PDB-initiated ensemble (1.55 Å).

Furthermore, we tested LassoHTP using seven lasso peptides whose structures have been determined by NMR, including benenodin-1 conformer 1,<sup>98</sup> benenodin-1 conformer 2,<sup>98</sup> citrocin,<sup>99</sup> the RGD variant of microcin J25,<sup>100</sup> streptomomicin,<sup>101</sup> ubonodin,<sup>102, 103</sup> and xanthomonin-II<sup>104</sup> (Figure 3d). They involve a wide range of structural constructs. The first structural model of each peptide's NMR-resolved structural ensemble was used to initiate the MD simulation and as a reference structure for RMSD calculations in both LHTP- and PDB-initiated MD ensembles. The two ensembles are reasonably consistent: the difference of the RMSD values between the two

ensembles ranges from  $\sim 0.0$  Å for benenodin-1 conformer 1 and citrocin to  $\sim 1.2$  Å for streptomomicin and benenodin-1 conformer 2, with the average being 0.48 Å. The consistency between the two ensembles were also validated using principal component analysis (PCA). The benchmark shows that LassoHTP can generate reasonable lasso peptide structures and conformational ensembles from sequence. As such, LassoHTP provides a platform to build modules for high-throughput functional predictions including binding affinity to drug target, thermostability against harsh conditions, and permeability across membrane transport proteins.

Nonetheless, we should note some technical limitations that we would like to address in LassoHTP. For one, the isopeptide bonds with a *cis*-configuration, which populate with high abundance in benenodin-1,<sup>105</sup> have not been constructed in the scaffold library. Additionally, enhanced sampling methods have yet to be used for navigating the conformational space of lasso peptides. Both aspects are expected to be addressed in the next version of LassoHTP.





**Figure 3.** The design framework and application of LassoHTP. (a) A schematic of LassoHTP's workflow, which involves three modules: scaffold constructor, peptide mutant generator, and MD simulator, to transform a user-input sequence into a conformational ensemble. (b) Application of the mutant generator module to convert the poly-alanine lasso peptide scaffold into the lasso peptide that is consistent with the user-input sequence. Sequence shown is for xanthomonin-II<sup>104</sup> (PDB ID: 2MFV). (c) Distribution of root mean square deviation (RMSD) for LassoHTP-initiated and PDB-initiated MD conformational ensemble for caulosegnin-II<sup>97</sup>. (d) Average RMSD values of LassoHTP (LHTP)-initiated (colored in blue) and PDB-initiated (colored in orange) MD ensembles for eight lasso peptides involved in the benchmark. The structures of the lasso peptides were determined mostly by NMR except for caulosegnin-II by X-ray crystallography (PDB ID: 5D9E). For (c) and (d), the RMSD was calculated using backbone atoms (i.e., C<sub>α</sub>, N, C, and O) with reference to the reference crystal and NMR structure.

### 3.3 ARMer: A Python Library for Adaptive Resource Allocation of Molecular Modeling

#### Workflows on High Performance Computing Clusters

High-throughput computation emerges as a new paradigm to facilitate mechanistic study,<sup>106</sup> catalyst screening,<sup>107</sup> functional material design,<sup>108, 109</sup> drug discovery,<sup>110, 111</sup> and enzyme

modeling.<sup>56</sup> Our lab has developed EnzyHTP<sup>56</sup> and LassoHTP<sup>55</sup> as open-access software packages to enable the high-throughput modeling of enzymes and lasso peptides, respectively. High-throughput computation needs to allocate different types of computing resources (e.g., CPU, GPU, etc.) for multiple sub-tasks in high-performance computing (HPC) clusters. Resource allocation in the workflow to minimize resource and time waste remains a technical challenge in the computational community. To address this challenge, we developed a new Python library, adaptive resource manager (ARMer), to dynamically request computing resources based on the need of a specific modeling sub-task in the workflow. Using commands implemented in the ARMer library, a Python “workflow script” is prepared that runs on a single-CPU thread to configure, submit, and monitor molecular simulation jobs for a high-throughput workflow in HPC clusters. This is in sharp contrast to the traditional resource allocation scheme where a fixed amount of computing resources is requested for all types of molecular modeling tasks.

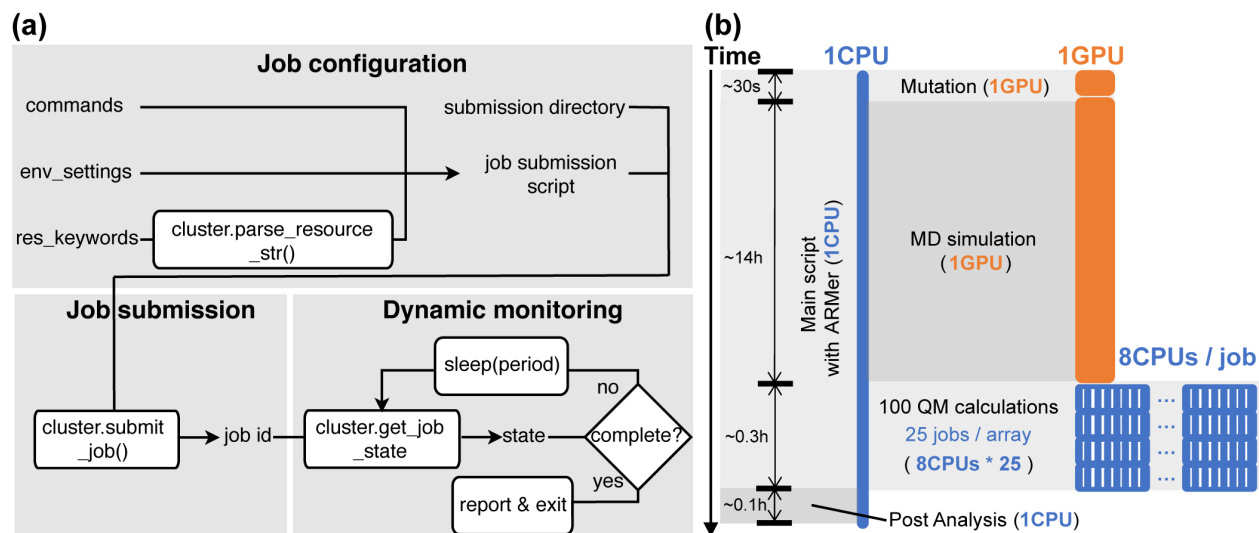
The ARMer Python library contains a Job class that defines variables and functions associated with job configuration, submission, and dynamic monitoring of job completion (Figure 4a). ARMer also contains an HPC class that supports the Job class with variables and functions to mediate external input/output in a local HPC cluster where ARMer is deployed. In the Job class, a job object is instantiated based on information provided by the user through the arguments: *commands*, *cluster*, *env\_settings*, and *res\_keywords*. With the Job object created, a script for the required tasks can be generated and then submitted by the *submit()* method (Figure 4a). A job ID is added to the object by the function. By tracing the job ID, the “workflow script” can monitor the status of a job object in the queue, and mediate the status by killing, holding, or releasing the job. The “workflow script” will dynamically detect the timing of the job completion by retrieving error or completion messages from the output file. Notably, dynamic monitoring of job completion

status is critical to a high-throughput modeling workflow because multiple types of simulation sub-tasks must be sequentially operated in the process. In the case of high-throughput enzyme modeling, after submitting an MD sampling task, the “workflow script” must put the rest of the sub-tasks on hold and wait for the conformational ensemble to generate before submission of the subsequent QM calculations.

We tested the resource and time consumption of the high-throughput molecular modeling workflow enabled by adaptive resource allocation on our local HPC at Vanderbilt’s advanced computing center for research and education (ACCRE). A single-CPU job was submitted to execute a “workflow script” that employs built-in commands from the ARMer library to manage computing resources for molecular simulation tasks involved in the high-throughput modeling of FAcD mutants (Figure 4b). Compared to traditional allocation strategy that directly execute sub-tasks using a fixed amount of allocated CPU or GPU nodes, this Python script configures resource-demanding sub-tasks (i.e., needing  $>1$  CPU or  $\geq 1$  GPU) in a new job script and then submits the job to the queue (i.e., setting *ifcluster* = ‘True’ in the code). This job was set with a 96-hour wall-clock running time so that it can oversee the entire workflow.

For the MD simulation task, the workflow script configures shell commands in a job script to request GPU resource, set environment variables, and conduct MD modeling using AMBER. The workflow script then submits the job and regularly monitors the completion status of the job. After receiving the signal of completion, the workflow script will continue operating the QM calculation sub-tasks in the workflow. Due to the independence of individual QM tasks, the workflow script can submit multiple QM jobs (8 CPU for each QM job) simultaneously to the job array so that they can run in parallel up to the size limit of job array (i.e., 25 jobs) in local HPC cluster (Figure 4b). New jobs will be submitted once the “workflow script” detects open slots on

the array. With an array size of 25 jobs, one would expect an approximate time acceleration by a factor of 25 if no major time is spent on job queueing. As such, the ARMer library makes it possible to adaptively allocate computing resources to effectively accomplish a high-throughput molecular modeling workflow. This is different from the traditional resource allocation strategy in which one relies on the initially requested/assigned GPU or CPU nodes for the entire computational workflow.



**Figure 4.** The framework and application of adaptive resource manager (i.e., ARMer), a Python library used for adaptive computing resource allocation on high-performance computing cluster. (a) Variables and functions used by ARMer for configuration, submission, and dynamic monitoring of computational tasks. The variables and functions are encapsulated in a Job class. They can be called by a user to prepare a Python script that enables the construction of a high-throughput molecular modeling workflow with effective allocation of computing resources (e.g., CPU and GPU). (b) An exemplary application of ARMer to construct a workflow for high-throughput modeling of fluoroacetate dehalogenase (FACD) mutants. In the workflow, a Python script that runs on a single-CPU thread leverages function and variables from the Job class to manage the modeling sub-tasks (i.e., mutation, molecular dynamics, and quantum mechanics simulations) by configuring, submitting, and monitoring new job scripts. The MD job requests 1 GPU (in orange) and each QM job requests 8 CPUs (in blue). To submit and run individual QM calculations in parallel, a job array with a size of 25 is employed. The type of modeling sub-tasks, time usage, and resource cost are noted on the Figure.

#### 4. Scoring Functions that Describes Sequence-Structure-Function Relationships for Protein Engineering

#### 4.1 A Molecular Dynamics-Derived Descriptor for Enzyme Catalysis

To guide predictive protein engineering, physical descriptors have been identified that correlate with enzyme catalytic efficiency, including enzyme electrostatics in ketosteroid isomerase,<sup>112</sup> Kemp eliminase,<sup>113, 114</sup> methyltransferase,<sup>115</sup> and P450 enzymes;<sup>116</sup> and binding affinity in endoglucanases and cellobiohydrolases<sup>117-119</sup>. Protein dynamics have been proposed as a critical factor to favor substrate positioning<sup>120-127</sup>, control reaction dynamics,<sup>128-133</sup> regulate dynamic network for thermal activation,<sup>134</sup> and tune protein thermal capacity.<sup>135</sup> However, the descriptors that represent the impact of protein dynamics on catalysis remain largely unexplored. Here, we employed a statistical modeling with PCA to identify molecular dynamics-derived descriptors that guide the search of enzyme variants that accommodate non-native substrates with optimal substrate positioning dynamics.

