



Research review paper

Unlocking the potential of enzyme engineering *via* rational computational design strategies

Lei Zhou¹, Chunmeng Tao¹, Xiaolin Shen, Xinxiao Sun, Jia Wang^{*}, Qipeng Yuan^{*}

State Key Laboratory of Chemical Resource Engineering, Beijing University of Chemical Technology, Beijing 100029, China

ARTICLE INFO

Keywords:

Enzyme engineering
Rational computational design
Structure-based computational design
Sequence-based computational design
Data-driven machine learning

ABSTRACT

Enzymes play a pivotal role in various industries by enabling efficient, eco-friendly, and sustainable chemical processes. However, the low turnover rates and poor substrate selectivity of enzymes limit their large-scale applications. Rational computational enzyme design, facilitated by computational algorithms, offers a more targeted and less labor-intensive approach. There has been notable advancement in employing rational computational protein engineering strategies to overcome these issues, it has not been comprehensively reviewed so far. This article reviews recent developments in rational computational enzyme design, categorizing them into three types: structure-based, sequence-based, and data-driven machine learning computational design. Case studies are presented to demonstrate successful enhancements in catalytic activity, stability, and substrate selectivity. Lastly, the article provides a thorough analysis of these approaches, highlights existing challenges and potential solutions, and offers insights into future development directions.

1. Introduction

Remarkably efficient metabolic enzymes play a pivotal role in converting readily available, simple starting materials into valuable products within microbial cells. Furthermore, these enzymes facilitate a diverse array of chemical reactions, typically operating under mild reaction conditions and exhibiting high selectivity (Hollmann and Fernandez-Lafuente, 2021). These advantages have led to the incorporation of enzymes in various industries, including food (Flynn et al., 2021; Lin et al., 2022; Punia, 2020; Rentschler et al., 2015; van Donkelaar et al., 2016; Xu et al., 2020), agricultural (Costa-Silva et al., 2021; Sijnamanoj et al., 2021; Tingley et al., 2021), cosmetics (Almeida et al., 2021; Fournière et al., 2021) and pharmaceuticals (Meghwanshi et al., 2020; Park et al., 2017; Rosenthal and Lütz, 2018). Regrettably, enzymes come with certain drawbacks such as low turnover rates, limited stability and a narrow substrate scope. Consequently, there is a necessity for enzyme engineering to tailor them for diverse industrial applications (Li et al., 2018b).

Utilizing both directed evolution and semi-rational approaches, enzyme engineering involves iterative steps and mutation libraries, with the process repeated until the desired variant is obtained (Dinmukhamed et al., 2021; Gargiulo and Soumillion, 2021). However, the

unknown structure-function relationship introduces uncertainty in discovering desired protein variants through a few rounds of iterative mutagenesis. It's more likely that the researchers may never come across it even with a comprehensive library, leading them to rely on insights from rational design and evolutionary analysis (Blazeck et al., 2022). In contrast, rational enzyme design stands out for its ability to provide higher predictive accuracy and streamline the screening library (Jumper et al., 2021; Kuhlman and Bradley, 2019). Proposed mutations are introduced after evaluation and design, increasing the likelihood of beneficial mutations while saving time and labor. When high-throughput screening is not feasible, this approach stands out as particularly useful (Cui et al., 2022; Steiner and Schwab, 2012).

Based on the design principles, rational computational enzyme design is divided into three classes, structure-based computational design, sequence-based computational design, and data-driven machine learning computational design. Here, we assess the recent advancements in rational computational enzyme design, emphasizing the effectiveness and versatility of different approaches in enhancing catalytic activity, stability and substrate selectivity (Table 1). Successful cases are cited to illustrate these improvements. Lastly, the article provides a thorough analysis of these approaches, highlights existing challenges and potential solutions, and offers insights into future development directions.

* Corresponding authors.

E-mail addresses: wangjia@buct.edu.cn (J. Wang), yuanqp@buct.edu.cn (Q. Yuan).

¹ Lei Zhou and Chunmeng Tao contributed equally.