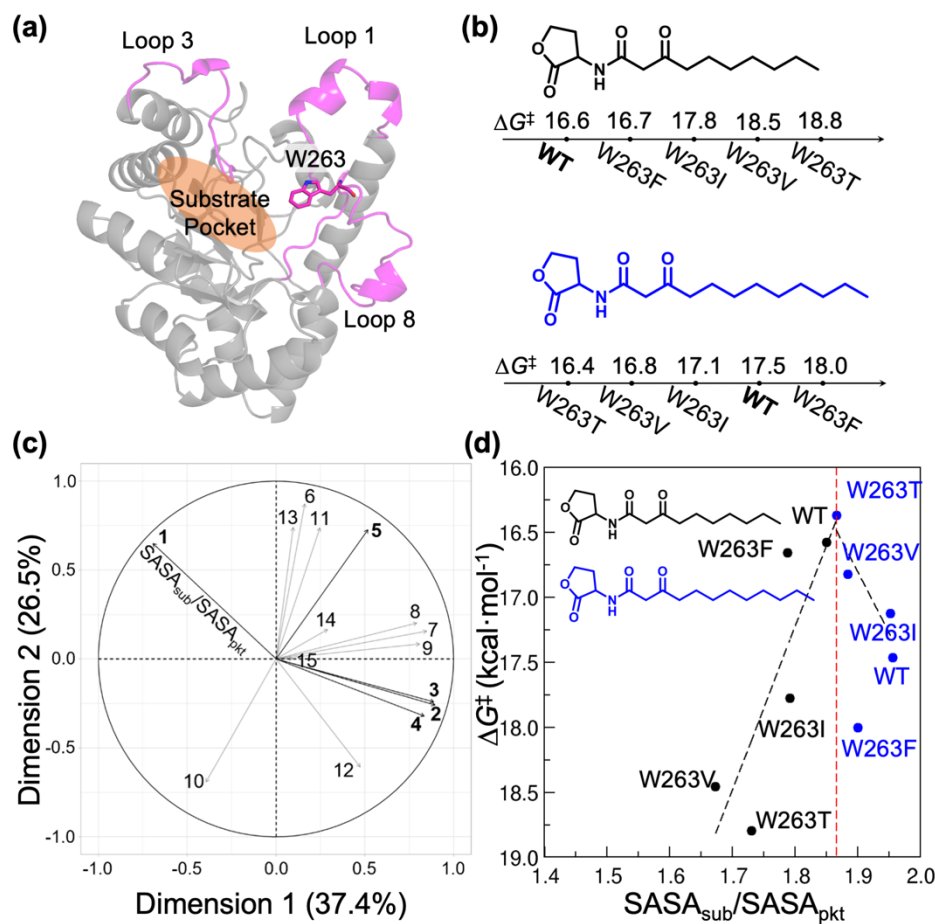
We used lactonase *SsoPox* as a model system (Figure 5a), which catalyzes the hydrolysis of 3-oxo-CX acyl-homoserine lactone (X=10 or 12).<sup>136-141</sup> The WT *SsoPox* is most reactive for the C10 substrate, while the W263T mutant for the C12 substrate (Figure 5b).<sup>137</sup> This enzyme system was chosen primarily because kinetic turnover numbers have been characterized for both C10 and C12 substrates combined with the same set of *SsoPox* variants (i.e., WT, W263F, W263T, W263I, and W263V). This allows us to identify physical descriptors that inform distinct substrate-positioning behaviors of the same enzyme scaffold towards different substrates.

Using molecular dynamics trajectories of each substrate-enzyme variant complex, we calculated fifteen molecular features that are associated with the structural and dynamics characteristics. The descriptors are classified into four groups: 1) solvent accessible surface area; 2) electric field; 3) root-mean-square deviation; and 4) functionally important substrate-residue, residue-residue, and loop-loop distances. We utilized a PCA loading plot to rank the importance

of these descriptors – a higher importance rank indicates that the descriptor contains more information to predict the change of experimental activation free energies (Figure 5c). The PCA analysis detects the SASA ratio of substrate to active-site pocket (i.e.,  $SASA_{\text{sub}}/SASA_{\text{pkt}}$ ) as the top predictor for the mutation effect on activation free energy. Notably, this descriptor was defined to be the substrate-positioning index in our later study.<sup>62</sup>

We further investigated the distribution of  $\Delta G^\ddagger$  values versus  $SASA_{\text{sub}}/SASA_{\text{pkt}}$  for C10 and C12 substrates combined with different enzyme variants. The distribution conforms to a two-segment, piecewise linear correlation plot with a volcano shape (Figure 5d). This quantitative relationship is very similar to the Sabatier principle observed for cellobiohydrolases by Jeppe et al.<sup>117-119</sup> The value of  $SASA_{\text{sub}}/SASA_{\text{pkt}}$  ranges from 1.67 to 1.96, and the activation free energy reaches the minimum ( $\sim 16.5 \text{ kcal}\cdot\text{mol}^{-1}$ ) under an optimal SASA ratio. For the C10 substrate, WT *SsoPox* is most favorable with an SASA ratio of 1.85. In contrast, for C12 substrate, the SASA ratio for WT drifts to 1.96. The  $\Delta G^\ddagger$  reaches the minimum value of  $16.4 \text{ kcal}\cdot\text{mol}^{-1}$  in W263T, where the reaction turnover number for C12 is comparable to the native reaction for C10 in the WT enzyme ( $16.6 \text{ kcal}\cdot\text{mol}^{-1}$ ). The shift of the SASA ratio upon mutation is dominated by the size variation of the active-site pocket. As such, the optimal SASA ratio of substrate to active-site pocket shown in Figure 5a likely reflects the desired enzyme cavity that best accommodates a substrate to achieve efficient catalysis. Replacing the native substrate C10 by C12 leads to an increase of substrate size, which is beyond the accommodation capacity of the WT enzyme but presents a good fit in the W263T variant that has a larger active-site pocket. The results show that the SASA ratio can be employed as a predictive descriptor to guide the search for optimal enzyme mutants for catalyzing non-native substrate. To achieve efficient hydrolysis, a non-native

substrate-bound enzyme variant needs to involve a similar range of SASA ratio to the native substrate-bound WT enzyme.



**Figure 5.** A molecular dynamics-derived descriptor for representing the impact of mutation on enzyme catalysis. (a) The crystal structure for the model enzyme used in the study: lactonase SsoPox (PDB ID: 2VC7). Flexible loops are colored in pink. Substrate binding pocket is indicated by an orange oval. W263 is the spot in which mutations have been performed to investigate the role of mutation on enzyme kinetics. (b) The reaction activation free energies ( $\Delta G^\ddagger$ ) for 3-oxo-CX acyl-homoserine lactone substrates (X = 10 or 12, colored in black and blue, respectively) combined with different enzyme variants (WT, W263F, W263T, W263I, and W263V).  $\Delta G^\ddagger$  value is converted from the turnover rate using Eyring's equation. (c) The PCA loading plot for the descriptors tested in the study. The descriptor is ranked based on its contribution in principal components (from major to minor): 1. SASA<sub>sub</sub>/SASA<sub>pkt</sub>; 2. SASA<sub>pkt</sub>; 3. RMSD<sub>pro</sub>; 4. d<sub>loop1-3</sub>; 5. RMSD<sub>pkt\_sub</sub>; 6. RMSD<sub>sub</sub>; 7. RMSD<sub>pkt</sub>; 8. d<sub>99-229</sub>; 9. d<sub>258-sub</sub>; 10. EF<sub>all</sub>; 11. SASA<sub>sub</sub>; 12. d<sub>97-sub</sub>; 13. EF<sub>noion</sub>; 14. d<sub>223-256</sub>; 15. d<sub>tail-loop8</sub>. The percentage in each axis label indicates the contribution of the principal component to the total variation. (d) Distribution of  $\Delta G^\ddagger$  values versus SASA ratio (i.e., SASA<sub>sub</sub>/SASA<sub>pkt</sub>) in enzyme variants across C10 and C12 substrates. The red dashed lines

indicate the optimal point of  $\Delta G^\ddagger$  and the black dashed lines are the linear fitting of data points on each side of the optimal point.

## 4.2 Deep Learning Models for Protein Function Prediction

In this section, we will introduce two deep learning models that our group recently developed for engineering of enantioselective biocatalysts (i.e., EnzyKR<sup>58</sup>) and antimicrobial peptides (i.e., DeepLasso<sup>142</sup>). EnzyKR was developed to predict the enantiomeric outcome of hydrolase-catalyzed kinetic resolution reactions. DeepLasso was built to predict the antimicrobial activity for ubonodin variants.

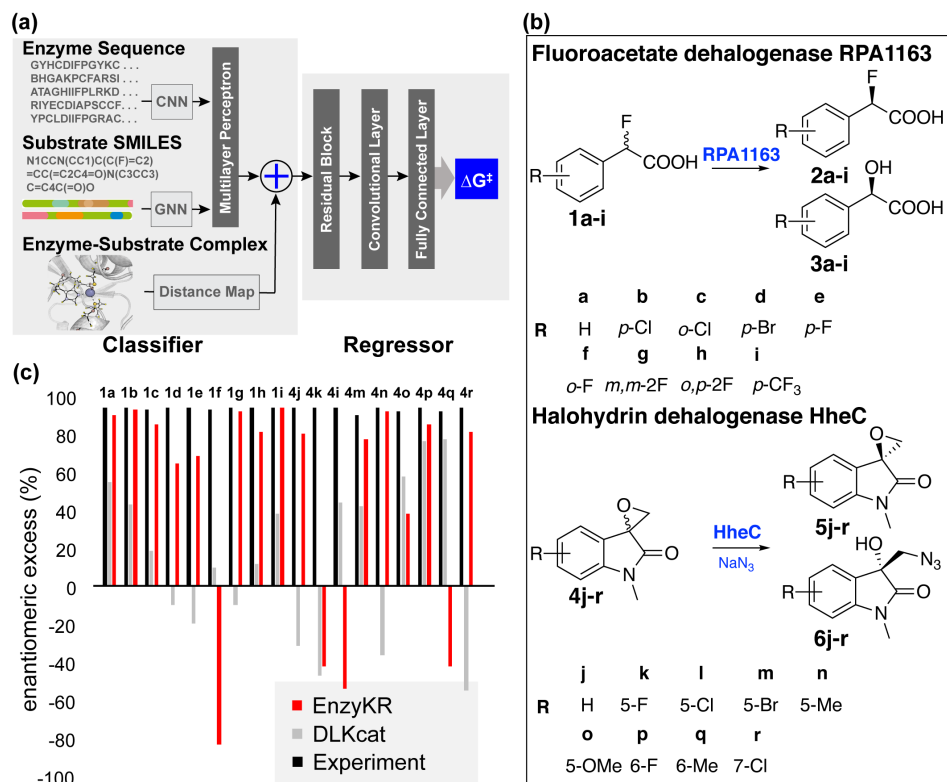
Hydrolases, such as lipases, esterases, and dehalogenases, have been widely employed for kinetic resolution in synthetic reactions in the chemical and pharmaceutical industries.<sup>143-146</sup> Despite the development of empirical,<sup>147</sup> statistical,<sup>148</sup> machine learning,<sup>149</sup> and deep learning models,<sup>150, 151</sup> the “generalist” models that can predict enantioselectivity across a broad spectrum of hydrolase scaffolds, mechanisms, and substrate types remain undeveloped.<sup>152</sup> To address the challenge, our group developed a deep learning model, EnzyKR, to predict the activation free energy of a hydrolase-substrate enantiomer complex. The training and test data include a total of 224 hydrolase-substrate complexes curated from 13 enzyme commission subclasses under the category of hydrolases, which are curated from our integrated enzyme structure-kinetic database IntEnzyDB<sup>53</sup>.

The model consists of a classifier that distinguishes reactive hydrolase-enantiomer complexes from unreactive binding poses, while the regressor predicts the hydrolytic activation free energy (i.e.,  $\Delta G^\ddagger$ ) for the reactive complex. The classifier employs convolutional and graph neural networks to separately encode three types of input: enzyme sequences, substrate SMILES strings, and the distance maps for the hydrolase-substrate complex (Figure 6a). The regressor of EnzyKR takes input from both the classifier embedding and substrate-enzyme interaction maps (a



stacked form of atomic distance map). Notably, the atomic distance map and interaction map differentiate substrate chirality, allowing the model to effectively learn the enantiomeric preference of hydrolases. EnzyKR exhibits a decent prediction accuracy with a Pearson R of 0.91, Spearman R of 0.86, and a mean absolute error (MAE) of 0.8 kcal/mol on the training set (204 data points). EnzyKR also achieves a Pearson R of 0.66, Spearman R of 0.70, and MAE of 1.5 kcal/mol on the test set (20 data points). For both training and test sets, the value of Spearman R resembles that of Pearson R. This indicates that EnzyKR balances the regression of target values or ranking without overfitting.

We further tested EnzyKR for its ability to differentiate enantiomeric reactions using 18 separately-curated hydrolytic reactions catalyzed by FAcD RPA1163<sup>153</sup> and halohydrin HheC<sup>154</sup> (Figure 6b). The performance of EnzyKR was compared against DLKcat, a deep learning  $k_{\text{cat}}$  predictor.<sup>151</sup> Figure 6c shows that compared to the experimental results (black), EnzyKR (red) correctly predicts the favored enantiomer and outperforms DLKcat (grey) in 13 out of 18 reactions (i.e., 1a-e, 1g-i, 4j, 4m-n, 4p, and 4r). In more than half of the test cases, DLKcat predicts an *ee*% value lower than 50%. Due to the lack of chirality encoding in the model, the overall predictive performance of DLKcat appears to be similar to a random guess. Despite a decent performance, we should note that the limitation of EnzyKR lies in the small size of dataset and potentially inadequate representation of chirality. We plan to address these issues in our ongoing works.



**Figure 6.** Design and application of EnzyKR, a deep learning model for predicting the enantiomeric outcome of hydrolase-catalyzed kinetic resolution. (a) EnzyKR consists of a classifier and a regressor. Three types of input data for the classifier involve the complex structure, enzyme sequence, and simplified molecular-input line-entry system (SMILES) string. The distance map derived from the complex structure is encoded using a 2D convolutional neural network (CNN). The multiple sequence alignments (MSA) of the enzyme sequences are also encoded by a 2D CNN model. The substrate SMILES strings are encoded by a graph neural network (GNN) model. The embeddings from the classifier and the interaction maps are used as input for the regressor. The regressor involves one module of cross-attention, followed by residual blocks consisting of three 2D dilated convolution layers, one 2D batch norm layer, and one ReLU layer. Two layers a of fully connected neural network (i.e., multiple-layer perceptron) are employed to conduct regression between the extracted feature and the activation free energy. (b) The test reactions used to assess the ability of EnzyKR to predict the outcomes of kinetic resolution. The test set involves 18 enantioselective hydrolytic reactions catalyzed by two hydrolases. RPA1163 is a fluoroacetate dehalogenase that catalyzes the C–F bond hydrolysis in 9 fluoroacetic acid derivatives labeled using a to i. HheC is a halohydrin dehalogenase that catalyzes the stereoselective epoxide ring-opening in 9 spiro-epoxyoxindoles derivatives labeled using j to r. (c) The predicted enantiomeric excess (*ee*%) values of EnzyKR (red) and the baseline model DLKcat (grey) for 18 enantiomer pairs in hydrolase-catalyzed reactions. The experimental *ee*% value is shown in black.

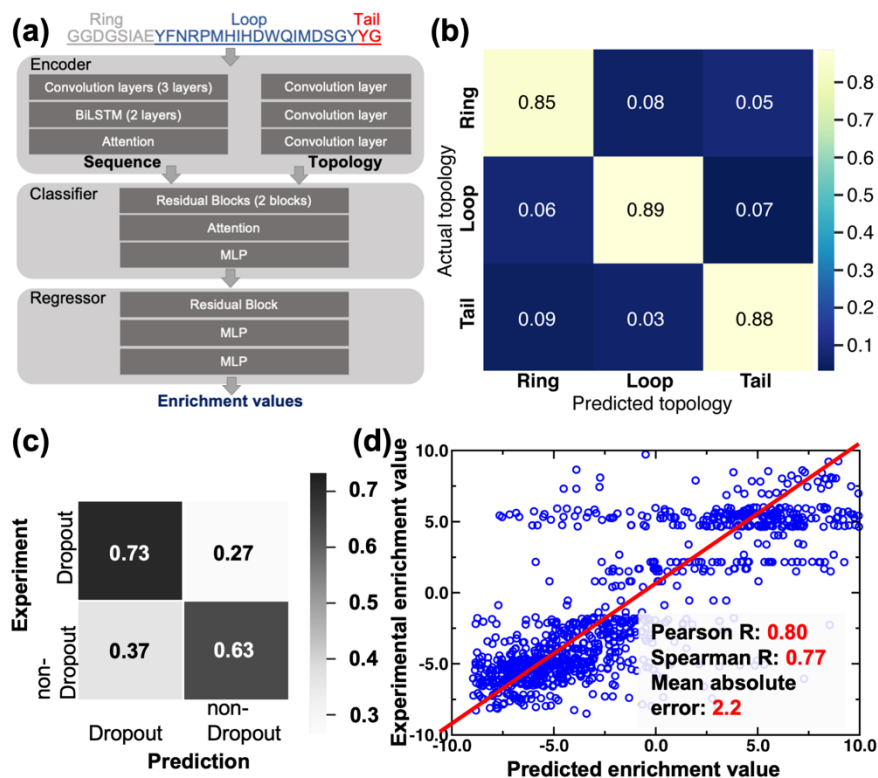
The second deep learning model we intend to introduce is DeepLasso. In recent decades, lasso peptides, such as microcin J25,<sup>155</sup> ubonodin,<sup>156</sup> cloacaenodin,<sup>157</sup> and so on,<sup>18</sup> have emerged

as promising candidates for stemming the tide of antimicrobial crisis. However, the development of computational models to facilitate the engineering of lasso peptide mutants with enhanced antimicrobial activities lag far behind the pace of lasso peptide discovery. To fill in the void, we collaborated with the Link lab and developed DeepLasso to predict the antibiotic activity (i.e., enrichment value) for ubonodin variants (Figure 7). The training and test data involve ~90,000 mutants of lasso peptide ubonodin that were collected from experimental high-throughput screening and next-generation sequencing of single and double mutant library constructed by site-saturation mutagenesis. The antimicrobial activity of a ubonodin mutant is represented by an enrichment value, which is the base-2 logarithm of the ratio of the mutant's frequency at a specific step of the screen relative to the mutant's frequency in the cloning transformation library.<sup>142</sup> Negative enrichment values indicate that the variant likely inhibits RNAP. Dropout mutants are those with a super strong RNAP inhibition activity and their the enrichment values are annotated as “not available”.

Similar to EnzyKR, DeepLasso also adopts a classifier-regressor architecture (Figure 7a). With an input of an ubonodin variant sequence, the classifier first predicts whether the variant likely is a dropout variant. If determined as a non-dropout variant, the regressor is used to predict an enrichment value for the variant. DeepLasso employs a sequence encoder to learn the pattern of the ubonodin amino acid sequence as well as a topology encoder to represent the sequence regions for the ring, loop, and tail of the lasso peptide (Figure 7b). Different from existing deep learning models for prediction of antimicrobial peptides,<sup>158</sup> the topology encoder we implemented in DeepLasso can potentially improve the learning efficiency because the topology of lasso peptides is known to be essential in the inhibition of RNAP.<sup>159-161</sup> To evaluate the accuracy of DeepLasso, we performed confusion matrix analysis for the classifier (Figure 7c) and linear

regression analysis for the regressor (Figure 7d). The results show that DeepLasso achieves a 73% hit rate for the dropout variants and 63% for non-dropout variants (Figure 7c). The higher accuracy for predicting dropout variants is desired because these variants are the most likely to exhibit strong antibiotic activity. For non-dropout variants, the predicted enrichment values are correlated to the experimental value with a Pearson correlation R of 0.80, a Spearman rank correlation R of 0.77, and a MAE of 2.2. The regressor allows us to score the non-dropout variants for their RNAP inhibition activity.

DeepLasso provides a computational tool to map out the fitness landscape of ubonodin variants as potential antibiotics. Though trained with mostly single and double mutants, DeepLasso is capable of identifying higher order ubonodin mutants with enhanced antimicrobial activity. One critical aspect that has yet to be considered here is the ability to predict permeability of ubonodin variants through the membrane of target bacteria. The permeability through cell membrane is independent from RNAP inhibition but should weigh in as an important factor for development of the next version of DeepLasso. Besides, the magnitude to which we can generalize DeepLasso for the antimicrobial prediction of other types of lasso peptides remains a valuable question for investigation.



**Figure 7.** Design architecture of DeepLasso, a deep learning model for predicting the antibiotic activity of lasso peptide ubonodin mutants. (a) The architecture of DeepLasso consists of an encoder, classifier, and regressor. The sequence encoder is constructed by three layers of a convolutional neural network (CNN), two layers of a bidirectional long-short term memory network, and one attention layer. The topology encoder is constructed by three layers of CNN with each layer used to learn a specific topological region of lasso peptide sequence (i.e., ring, loop, or tail). The classifier involves a sequential layout of two residual blocks, one attention layer, and one layer of multilayer perceptron (MLP). The regressor involves a sequential layout of one residual block and two layers of MLP. The tensors derived from the encoder are concatenated and fed into the classifier for prediction; the resulting tensor from the classifier is then used in the regressor for prediction. (b) Confusion matrix analysis for the classifier of DeepLasso. The matrix shows classification of sequence regions of ubonodin variants (ring, loop, and tail). The color scale is used to represent the magnitude of hit rate (i.e., high: yellow; low: dark blue). (c) Confusion matrix analysis for the classifier of DeepLasso. The matrix shows binary classification of dropout versus non-dropout variants. Grayscale is used to represent the magnitude of hit rate (i.e., high: black; low: white). (d) Regression analysis for the non-dropout variants with enrichment values. The linear correlation between experimental vs. predicted enrichment values is shown along with Pearson correlation coefficient, Spearman correlation coefficient, and mean absolute error.

## 5. Applications

The preceding sections present the core technical components underlying Mutexa, including an integrated structure-function database (Section 2), software packages for high-

throughput modeling of protein mutants (Section 3), and scoring functions for predicting the sequence-structure-function relationships (Section 4). In this section, we will demonstrate three applications where these new computational tools are leveraged to investigate the conditions for computational convergence in enzyme modeling,<sup>60</sup> to gain a statistical view across members of methyltransferase family,<sup>61</sup> and to deepen the understanding of dynamic effects in enzyme catalysis.<sup>62</sup>

The first case applies the high-throughput enzyme modeling workflow of Mutexa (i.e., EnzyHTP<sup>56</sup>) to investigate the boundary conditions that should be used in enzyme modeling for a reliable description of mutation effects. In computational protein engineering, functional descriptors have been calculated from molecular simulations to aid the search for beneficial enzyme variants.<sup>162-165</sup> However, the optimal size of the active-site region for computing these descriptors across multiple enzyme variants has not yet been investigated. Using EnzyHTP, we conducted convergence tests on 18 Kemp eliminase variants,<sup>166, 167</sup> evaluating functional descriptors in six active-site regions with varying distances from the substrate. The assessed descriptors include the dynamic fluctuation of the active-site (represented by root-mean-square deviation, or RMSD), the substrate positioning index (represented by the SASA ratio between the substrate and the active site), and the electric field index (represented by the projection of the electric field on the reacting C–H bond). Both molecular mechanics and multiscale quantum mechanics/molecular mechanics methods have been used to compute the descriptors. The descriptor values were determined for each of the eighteen Kemp eliminase variants. Spearman correlation matrices were employed to identify the condition for the region size beyond which further expansion of the boundary does not significantly alter the ranking of descriptor values. Our results show that dynamics-derived descriptors, specifically the dynamic fluctuation and substrate

positioning index, reached convergence at a distance cutoff of 5 Å from the substrate. The electric field descriptor exhibits convergence at 6 Å when employing molecular mechanics methods with truncated enzyme models, and at 4 Å when utilizing quantum mechanics/molecular mechanics methods with the entire enzyme model. This study serves as a reference for selecting descriptors in predictive modeling of enzyme engineering.

The second case combines our integrated structure-function database IntEnzyDB with the workflow software EnzyHTP to study the convergent catalytic behaviors of *S*-adenosyl methionine (SAM)-dependent methyl transferases (MTases). MTases are a ubiquitous class of enzymes catalyzing dozens of reactions in the life processes.<sup>168-171</sup> Despite targeting a large variety of substrates with diverse intrinsic reactivity, MTases demonstrate similar catalytic efficiency.<sup>53, 54, 172</sup> To elucidate the evolutionary adaptation that allows MTases to accommodate the diverse chemical features of their respective substrates, we curated 91 SAM MTases from IntEnzyDB and conducted a comprehensive computational analysis using EnzyHTP to gain insights into how specific properties, such as electric field strength and active site volumes, contribute to achieving similar catalytic efficiency across substrates with different reactivity levels. When looking at *O*-, *N*- and even *C*-targeting MTases, we found that there was not a significant difference in cavity volumes but the electric field strengths have largely adjusted to enhance the target atom's ability to accept a methyl group. For MTases targeting RNA/DNA and histone proteins, the electric field strength accommodates the formal hybridization state. Our study also shows that metal ions in MTases contribute negatively to electric field strength for methyl donation and enzyme scaffolds likely offset these contributions.

The last case integrates the workflow software EnzyHTP with a scoring function of Mutexa, substrate positioning index (SPI, discussed in the Section 4.1), to investigate the behavior of non-

electrostatic dynamics in enzyme catalysis. The dynamic positioning of substrates within the active site, known as substrate positioning dynamics (SPD), plays a crucial role in facilitating enzyme catalysis by aligning the substrate in a reactive conformation.<sup>124, 163, 173-185</sup> However, as conformational changes often coincide with alterations in the electrostatic environment inside the enzyme, it remains unclear whether SPD involves a non-electrostatic component that independently influences catalysis or primarily arises from perturbations in the enzyme's internal electrostatics.<sup>183, 186, 187</sup> To answer this question, we integrated computational and experimental approaches to investigate the non-electrostatic component of SPD using Kemp eliminase as a model enzyme. We employed substrate positioning index to quantify the impact of protein dynamics on substrate positioning. Using EnzyHTP, we selected seven variants for kinetic evaluation, which exhibited significantly different SPD while maintaining similar enzyme interior electrostatics. Our analysis revealed a valley-shaped, two-segment piecewise linear correlation between experimentally determined activation free energies and substrate positioning index values. This trend was further validated using previously reported kinetic data from the Head-Gordon group.<sup>188</sup> Notably, an optimal SPI value, corresponding to the lowest activation free energy, was observed for the R154W variant, a surface mutation located distantly from the active site. Compared to the wild type, the R154W variant displayed favorable SPD, resulting in an increased proportion of reactive conformations for substrate deprotonation. These findings indicate the existence of a non-electrostatic component in SPD, serving as a factor that mediates catalysis by modulating the population of reactive conformations.

## **6. Next Steps**

In this review, we have discussed the construction and applications of Mutexa as a computational ecosystem to facilitate protein engineering. To further progress towards intelligent



protein engineering, we expect to continue developing Mutexa as tools to address real-life problems encountered in protein engineering. The immediate next step is to build a selector of beneficial mutants to enhance catalytic efficiency, mediate selectivity, and expand substrate scope in enzyme engineering. We expect the selector to contain three computational modules that separately evaluate the impact of mutations on 1) enzyme biophysics (i.e., thermal stability, solubility, etc.), 2) enzyme-substrate binding affinity, and 3) enzyme specificity and selectivity. For each of the modules, the proper computational readouts, either derived from data-driven modeling or physics-based simulations, remain a question of investigation. Besides being predictive about the functions, these readouts must be computed with a balanced accuracy and efficiency for better compatibility with a high-throughput computational workflow. How to design a computational protocol within these constraints present the first challenge.

Another challenge faced in the community is how to achieve the modeling or prediction of complex mutants that go beyond single amino acid substitution. Complex mutants with multiple mutation, insertion mutation, or deletion mutation are commonly seen in protein engineering. However, the data-driven and molecular modeling approaches for describing and predicting complex mutants are significantly underdeveloped. With increasing joint efforts in computational and experimental protein engineering, we are hopeful that more and more solutions will be proposed to predict the mutation effects for complex mutations.

## **AUTHOR INFORMATION**

### **Corresponding Author**

\*Email: zhongyue.yang@vanderbilt.edu Phone: 615-343-9849

### **Notes**

The authors declare no competing financial interest.

## ACKNOWLEDGMENTS

This research was supported by the startup grant from Vanderbilt University and the National Institute of General Medical Sciences of the National Institutes of Health under award number R35GM146982. Z. J. Yang thanks the sponsorship from the Dean's Faculty Fellowship in the College of Arts and Science at Vanderbilt. S. L. Stull acknowledges financial support from the Vanderbilt Undergraduate Summer Research Program and the Department of Computer Science. R.J.J. thanks the financial support from the National Institutes of Health Molecular Biophysics Training Grant (MBTP T32 GM008320).

## References

1. Arnold, F. H.; Volkov, A. A., Directed evolution of biocatalysts. *Curr Opin Chem Biol* **1999**, *3* (1), 54-9.
2. Packer, M. S.; Liu, D. R., Methods for the directed evolution of proteins. *Nat Rev Genet* **2015**, *16* (7), 379-94.
3. Arnold, F. H., Directed Evolution: Bringing New Chemistry to Life. *Angew Chem Int Ed Engl* **2018**, *57* (16), 4143-4148.
4. Wang, Y.; Xue, P.; Cao, M.; Yu, T.; Lane, S. T.; Zhao, H., Directed Evolution: Methodologies and Applications. *Chem. Rev.* **2021**, *121* (20), 12384-12444.
5. Kolkman, J. A.; Stemmer, W. P., Directed evolution of proteins by exon shuffling. *Nat Biotechnol* **2001**, *19* (5), 423-8.
6. Akbulut, N.; Tuzlakoglu Ozturk, M.; Pijning, T.; Issever Ozturk, S.; Gumusel, F., Improved activity and thermostability of *Bacillus pumilus* lipase by directed evolution. *J Biotechnol* **2013**, *164* (1), 123-9.
7. Reetz, M. T.; Bocola, M.; Carballeira, J. D.; Zha, D.; Vogel, A., Expanding the range of substrate acceptance of enzymes: combinatorial active-site saturation test. *Angew Chem Int Ed Engl* **2005**, *44* (27), 4192-6.
8. Reetz, M. T.; Carballeira, J. D.; Peyralans, J.; Hobenreich, H.; Maichele, A.; Vogel, A., Expanding the substrate scope of enzymes: combining mutations obtained by CASTing. *Chemistry* **2006**, *12* (23), 6031-8.
9. Yi, D.; Bayer, T.; Badenhorst, C. P. S.; Wu, S.; Doerr, M.; Hohne, M.; Bornscheuer, U. T., Recent trends in biocatalysis. *Chem Soc Rev* **2021**, *50* (14), 8003-8049.
10. Ali, M.; Ishqi, H. M.; Husain, Q., Enzyme engineering: Reshaping the biocatalytic functions. *Biotechnol Bioeng* **2020**, *117* (6), 1877-1894.
11. Min, K.; Kim, H.; Park, H. J.; Lee, S.; Jung, Y. J.; Yoon, J. H.; Lee, J. S.; Park, K.; Yoo, Y. J.; Joo, J. C., Improving the catalytic performance of xylanase from *Bacillus circulans* through structure-based rational design. *Bioresour Technol* **2021**, *340*, 125737.

12. Cecchini, D. A.; Pepe, O.; Pennacchio, A.; Fagnano, M.; Faraco, V., Directed evolution of the bacterial endo-beta-1,4-glucanase from *Streptomyces* sp. G12 towards improved catalysts for lignocellulose conversion. *AMB Express* **2018**, *8* (1), 74.
13. DelRe, C.; Jiang, Y.; Kang, P.; Kwon, J.; Hall, A.; Jayapurna, I.; Ruan, Z.; Ma, L.; Zolkin, K.; Li, T.; Scown, C. D.; Ritchie, R. O.; Russell, T. P.; Xu, T., Near-complete depolymerization of polyesters with nano-dispersed enzymes. *Nature* **2021**, *592* (7855), 558-563.
14. Tournier, V.; Topham, C. M.; Gilles, A.; David, B.; Folgoas, C.; Moya-Leclair, E.; Kamionka, E.; Desrousseaux, M. L.; Texier, H.; Gavalda, S.; Cot, M.; Guemard, E.; Dalibey, M.; Nomme, J.; Cioci, G.; Barbe, S.; Chateau, M.; Andre, I.; Duquesne, S.; Marty, A., An engineered PET depolymerase to break down and recycle plastic bottles. *Nature* **2020**, *580* (7802), 216-219.
15. Li, Z.; Jiang, Y.; Guengerich, F. P.; Ma, L.; Li, S.; Zhang, W., Engineering cytochrome P450 enzyme systems for biomedical and biotechnological applications. *J Biol Chem* **2020**, *295* (3), 833-849.
16. Tang, Q.; Grathwol, C. W.; Aslan-Uzel, A. S.; Wu, S.; Link, A.; Pavlidis, I. V.; Badenhorst, C. P. S.; Bornscheuer, U. T., Directed Evolution of a Halide Methyltransferase Enables Biocatalytic Synthesis of Diverse SAM Analogs. *Angew Chem Int Ed Engl* **2021**, *60* (3), 1524-1527.
17. Lau, J. L.; Dunn, M. K., Therapeutic peptides: Historical perspectives, current development trends, and future directions. *Bioorganic & Medicinal Chemistry* **2018**, *26* (10), 2700-2707.
18. Cheung-Lee, W. L.; Link, A. J., Genome mining for lasso peptides: past, present, and future. *Journal of Industrial Microbiology and Biotechnology* **2019**, *46* (9-10), 1371-1379.
19. Wang, X.; Li, F.; Qiu, W.; Xu, B.; Li, Y.; Lian, X.; Yu, H.; Zhang, Z.; Wang, J.; Li, Z.; Xue, W.; Zhu, F., SYNBP: synthetic binding proteins for research, diagnosis and therapy. *Nucleic Acids Research* **2022**, *50* (D1), D560-D570.
20. Xu, H.; Diolintzi, A.; Storch, J., Fatty acid-binding proteins: functional understanding and diagnostic implications. *Current Opinion in Clinical Nutrition & Metabolic Care* **2019**, *22* (6).
21. Bensing, B. A.; Stubbs, H. E.; Agarwal, R.; Yamakawa, I.; Luong, K.; Solakyildirim, K.; Yu, H.; Hadadianpour, A.; Castro, M. A.; Fialkowski, K. P.; Morrison, K. M.; Wawrzak, Z.; Chen, X.; Lebrilla, C. B.; Baudry, J.; Smith, J. C.; Sullam, P. M.; Iverson, T. M., Origins of glycan selectivity in streptococcal Siglec-like adhesins suggest mechanisms of receptor adaptation. *Nature Communications* **2022**, *13* (1), 2753.
22. Narayanan, H.; Dingfelder, F.; Butté, A.; Lorenzen, N.; Sokolov, M.; Arosio, P., Machine Learning for Biologics: Opportunities for Protein Engineering, Developability, and Formulation. *Trends in Pharmacological Sciences* **2021**, *42* (3), 151-165.
23. Jiang, Y.; Ran, X.; Yang, Z. J., Data-driven enzyme engineering to identify function-enhancing enzymes. *Protein Engineering, Design and Selection* **2023**, *36*, gzac009.
24. Pavelka, A.; Chovancova, E.; Damborsky, J., HotSpot Wizard: a web server for identification of hot spots in protein engineering. *Nucleic Acids Research* **2009**, *37* (suppl\_2), W376-W383.
25. Damborsky, J.; Brezovsky, J., Computational tools for designing and engineering enzymes. *Current Opinion in Chemical Biology* **2014**, *19*, 8-16.
26. Romero, P. A.; Arnold, F. H., Exploring protein fitness landscapes by directed evolution. *Nature Reviews Molecular Cell Biology* **2009**, *10* (12), 866-876.

27. Melnikov, A.; Rogov, P.; Wang, L.; Gnirke, A.; Mikkelsen, T. S., Comprehensive mutational scanning of a kinase in vivo reveals substrate-dependent fitness landscapes. *bioRxiv* **2014**, 004317.
28. Fowler, D. M.; Fields, S., Deep mutational scanning: a new style of protein science. *Nature Methods* **2014**, *11* (8), 801-807.
29. Araya, C. L.; Fowler, D. M., Deep mutational scanning: assessing protein function on a massive scale. *Trends in Biotechnology* **2011**, *29* (9), 435-442.
30. Kries, H.; Blomberg, R.; Hilvert, D., De novo enzymes by computational design. *Current Opinion in Chemical Biology* **2013**, *17* (2), 221-228.
31. Hilvert, D., Design of Protein Catalysts. *Annual Review of Biochemistry* **2013**, *82* (1), 447-470.
32. Bunzel, H. A.; Garrabou, X.; Pott, M.; Hilvert, D., Speeding up enzyme discovery and engineering with ultrahigh-throughput methods. *Current Opinion in Structural Biology* **2018**, *48*, 149-156.
33. Zeymer, C.; Hilvert, D., Directed Evolution of Protein Catalysts. *Annual Review of Biochemistry* **2018**, *87* (1), 131-157.
34. Yeh, A. H.-W.; Norn, C.; Kipnis, Y.; Tischer, D.; Pellock, S. J.; Evans, D.; Ma, P.; Lee, G. R.; Zhang, J. Z.; Anishchenko, I.; Coventry, B.; Cao, L.; Dauparas, J.; Halabiya, S.; DeWitt, M.; Carter, L.; Houk, K. N.; Baker, D., De novo design of luciferases using deep learning. *Nature* **2023**, *614* (7949), 774-780.
35. Pan, X.; Kortemme, T., Recent advances in *de novo* protein design: Principles, methods, and applications. *Journal of Biological Chemistry* **2021**, *296*.
36. Korendovych, I. V.; DeGrado, W. F., De novo protein design, a retrospective. *Quarterly Reviews of Biophysics* **2020**, *53*, e3.
37. Huang, P.-S.; Boyken, S. E.; Baker, D., The coming of age of de novo protein design. *Nature* **2016**, *537* (7620), 320-327.
38. Liu, Q.; Xun, G.; Feng, Y., The state-of-the-art strategies of protein engineering for enzyme stabilization. *Biotechnology Advances* **2019**, *37* (4), 530-537.
39. Rosenfeld, L.; Heyne, M.; Shifman, J. M.; Papo, N., Protein Engineering by Combined Computational and In Vitro Evolution Approaches. *Trends in Biochemical Sciences* **2016**, *41* (5), 421-433.
40. Vaissier Welborn, V.; Head-Gordon, T., Computational Design of Synthetic Enzymes. *Chemical Reviews* **2019**, *119* (11), 6613-6630.
41. Khersonsky, O.; Lipsh, R.; Avizemer, Z.; Ashani, Y.; Goldsmith, M.; Leader, H.; Dym, O.; Rogotner, S.; Trudeau, D. L.; Prilusky, J.; Amengual-Rigo, P.; Guallar, V.; Tawfik, D. S.; Fleishman, S. J., Automated Design of Efficient and Functionally Diverse Enzyme Repertoires. *Molecular Cell* **2018**, *72* (1), 178-186.e5.
42. Risso, V. A.; Romero-Rivera, A.; Gutierrez-Rus, L. I.; Ortega-Muñoz, M.; Santoyo-Gonzalez, F.; Gavira, J. A.; Sanchez-Ruiz, J. M.; Kamerlin, S. C. L., Enhancing a de novo enzyme activity by computationally-focused ultra-low-throughput screening. *Chemical Science* **2020**, *11* (24), 6134-6148.
43. Petřek, M.; Otyepka, M.; Banáš, P.; Košinová, P.; Koča, J.; Damborský, J., CAVER: a new tool to explore routes from protein clefts, pockets and cavities. *BMC Bioinformatics* **2006**, *7* (1), 316.

44. Yang, Z.; Mehmood, R.; Wang, M.; Qi, H. W.; Steeves, A. H.; Kulik, H. J., Revealing quantum mechanical effects in enzyme catalysis with large-scale electronic structure simulation. *Reaction Chemistry & Engineering* **2019**, *4* (2), 298-315.
45. Sheng, X.; Kazemi, M.; Planas, F.; Himo, F., Modeling Enzymatic Enantioselectivity using Quantum Chemical Methodology. *ACS Catalysis* **2020**, *10* (11), 6430-6449.
46. Kiss, G.; Çelebi-Ölçüm, N.; Moretti, R.; Baker, D.; Houk, K. N., Computational Enzyme Design. *Angewandte Chemie International Edition* **2013**, *52* (22), 5700-5725.
47. Sheng, X.; Himo, F., The Quantum Chemical Cluster Approach in Biocatalysis. *Accounts of Chemical Research* **2023**, *56* (8), 938-947.
48. Sequeiros-Borja, C. E.; Surpeta, B.; Brezovsky, J., Recent advances in user-friendly computational tools to engineer protein function. *Briefings in Bioinformatics* **2021**, *22* (3), bbaa150.
49. Greenhalgh, J.; Saraogee, A.; Romero, P. A., Data-driven Protein Engineering. In *Protein Engineering*, 2021; pp 133-151.
50. Yang, K. K.; Wu, Z.; Arnold, F. H., Machine-learning-guided directed evolution for protein engineering. *Nature Methods* **2019**, *16* (8), 687-694.
51. Xu, Y.; Verma, D.; Sheridan, R. P.; Liaw, A.; Ma, J.; Marshall, N. M.; McIntosh, J.; Sherer, E. C.; Svetnik, V.; Johnston, J. M., Deep Dive into Machine Learning Models for Protein Engineering. *Journal of Chemical Information and Modeling* **2020**, *60* (6), 2773-2790.
52. Alley, E. C.; Khimulya, G.; Biswas, S.; AlQuraishi, M.; Church, G. M., Unified rational protein engineering with sequence-based deep representation learning. *Nature Methods* **2019**, *16* (12), 1315-1322.
53. Yan, B.; Ran, X.; Gollu, A.; Cheng, Z.; Zhou, X.; Chen, Y.; Yang, Z. J., IntEnzyDB: an Integrated Structure–Kinetics Enzymology Database. *Journal of Chemical Information and Modeling* **2022**, *62* (22), 5841-5848.
54. Yan, B.; Ran, X.; Jiang, Y.; Torrence, S. K.; Yuan, L.; Shao, Q.; Yang, Z. J., Rate-Perturbing Single Amino Acid Mutation for Hydrolases: A Statistical Profiling. *The Journal of Physical Chemistry B* **2021**, *125* (38), 10682-10691.
55. Juarez, R. J.; Jiang, Y.; Tremblay, M.; Shao, Q.; Link, A. J.; Yang, Z. J., LassoHTP: A High-Throughput Computational Tool for Lasso Peptide Structure Construction and Modeling. *Journal of Chemical Information and Modeling* **2023**, *63* (2), 522-530.
56. Shao, Q.; Jiang, Y.; Yang, Z. J., EnzyHTP: A High-Throughput Computational Platform for Enzyme Modeling. *Journal of Chemical Information and Modeling* **2022**, *62* (3), 647-655.
57. Jiang, Y.; Yan, B.; Chen, Y.; Juarez, R. J.; Yang, Z. J., Molecular Dynamics-Derived Descriptor Informs the Impact of Mutation on the Catalytic Turnover Number in Lactonase Across Substrates. *The Journal of Physical Chemistry B* **2022**, *126* (13), 2486-2495.
58. Ran, X.; Jiang, Y.; Shao, Q.; Yang, Z. J., EnzyKR: A Chirality-Aware Deep Learning Model for Predicting the Outcomes of the Hydrolase-Catalyzed Kinetic Resolution. **2023**.
59. Thokkadam, A.; Do, T.; Ran, X.; Brynildsen, M. P.; Yang, Z. J.; Link, A. J., A High-Throughput Screen Reveals the Structure-Activity Relationship of the Antimicrobial Lasso Peptide Ubonodin. *bioRxiv* **2022**, 2022.12.13.520261.
60. Jiang, Y.; Stull, S. L.; Shao, Q.; Yang, Z. J., Convergence in determining enzyme functional descriptors across Kemp eliminase variants. *Electronic Structure* **2022**, *4* (4), 044007.
61. Jurich, C.; Yang, Z., High-Throughput Computational Investigation of Protein Electrostatics and Cavity for SAM-Dependent Methyltransferases. **2023**.
62. Jiang, Y.; Ding, N.; Shao, Q.; Stull, S.; Cheng, Z.; Yang, Z., Investigating the Non-Electrostatic Component of Substrate Positioning Dynamics. **2023**.

63. Berman, H. M.; Westbrook, J.; Feng, Z.; Gilliland, G.; Bhat, T. N.; Weissig, H.; Shindyalov, I. N.; Bourne, P. E., The Protein Data Bank. *Nucleic Acids Research* **2000**, *28* (1), 235-242.
64. Fleischmann, A.; Darsow, M.; Degtyarenko, K.; Fleischmann, W.; Boyce, S.; Axelsen, K. B.; Bairoch, A.; Schomburg, D.; Tipton, K. F.; Apweiler, R., IntEnz, the integrated relational enzyme database. *Nucleic Acids Research* **2004**, *32* (suppl\_1), D434-D437.
65. Chang, A.; Jeske, L.; Ulbrich, S.; Hofmann, J.; Koblitz, J.; Schomburg, I.; Neumann-Schaal, M.; Jahn, D.; Schomburg, D., BRENDA, the ELIXIR core data resource in 2021: new developments and updates. *Nucleic Acids Research* **2021**, *49* (D1), D498-D508.
66. Wittig, U.; Kania, R.; Golebiewski, M.; Rey, M.; Shi, L.; Jong, L.; Alga, E.; Weidemann, A.; Sauer-Danzwith, H.; Mir, S.; Krebs, O.; Bittkowski, M.; Wetsch, E.; Rojas, I.; Müller, W., SABIO-RK—database for biochemical reaction kinetics. *Nucleic Acids Research* **2012**, *40* (D1), D790-D796.
67. Wang, C. Y.; Chang, P. M.; Ary, M. L.; Allen, B. D.; Chica, R. A.; Mayo, S. L.; Olafson, B. D., ProtBank: A repository for protein design and engineering data. *Protein Science* **2018**, *27* (6), 1113-1124.
68. Huang, X.; Kim, D.; Huang, P.; Vater, A.; Siegel, J. B., Design to Data for mutants of  $\beta$ -glucosidase B from *Paenibacillus polymyxa*: Q22T, W123R, F155G, Y169M, W438D, V401A. *bioRxiv* **2020**, 2020.11.17.387829.
69. Li, F.; Zhang, X.; Renata, H., Enzymatic CH functionalizations for natural product synthesis. *Current Opinion in Chemical Biology* **2019**, *49*, 25-32.
70. Knott, B. C.; Erickson, E.; Allen, M. D.; Gado, J. E.; Graham, R.; Kearns, F. L.; Pardo, I.; Topuzlu, E.; Anderson, J. J.; Austin, H. P.; Dominick, G.; Johnson, C. W.; Rorrer, N. A.; Szostkiewicz, C. J.; Copié, V.; Payne, C. M.; Woodcock, H. L.; Donohoe, B. S.; Beckham, G. T.; McGeehan, J. E., Characterization and engineering of a two-enzyme system for plastics depolymerization. *Proceedings of the National Academy of Sciences* **2020**, *117* (41), 25476-25485.
71. Rorrer, N. A.; Nicholson, S.; Carpenter, A.; Bidy, M. J.; Grundl, N. J.; Beckham, G. T., Combining Reclaimed PET with Bio-based Monomers Enables Plastics Upcycling. *Joule* **2019**, *3* (4), 1006-1027.
72. Wang, J.-B.; Ilie, A.; Yuan, S.; Reetz, M. T., Investigating Substrate Scope and Enantioselectivity of a Defluorinase by a Stereochemical Probe. *Journal of the American Chemical Society* **2017**, *139* (32), 11241-11247.
73. Goldman, P., The Carbon-Fluorine Bond in Compounds of Biological Interest. *Science* **1969**, *164* (3884), 1123-1130.
74. Gan, J.; Siegel, J. B.; German, J. B., Molecular annotation of food – Towards personalized diet and precision health. *Trends in Food Science & Technology* **2019**, *91*, 675-680.
75. Doerr, S.; Harvey, M. J.; Noé, F.; De Fabritiis, G., HTMD: High-Throughput Molecular Dynamics for Molecular Discovery. *Journal of Chemical Theory and Computation* **2016**, *12* (4), 1845-1852.
76. Godwin, R. C.; Melvin, R.; Salsbury, F. R., *Molecular Dynamics Simulations and Computer-Aided Drug Discovery*. Springer New York: 2015; pp 1-30.
77. Parton, D. L.; Grinaway, P. B.; Hanson, S. M.; Beauchamp, K. A.; Chodera, J. D., Ensembler: Enabling High-Throughput Molecular Simulations at the Superfamily Scale. *PLoS Comput Biol* **2016**, *12* (6), e1004728.
78. Amrein, B. A.; Steffen-Munsberg, F.; Szeler, I.; Purg, M.; Kulkarni, Y.; Kamerlin, S. C. L., CADEE: Computer-Aided Directed Evolution of Enzymes. *IUCrJ* **2017**, *4* (1), 50-64.

79. Kim, T. H.; Mehrabi, P.; Ren, Z.; Sljoka, A.; Ing, C.; Bezginov, A.; Ye, L.; Pomès, R.; Prosser, R. S.; Pai, E. F., The role of dimer asymmetry and protomer dynamics in enzyme catalysis. *Science* **2017**, *355* (6322), eaag2355.
80. Makinen, M. W.; Fink, A. L., Reactivity and Cryoenzymology of Enzymes in the Crystalline State. *Annual Review of Biophysics and Bioengineering* **1977**, *6* (1), 301-343.
81. Mehrabi, P.; Di Pietrantonio, C.; Kim, T. H.; Sljoka, A.; Taverner, K.; Ing, C.; Kruglyak, N.; Pomès, R.; Pai, E. F.; Prosser, R. S., Substrate-Based Allosteric Regulation of a Homodimeric Enzyme. *Journal of the American Chemical Society* **2019**, *141* (29), 11540-11556.
82. Mehrabi, P.; Schulz, E. C.; Dsouza, R.; Müller-Werkmeister, H. M.; Tellkamp, F.; Miller, R. J. D.; Pai, E. F., Time-resolved crystallography reveals allosteric communication aligned with molecular breathing. *Science* **2019**, *365* (6458), 1167-1170.
83. Schulz, E. C.; Mehrabi, P.; Müller-Werkmeister, H. M.; Tellkamp, F.; Jha, A.; Stuart, W.; Persch, E.; De Gasparo, R.; Diederich, F.; Pai, E. F.; Miller, R. J. D., The hit-and-return system enables efficient time-resolved serial synchrotron crystallography. *Nature Methods* **2018**, *15* (11), 901-904.
84. Yue, Y.; Fan, J.; Xin, G.; Huang, Q.; Wang, J.-B.; Li, Y.; Zhang, Q.; Wang, W., Comprehensive Understanding of Fluoroacetate Dehalogenase-Catalyzed Degradation of Fluorocarboxylic Acids: A QM/MM Approach. *Environmental Science & Technology* **2021**, *55* (14), 9817-9825.
85. Weber, W.; Fischli, W.; Hochuli, E.; Kupfer, E.; Weibel, E. K., Anantin-a peptide antagonist of the atrial natriuretic factor (ANF). *The Journal of antibiotics* **1991**, *44* (2), 164-171.
86. Salomon, R. A.; Fariás, R. N., Microcin 25, a novel antimicrobial peptide produced by *Escherichia coli*. *Journal of bacteriology* **1992**, *174* (22), 7428-7435.
87. Constantine, K. L.; Friedrichs, M. S.; Detlefsen, D.; Nishio, M.; Tsunakawa, M.; Furumai, T.; Ohkuma, H.; Oki, T.; Hill, S.; Bruccoleri, R. E., High-resolution solution structure of siamycin II: novel amphipathic character of a 21-residue peptide that inhibits HIV fusion. *Journal of biomolecular NMR* **1995**, *5* (3), 271-286.
88. Tsunakawa, M.; Hu, S.-L.; Hoshino, Y.; Detlefsen, D. J.; Hill, S. E.; Furumai, T.; White, R. J.; Nishio, M.; Kawano, K.; Yamamoto, S., Siamycins I and II, new anti-HIV peptides: I. Fermentation, isolation, biological activity and initial characterization. *The Journal of antibiotics* **1995**, *48* (5), 433-434.
89. Pan, S. J.; Cheung, W. L.; Fung, H. K.; Floudas, C. A.; Link, A. J., Computational design of the lasso peptide antibiotic microcin J25. *Protein Engineering, Design and Selection* **2011**, *24* (3), 275-282.
90. Braffman, N. R.; Piscotta, F. J.; Hauver, J.; Campbell, E. A.; Link, A. J.; Darst, S. A., Structural mechanism of transcription inhibition by lasso peptides microcin J25 and capistrain. *Proc. Natl. Acad. Sci.* **2019**, *116* (4), 1273.
91. Solbiati, J. O.; Ciaccio, M.; Fariás, R. N.; González-Pastor, J. E.; Moreno, F.; Salomón, R. A., Sequence analysis of the four plasmid genes required to produce the circular peptide antibiotic microcin J25. *Journal of bacteriology* **1999**, *181* (8), 2659-2662.
92. Cheung-Lee, W. L.; Link, A. J., Genome mining for lasso peptides: past, present, and future. *J. Ind. Microbiol. Biotechnol.* **2019**, 1-9.
93. Do, T.; Link, A. J., Protein Engineering in Ribosomally Synthesized and Post-translationally Modified Peptides (RiPPs). *Biochemistry* **2022**.
94. Jumper, J.; Evans, R.; Pritzel, A.; Green, T.; Figurnov, M.; Ronneberger, O.; Tunyasuvunakool, K.; Bates, R.; Židek, A.; Potapenko, A.; Bridgland, A.; Meyer, C.; Kohl, S.

- A. A.; Ballard, A. J.; Cowie, A.; Romera-Paredes, B.; Nikolov, S.; Jain, R.; Adler, J.; Back, T.; Petersen, S.; Reiman, D.; Clancy, E.; Zielinski, M.; Steinegger, M.; Pacholska, M.; Berghammer, T.; Bodenstein, S.; Silver, D.; Vinyals, O.; Senior, A. W.; Kavukcuoglu, K.; Kohli, P.; Hassabis, D., Highly accurate protein structure prediction with AlphaFold. *Nature* **2021**, *596* (7873), 583-589.
95. Stephen, A. R.; Katelyn, V. C.; Asim, K. B.; Alex, K.; Simon, K.; Joshmyn De La, C.; Victor, A.; Guangfeng, Z.; Frank, D.; Sergey, O.; Gaurav, B., Cyclic peptide structure prediction and design using AlphaFold. *bioRxiv* **2023**, 2023.02.25.529956.
96. D.A. Case, D. S. C., T.E. Cheatham, III, T.A. Darden, R.E. Duke, T.J. Giese, H. Gohlke, A.W. Goetz, D.; Greene, N. H., S. Izadi, A. Kovalenko, T.S. Lee, S. LeGrand, P. Li, C. Lin, J. Liu, T. Luchko, R. Luo; D. Mermelstein, K. M. M., G. Monard, H. Nguyen, I. Omelyan, A. Onufriev, F. Pan, R. Qi, D.R. Roe, A.; Roitberg, C. S., C.L. Simmerling, W.M. Botello-Smith, J. Swails, R.C. Walker, J. Wang, R.M. Wolf, X.; Wu, L. X., D.M. York and P.A. Kollman *AMBER 2017*, University of California, San Francisco: 2017.
97. Hegemann, J. D.; Fage, C. D.; Zhu, S.; Harms, K.; Di Leva, F. S.; Novellino, E.; Marinelli, L.; Marahiel, M. A., The ring residue proline 8 is crucial for the thermal stability of the lasso peptide caulosegnin II. *Molecular BioSystems* **2016**, *12* (4), 1106-1109.
98. Zong, C.; Wu, M. J.; Qin, J. Z.; Link, A. J., Lasso Peptide Benenodin-1 Is a Thermally Actuated [1]Rotaxane Switch. *Journal of the American Chemical Society* **2017**, *139* (30), 10403-10409.
99. Cheung-Lee, W. L.; Parry, M. E.; Jaramillo Cartagena, A.; Darst, S. A.; Link, A. J., Discovery and structure of the antimicrobial lasso peptide citrocin. *Journal of Biological Chemistry* **2019**, *294* (17), 6822-6830.
100. Knappe, T. A.; Manzenrieder, F.; Mas-Moruno, C.; Linne, U.; Sasse, F.; Kessler, H.; Xie, X.; Marahiel, M. A., Introducing lasso peptides as molecular scaffolds for drug design: engineering of an integrin antagonist. *Angewandte Chemie International Edition* **2011**, *50* (37), 8714-8717.
101. Metelev, M.; Tietz, Jonathan I.; Melby, Joel O.; Blair, Patricia M.; Zhu, L.; Livnat, I.; Severinov, K.; Mitchell, Douglas A., Structure, Bioactivity, and Resistance Mechanism of Streptomycin, an Unusual Lasso Peptide from an Understudied Halophilic Actinomycete. *Chemistry & Biology* **2015**, *22* (2), 241-250.
102. Do, T.; Thokkadam, A.; Leach, R.; Link, A. J., Phenotype-Guided Comparative Genomics Identifies the Complete Transport Pathway of the Antimicrobial Lasso Peptide Ubonodin in Burkholderia. *ACS Chemical Biology* **2022**.
103. Cheung-Lee, W. L.; Parry, M. E.; Zong, C.; Cartagena, A. J.; Darst, S. A.; Connell, N. D.; Russo, R.; Link, A. J., Discovery of ubonodin, an antimicrobial lasso peptide active against members of the Burkholderia cepacia complex. *ChemBioChem* **2020**, *21* (9), 1335-1340.
104. Hegemann, J. D.; Zimmermann, M.; Zhu, S.; Steuber, H.; Harms, K.; Xie, X.; Marahiel, M. A., Xanthomonins I–III: A New Class of Lasso Peptides with a Seven-Residue Macrolactam Ring. *Angewandte Chemie International Edition* **2014**, *53* (8), 2230-2234.
105. Yang, Z.; Hajlasz, N.; Kulik, H. J., Computational Modeling of Conformer Stability in Benenodin-1, a Thermally Actuated Lasso Peptide Switch. *The Journal of Physical Chemistry B* **2022**, *126* (18), 3398-3406.
106. Young, T. A.; Silcock, J. J.; Sterling, A. J.; Duarte, F., autodE: Automated Calculation of Reaction Energy Profiles—Application to Organic and Organometallic Reactions. *Angewandte Chemie* **2021**, *133* (8), 4312-4320.



107. An, Q.; Shen, Y.; Fortunelli, A.; Goddard, W. A., QM-Mechanism-Based Hierarchical High-Throughput in Silico Screening Catalyst Design for Ammonia Synthesis. *Journal of the American Chemical Society* **2018**, *140* (50), 17702-17710.
108. Colón, Y. J.; Snurr, R. Q., High-throughput computational screening of metal–organic frameworks. *Chem. Soc. Rev.* **2014**, *43* (16), 5735-5749.
109. Gan, Y.; Miao, N.; Lan, P.; Zhou, J.; Elliott, S. R.; Sun, Z., Robust Design of High-Performance Optoelectronic Chalcogenide Crystals from High-Throughput Computation. *Journal of the American Chemical Society* **2022**, *144* (13), 5878-5886.
110. McInnes, C., Virtual screening strategies in drug discovery. *Current Opinion in Chemical Biology* **2007**, *11* (5), 494-502.
111. Li, Z.; Li, X.; Huang, Y.-Y.; Wu, Y.; Liu, R.; Zhou, L.; Lin, Y.; Wu, D.; Zhang, L.; Liu, H.; Xu, X.; Yu, K.; Zhang, Y.; Cui, J.; Zhan, C.-G.; Wang, X.; Luo, H.-B., Identify potent SARS-CoV-2 main protease inhibitors via accelerated free energy perturbation-based virtual screening of existing drugs. *Proceedings of the National Academy of Sciences* **2020**, *117* (44), 27381-27387.
112. Welborn, V. V.; Head-Gordon, T., Fluctuations of Electric Fields in the Active Site of the Enzyme Ketosteroid Isomerase. *J Am Chem Soc* **2019**, *141* (32), 12487-12492.
113. Bhowmick, A.; Sharma, S. C.; Head-Gordon, T., The Importance of the Scaffold for de Novo Enzymes: A Case Study with Kemp Eliminase. *J Am Chem Soc* **2017**, *139* (16), 5793-5800.
114. Vaissier, V.; Sharma, S. C.; Schaettle, K.; Zhang, T.; Head-Gordon, T., Computational Optimization of Electric Fields for Improving Catalysis of a Designed Kemp Eliminase. *ACS Catalysis* **2018**, *8* (1), 219-227.
115. Yang, Z.; Liu, F.; Steeves, A. H.; Kulik, H. J., Quantum Mechanical Description of Electrostatics Provides a Unified Picture of Catalytic Action Across Methyltransferases. *J Phys Chem Lett* **2019**, *10* (13), 3779-3787.
116. Bím, D.; Alexandrova, A. N., Local Electric Fields As a Natural Switch of Heme-Iron Protein Reactivity. *ACS Catalysis* **2021**, *11* (11), 6534-6546.
117. Kari, J.; Schaller, K.; Molina, G. A.; Borch, K.; Westh, P., The Sabatier principle as a tool for discovery and engineering of industrial enzymes. *Current Opinion in Biotechnology* **2022**, *78*, 102843.
118. Arnlung Bååth, J.; Jensen, K.; Borch, K.; Westh, P.; Kari, J., Sabatier Principle for Rationalizing Enzymatic Hydrolysis of a Synthetic Polyester. *JACS Au* **2022**, *2* (5), 1223-1231.
119. Schaller, K. S.; Molina, G. A.; Kari, J.; Schiano-di-Cola, C.; Sørensen, T. H.; Borch, K.; Peters, G. H. J.; Westh, P., Virtual Bioprospecting of Interfacial Enzymes: Relating Sequence and Kinetics. *ACS Catalysis* **2022**, *12* (12), 7427-7435.
120. Vaissier Welborn, V.; Head-Gordon, T., Computational Design of Synthetic Enzymes. *Chem. Rev.* **2019**, *119* (11), 6613-6630.
121. Mehmood, R.; Vennelakanti, V.; Kulik, H. J., Spectroscopically Guided Simulations Reveal Distinct Strategies for Positioning Substrates to Achieve Selectivity in Nonheme Fe(II)/ $\alpha$ -Ketoglutarate-Dependent Halogenases. *ACS Catalysis* **2021**, *11* (19), 12394-12408.
122. Broom, A.; Rakotoharisoa, R. V.; Thompson, M. C.; Zarifi, N.; Nguyen, E.; Mukhametzhanov, N.; Liu, L.; Fraser, J. S.; Chica, R. A., Ensemble-based enzyme design can recapitulate the effects of laboratory directed evolution in silico. *Nature Communications* **2020**, *11* (1), 4808.
123. Siegel, J. B.; Zanghellini, A.; Lovick, H. M.; Kiss, G.; Lambert, A. R.; St Clair, J. L.; Gallaher, J. L.; Hilvert, D.; Gelb, M. H.; Stoddard, B. L.; Houk, K. N.; Michael, F. E.; Baker,

- D., Computational design of an enzyme catalyst for a stereoselective bimolecular Diels-Alder reaction. *Science* **2010**, *329* (5989), 309-13.
124. Khersonsky, O.; Kiss, G.; Rothlisberger, D.; Dym, O.; Albeck, S.; Houk, K. N.; Baker, D.; Tawfik, D. S., Bridging the gaps in design methodologies by evolutionary optimization of the stability and proficiency of designed Kemp eliminase KE59. *Proc Natl Acad Sci U S A* **2012**, *109* (26), 10358-63.
125. Hur, S.; Bruice, T. C., The near attack conformation approach to the study of the chorismate to prephenate reaction. *Proceedings of the National Academy of Sciences* **2003**, *100* (21), 12015-12020.
126. Kurkcuoglu, Z.; Bakan, A.; Kocaman, D.; Bahar, I.; Doruker, P., Coupling between Catalytic Loop Motions and Enzyme Global Dynamics. *PLOS Computational Biology* **2012**, *8* (9), e1002705.
127. Liao, Q.; Kulkarni, Y.; Sengupta, U.; Petrović, D.; Mulholland, A. J.; van der Kamp, M. W.; Strodel, B.; Kamerlin, S. C. L., Loop Motion in Triosephosphate Isomerase Is Not a Simple Open and Shut Case. *Journal of the American Chemical Society* **2018**, *140* (46), 15889-15903.
128. Masterson, J. E.; Schwartz, S. D., Evolution alters the enzymatic reaction coordinate of dihydrofolate reductase. *J Phys Chem B* **2015**, *119* (3), 989-96.
129. Wang, Z.; Antoniou, D.; Schwartz, S. D.; Schramm, V. L., Hydride Transfer in DHFR by Transition Path Sampling, Kinetic Isotope Effects, and Heavy Enzyme Studies. *Biochemistry* **2016**, *55* (1), 157-66.
130. Liu, C. T.; Layfield, J. P.; Stewart, R. J., 3rd; French, J. B.; Hanoian, P.; Asbury, J. B.; Hammes-Schiffer, S.; Benkovic, S. J., Probing the electrostatics of active site microenvironments along the catalytic cycle for *Escherichia coli* dihydrofolate reductase. *J Am Chem Soc* **2014**, *136* (29), 10349-60.
131. Liu, C. T.; Francis, K.; Layfield, J. P.; Huang, X.; Hammes-Schiffer, S.; Kohen, A.; Benkovic, S. J., *Escherichia coli* dihydrofolate reductase catalyzed proton and hydride transfers: temporal order and the roles of Asp27 and Tyr100. *Proc Natl Acad Sci U S A* **2014**, *111* (51), 18231-6.
132. Venkitakrishnan, R. P.; Zaborowski, E.; McElheny, D.; Benkovic, S. J.; Dyson, H. J.; Wright, P. E., Conformational changes in the active site loops of dihydrofolate reductase during the catalytic cycle. *Biochemistry* **2004**, *43* (51), 16046-55.
133. Bhabha, G.; Ekiert, D. C.; Jennewein, M.; Zmasek, C. M.; Tuttle, L. M.; Kroon, G.; Dyson, H. J.; Godzik, A.; Wilson, I. A.; Wright, P. E., Divergent evolution of protein conformational dynamics in dihydrofolate reductase. *Nat Struct Mol Biol* **2013**, *20* (11), 1243-9.
134. Gao, S.; Thompson, E. J.; Barrow, S. L.; Zhang, W.; Iavarone, A. T.; Klinman, J. P., Hydrogen-Deuterium Exchange within Adenosine Deaminase, a TIM Barrel Hydrolase, Identifies Networks for Thermal Activation of Catalysis. *J Am Chem Soc* **2020**, *142* (47), 19936-19949.
135. Bunzel, H. A.; Kries, H.; Marchetti, L.; Zeymer, C.; Mittl, P. R. E.; Mulholland, A. J.; Hilvert, D., Emergence of a Negative Activation Heat Capacity during Evolution of a Designed Enzyme. *J Am Chem Soc* **2019**, *141* (30), 11745-11748.
136. Afriat, L.; Roodveldt, C.; Manco, G.; Tawfik, D. S., The latent promiscuity of newly identified microbial lactonases is linked to a recently diverged phosphotriesterase. *Biochemistry* **2006**, *45* (46), 13677-86.
137. Hiblot, J.; Gotthard, G.; Elias, M.; Chabriere, E., Differential active site loop conformations mediate promiscuous activities in the lactonase *SsoPox*. *PLoS One* **2013**, *8* (9), e75272.

138. Ng, F. S.; Wright, D. M.; Seah, S. Y., Characterization of a phosphotriesterase-like lactonase from *Sulfolobus solfataricus* and its immobilization for disruption of quorum sensing. *Applied and Environmental Microbiology* **2011**, *77* (4), 1181-6.
139. Bzdrenga, J.; Daude, D.; Remy, B.; Jacquet, P.; Plener, L.; Elias, M.; Chabriere, E., Biotechnological applications of quorum quenching enzymes. *Chem Biol Interact* **2017**, *267*, 104-115.
140. Billot, R.; Plener, L.; Jacquet, P.; Elias, M.; Chabriere, E.; Daude, D., Engineering acyl-homoserine lactone-interfering enzymes toward bacterial control. *Journal of Biological Chemistry* **2020**, *295* (37), 12993-13007.
141. Sikdar, R.; Elias, M., Quorum quenching enzymes and their effects on virulence, biofilm, and microbiomes: a review of recent advances. *Expert Rev Anti-Infe* **2020**, *18* (12), 1221-1233.
142. Thokkadam, A.; Do, T.; Ran, X.; Brynildsen, M. P.; Yang, Z. J.; Link, A. J., High-Throughput Screen Reveals the Structure–Activity Relationship of the Antimicrobial Lasso Peptide Ubonodin. *ACS Central Science* **2023**, *9* (3), 540-550.
143. Pinheiro, M. P.; Rios, N. S.; Fonseca, T. d. S.; Bezerra, F. d. A.; Rodríguez-Castellón, E.; Fernandez-Lafuente, R.; Carlos de Mattos, M.; Dos Santos, J. C.; Gonçalves, L. R., Kinetic resolution of drug intermediates catalyzed by lipase B from *Candida antarctica* immobilized on immovead-350. *Biotechnology Progress* **2018**, *34* (4), 878-889.
144. Bornscheuer, U. T.; Kazlauskas, R. J., *Hydrolases in organic synthesis: regio-and stereoselective biotransformations*. John Wiley & Sons: 2006.
145. Lee, J.; Oh, Y.; Choi, Y. K.; Choi, E.; Kim, K.; Park, J.; Kim, M.-J., Dynamic kinetic resolution of diarylmethanols with an activated lipoprotein lipase. *ACS Catalysis* **2015**, *5* (2), 683-689.
146. Bassegoda, A.; Nguyen, G. S.; Schmidt, M.; Kourist, R.; Diaz, P.; Bornscheuer, U. T., Rational protein design of *Paenibacillus barcinonensis* esterase EstA for kinetic resolution of tertiary alcohols. *ChemCatChem* **2010**, *2* (8), 962-967.
147. Kazlauskas, R. J.; Weissfloch, A. N.; Rappaport, A. T.; Cuccia, L. A., A rule to predict which enantiomer of a secondary alcohol reacts faster in reactions catalyzed by cholesterol esterase, lipase from *Pseudomonas cepacia*, and lipase from *Candida rugosa*. *The Journal of Organic Chemistry* **1991**, *56* (8), 2656-2665.
148. Tomić, S.; Kojić-Prodić, B., A quantitative model for predicting enzyme enantioselectivity: application to *Burkholderia cepacia* lipase and 3-(aryloxy)-1, 2-propanediol derivatives. *Journal of Molecular Graphics and Modelling* **2002**, *21* (3), 241-252.
149. Cadet, F.; Fontaine, N.; Li, G.; Sanchis, J.; Ng Fuk Chong, M.; Pandjaitan, R.; Vetrivel, I.; Offmann, B.; Reetz, M. T., A machine learning approach for reliable prediction of amino acid interactions and its application in the directed evolution of enantioselective enzymes. *Scientific reports* **2018**, *8* (1), 16757.
150. Heckmann, D.; Lloyd, C. J.; Mih, N.; Ha, Y.; Zielinski, D. C.; Haiman, Z. B.; Desouki, A. A.; Lercher, M. J.; Palsson, B. O., Machine learning applied to enzyme turnover numbers reveals protein structural correlates and improves metabolic models. *Nature communications* **2018**, *9* (1), 5252.
151. Li, F.; Yuan, L.; Lu, H.; Li, G.; Chen, Y.; Engqvist, M. K.; Kerkhoven, E. J.; Nielsen, J., Deep learning-based  $k_{cat}$  prediction enables improved enzyme-constrained model reconstruction. *Nature Catalysis* **2022**, *5* (8), 662-672.
152. Jiang, Y.; Ran, X.; Yang, Z. J., Data-driven enzyme engineering to identify function-enhancing enzymes. *Protein Engineering, Design and Selection* **2022**, gzac009.

153. Zhang, H.; Tian, S.; Yue, Y.; Li, M.; Tong, W.; Xu, G.; Chen, B.; Ma, M.; Li, Y.; Wang, J.-b., Semirational design of fluoroacetate dehalogenase RPA1163 for kinetic resolution of  $\alpha$ -fluorocarboxylic acids on a gram scale. *ACS Catalysis* **2020**, *10* (5), 3143-3151.
154. Zhang, F.-R.; Wan, N.-W.; Ma, J.-M.; Cui, B.-D.; Han, W.-Y.; Chen, Y.-Z., Enzymatic Kinetic Resolution of Bulky Spiro-Epoxyoxindoles via Halohydrin Dehalogenase-Catalyzed Enantio- and Regioselective Azidolysis. *ACS Catalysis* **2021**, *11* (15), 9066-9072.
155. Braffman, N. R.; Piscotta, F. J.; Hauver, J.; Campbell, E. A.; Link, A. J.; Darst, S. A., Structural mechanism of transcription inhibition by lasso peptides microcin J25 and capistrain. *Proceedings of the National Academy of Sciences* **2019**, *116* (4), 1273-1278.
156. Cheung-Lee, W. L.; Parry, M. E.; Zong, C.; Cartagena, A. J.; Darst, S. A.; Connell, N. D.; Russo, R.; Link, A. J., Discovery of Ubonodin, an Antimicrobial Lasso Peptide Active against Members of the Burkholderia cepacia Complex. *ChemBioChem* **2020**, *21* (9), 1335-1340.
157. Carson, D. V.; Patiño, M.; Elashal, H. E.; Cartagena, A. J.; Zhang, Y.; Whitley, M. E.; So, L.; Kayser-Browne, A. K.; Earl, A. M.; Bhattacharyya, R. P.; Link, A. J., Cloacaenodin, an Antimicrobial Lasso Peptide with Activity against Enterobacter. *ACS Infectious Diseases* **2023**, *9* (1), 111-121.
158. Dean, S. N.; Alvarez, J. A. E.; Zabetakis, D.; Walper, S. A.; Malanoski, A. P., PepVAE: variational autoencoder framework for antimicrobial peptide generation and activity prediction. *Frontiers in microbiology* **2021**, 2764.
159. Mathavan, I.; Zirah, S.; Mehmood, S.; Choudhury, H. G.; Goulard, C.; Li, Y.; Robinson, C. V.; Rebuffat, S.; Beis, K., Structural basis for hijacking siderophore receptors by antimicrobial lasso peptides. *Nature chemical biology* **2014**, *10* (5), 340-342.
160. Braffman, N.; Piscotta, F. J.; Hauver, J.; Campbell, E. A.; Link, A. J.; Darst, S. A., Structural mechanism of transcription inhibition by lasso peptides microcin J25 and capistrain. *Proceedings of the National Academy of Sciences of the United States of America* **2019**, *116*, 1273-1278.
161. Wilson, K. A.; Kalkum, M.; Ottesen, J.; Yuzenkova, J.; Chait, B. T.; Landick, R.; Muir, T.; Severinov, K.; Darst, S. A., Structure of microcin J25, a peptide inhibitor of bacterial RNA polymerase, is a lassoed tail. *J Am Chem Soc* **2003**, *125* (41), 12475-83.
162. Fried, S. D.; Boxer, S. G., Electric Fields and Enzyme Catalysis. *Annu Rev Biochem* **2017**, *86*, 387-415.
163. Jiang, Y.; Yan, B.; Chen, Y.; Juarez, R. J.; Yang, Z. J., Molecular Dynamics-Derived Descriptor Informs the Impact of Mutation on the Catalytic Turnover Number in Lactonase Across Substrates. *J. Phys. Chem. B* **2022**, *126* (13), 2486-2495.
164. Bartlett, G. J.; Porter, C. T.; Borkakoti, N.; Thornton, J. M., Analysis of catalytic residues in enzyme active sites. *J Mol Biol* **2002**, *324* (1), 105-21.
165. Lodola, A.; Sirirak, J.; Fey, N.; Rivara, S.; Mor, M.; Mulholland, A. J., Structural Fluctuations in Enzyme-Catalyzed Reactions: Determinants of Reactivity in Fatty Acid Amide Hydrolase from Multivariate Statistical Analysis of Quantum Mechanics/Molecular Mechanics Paths. *J Chem Theory Comput* **2010**, *6* (9), 2948-60.
166. Rothlisberger, D.; Khersonsky, O.; Wollacott, A. M.; Jiang, L.; DeChancie, J.; Betker, J.; Gallaher, J. L.; Althoff, E. A.; Zanghellini, A.; Dym, O.; Albeck, S.; Houk, K. N.; Tawfik, D. S.; Baker, D., Kemp elimination catalysts by computational enzyme design. *Nature* **2008**, *453* (7192), 190-5.
167. Hong, N. S.; Petrovic, D.; Lee, R.; Gryn'ova, G.; Purg, M.; Saunders, J.; Bauer, P.; Carr, P. D.; Lin, C. Y.; Mabbitt, P. D.; Zhang, W.; Altamore, T.; Easton, C.; Coote, M. L.;

- Kamerlin, S. C. L.; Jackson, C. J., The evolution of multiple active site configurations in a designed enzyme. *Nat Commun* **2018**, *9* (1), 3900.
168. Bügl, H.; Fauman, E. B.; Staker, B. L.; Zheng, F.; Kushner, S. R.; Saper, M. A.; Bardwell, J. C.; Jakob, U., RNA methylation under heat shock control. *Mol Cell* **2000**, *6* (2), 349-60.
169. Nai, Y.-S.; Huang, Y.-C.; Yen, M.-R.; Chen, P.-Y., Diversity of Fungal DNA Methyltransferases and Their Association With DNA Methylation Patterns. *Frontiers in Microbiology* **2021**, *11*.
170. Dhe-Paganon, S.; Syeda, F.; Park, L., DNA methyl transferase 1: regulatory mechanisms and implications in health and disease. *Int J Biochem Mol Biol* **2011**, *2* (1), 58-66.
171. Zhang, H.; Lang, Z.; Zhu, J.-K., Dynamics and function of DNA methylation in plants. *Nature Reviews Molecular Cell Biology* **2018**, *19* (8), 489-506.
172. Sousa, S. F.; Calixto, A. R.; Ferreira, P.; Ramos, M. J.; Lim, C.; Fernandes, P. A., Activation Free Energy, Substrate Binding Free Energy, and Enzyme Efficiency Fall in a Very Narrow Range of Values for Most Enzymes. *ACS Catalysis* **2020**, *10* (15), 8444-8453.
173. Norberg, A. L.; Dybvik, A. I.; Zakariassen, H.; Mormann, M.; Peter-Katalinić, J.; Eijsink, V. G. H.; Sørli, M., Substrate positioning in chitinase A, a processive chito-biohydrolase from *Serratia marcescens*. *FEBS Letters* **2011**, *585* (14), 2339-2344.
174. Hamre, A. G.; Jana, S.; Reppert, N. K.; Payne, C. M.; Sorlie, M., Processivity, Substrate Positioning, and Binding: The Role of Polar Residues in a Family 18 Glycoside Hydrolase. *Biochemistry* **2015**, *54* (49), 7292-7306.
175. Patra, N.; Ioannidis, E. I.; Kulik, H. J., Computational Investigation of the Interplay of Substrate Positioning and Reactivity in Catechol O-Methyltransferase. *Plos One* **2016**, *11* (8).
176. Hu, S. S.; Offenbacher, A. R.; Thompson, E. M.; Gee, C. L.; Wilcoxon, J.; Carr, C. A. M.; Prigozhin, D. M.; Yang, V.; Alber, T.; Britt, R. D.; Fraser, J. S.; Klinman, J. P., Biophysical Characterization of a Disabled Double Mutant of Soybean Lipoxygenase: The "Undoing" of Precise Substrate Positioning Relative to Metal Cofactor and an Identified Dynamical Network. *Journal of the American Chemical Society* **2019**, *141* (4), 1555-1567.
177. Mehmood, R.; Qi, H. W.; Steeves, A. H.; Kulik, H. J., The Protein's Role in Substrate Positioning and Reactivity for Biosynthetic Enzyme Complexes: The Case of SyrB2/SyrB1. *Acs Catalysis* **2019**, *9* (6), 4930-4943.
178. Yabukarski, F.; Biel, J. T.; Pinney, M. M.; Doukov, T.; Powers, A. S.; Fraser, J. S.; Herschlag, D., Assessment of enzyme active site positioning and tests of catalytic mechanisms through X-ray-derived conformational ensembles. *Proc Natl Acad Sci U S A* **2020**, *117* (52), 33204-33215.
179. Mehmood, R.; Vennelakanti, V.; Kulik, H. J., Spectroscopically Guided Simulations Reveal Distinct Strategies for Positioning Substrates to Achieve Selectivity in Nonheme Fe(II)/alpha-Ketoglutarate-Dependent Halogenases. *Acs Catalysis* **2021**, *11* (19), 12394-12408.
180. Ruscio, J. Z.; Kohn, J. E.; Ball, K. A.; Head-Gordon, T., The Influence of Protein Dynamics on the Success of Computational Enzyme Design. *Journal of the American Chemical Society* **2009**, *131* (39), 14111-14115.
181. Siegel, J. B.; Zanghellini, A.; Lovick, H. M.; Kiss, G.; Lambert, A. R.; Clair, J. L. S.; Gallaher, J. L.; Hilvert, D.; Gelb, M. H.; Stoddard, B. L.; Houk, K. N.; Michael, F. E.; Baker, D., Computational Design of an Enzyme Catalyst for a Stereoselective Bimolecular Diels-Alder Reaction. *Science* **2010**, *329* (5989), 309-313.

182. Blomberg, R.; Kries, H.; Pinkas, D. M.; Mittl, P. R. E.; Grutter, M. G.; Privett, H. K.; Mayo, S. L.; Hilvert, D., Precision is essential for efficient catalysis in an evolved Kemp eliminase. *Nature* **2013**, *503* (7476), 418-+.
183. Broom, A.; Rakotoharisoa, R. V.; Thompson, M. C.; Zarifi, N.; Nguyen, E.; Mukhametzhonov, N.; Liu, L.; Fraser, J. S.; Chica, R. A., Ensemble-based enzyme design can recapitulate the effects of laboratory directed evolution in silico. *Nat. Commun.* **2020**, *11* (1).
184. Haataja, T.; Gado, J. E.; Nutt, A.; Anderson, N. T.; Nilsson, M.; Momeni, M. H.; Isaksson, R.; Valjamae, P.; Johansson, G.; Payne, C. M.; Stahlberg, J., Enzyme kinetics by GH7 cellobiohydrolases on chromogenic substrates is dictated by non-productive binding: insights from crystal structures and MD simulation. *Febs J* **2023**, *290* (2), 379-399.
185. Offenbacher, A. R.; Sharma, A.; Doan, P. E.; Klinman, J. P.; Hoffman, B. M., The Soybean Lipxygenase-Substrate Complex: Correlation between the Properties of Tunneling-Ready States and ENDOR-Detected Structures of Ground States. *Biochemistry* **2020**, *59* (7), 901-910.
186. Wu, Y. F.; Fried, S. D.; Boxer, S. G., A Preorganized Electric Field Leads to Minimal Geometrical Reorientation in the Catalytic Reaction of Ketosteroid Isomerase. *Journal of the American Chemical Society* **2020**, *142* (22), 9993-9998.
187. Otten, R.; Padua, R. A. P.; Bunzel, H. A.; Nguyen, V.; Pitsawong, W.; Patterson, M.; Sui, S.; Perry, S. L.; Cohen, A. E.; Hilvert, D.; Kern, D., How directed evolution reshapes the energy landscape in an enzyme to boost catalysis. *Science* **2020**, *370* (6523), 1442-1446.
188. Bhowmick, A.; Sharma, S. C.; Honma, H.; Head-Gordon, T., The role of side chain entropy and mutual information for improving the de Novo design of Kemp eliminases KE07 and KE70. *Phys Chem Chem Phys* **2016**, *18* (28), 19386-96.

## Entry for the Table of Contents