2. Structure-based computational enzyme engineering and design

Structure-based design necessitates the identification of critical residues within the substrate binding pocket, understanding the enzyme-catalyzed reaction's chemical mechanism and the key amino acid residues involved. This is achieved through chemical intuition and computational methods to precisely tailor the interplay between key amino acid residues and the substrate in the active pocket, potentially enabling *de novo* computational enzyme design. Although traditional

methods including cryo-electron microscopy or X-ray diffraction have significantly advanced in obtaining enzyme crystal structures, with 170,000 structures available out of 190 million proteins with known amino acid sequences, the process remains time-consuming and costly (Consortium, 2020) (Burley et al., 2020). To overcome these limitations, AlphaFold2 and RoseTTAfold employing machine learning algorithms to accurately predict protein structures using extensive protein information data (Baek et al., 2021; Jumper et al., 2021). AlphaFold2, in particular, has demonstrated high accuracy comparable to experimental results, exhibiting a root-mean-squared deviation (RMSD) of about 1.5 Å

Table 1

Examples of strategies for enzyme engineering via rational computational design.

Type	Enzyme/Reaction	Effect/Strategy	Targeted Property	Reference
Structure-based	Phosphoserine aminotransferase	Calculation of relative binding free energy	Substrate specificity	(Zhang et al., 2019)
Structure-based	Ferulic acid decarboxylase	Calculation of affinity between enzymes and substrates	Activity	(Mori et al., 2021)
Structure-based	Fatty acid photodecarboxylase	Modulating the steric hindrance of the binding pocket	Activity and selectivity	(Xu et al., 2022)
Structure-based	Fatty acid photodecarboxylase	Enhancing the electronic interaction between enzymes and substrates	Selectivity	(Li et al., 2021a)
Structure-based	4-hydroxyphenylpyruvate dioxygenase	Manipulating the hydrogen bond network	Product profile	(Lin et al., 2021)
Structure-based	Santalene synthase	Reducing the steric hindrance to increase pocket space	Product profile	(Zha et al., 2022)
Structure-based	The short-chain dehydrogenase/reductase	Tuning the size of the active pockets	Selectivity	(Su et al., 2020)
Structure-based	Amidase	Modifying electrostatic interactions	Activity and stability	(Galmés et al., 2022)
Structure-based	Formolase	Redesign of the binding pocket	Activity and specificity	(Siegel et al., 2015)
Structure-based	Limonene epoxide hydrolase	Redesign of the binding pocket	Selectivity	(Wijma et al., 2015)
Structure-based	C-N lyases	Redesign of the binding pocket	Activity and selectivity	(Cui et al., 2021)
Structure-based	Lipase	Design disulfide bonds	Stability	(Li et al., 2018a)
Structure-based	Glucuronidase	Truncating C-terminal region	Stability	(Han et al., 2018)
Structure-based	ω -transaminase	Calculating folding energy to identify stable point mutations	Stability	(Meng et al., 2020)
Structure-based	Diels-Alder	<i>De novo</i> design of enzyme by "inside-out"	New enzyme	(Siegel et al., 2010)
Structure-based	Kemp elimination	<i>De novo</i> design of enzyme by "inside-out"	New enzyme	(Röthlisberger et al., 2008)
Structure-based	Retro-Aldol	<i>De novo</i> design of enzyme by "inside-out"	New enzyme	(Jiang et al., 2008)
Structure-based	Ester hydrolysis	<i>De novo</i> design of enzyme by "inside-out"	New enzyme	(Richter et al., 2012)
Structure-based	Morita-Baylis-Hillman	<i>De novo</i> design of enzyme by "inside-out"	New enzyme	(Bjelic et al., 2013)
Sequence-based	Glycosyltransferases	Consensus design	Substrate specificity	(Teze et al., 2021)
Sequence-based	Triosephosphate isomerase	Consensus design	Stability	(Sullivan et al., 2012)
Sequence-based	Glucose-1-dehydrogenase	Consensus design	Stability	(Vázquez-Figueroa et al., 2007)
Sequence-based	Haloalkane dehalogenase	Ancestral sequence reconstruction	Stability	(Babkova et al., 2020)
Sequence-based	Cyclohexadienyl dehydratase	Ancestral sequence reconstruction	Activity	(Kaczmarek et al., 2020)
Sequence-based	β -lactamase/Kemp elimination	Ancestral sequence reconstruction	New enzyme	(Risso et al., 2017)
Data-driven	Halohydrin dehalogenase	Protein sequence activity relationships	Activity	(Fox et al., 2007)
Data-driven	Proteinase K	Linear regression algorithm	Activity and stability	(Liao et al., 2007)
Data-driven	Epoxide hydrolase	Support vector regression	Selectivity	(Zaugg et al., 2017)
Data-driven	Limonene epoxide hydrolase	Innovative sequence-activity relationship	Stability	(Li et al., 2021b)
Data-driven	Chorismate mutase	Evolution-based statistical approach	Activity	(Russ et al., 2020)
Data-driven	PET hydrolase	Three-dimensional self-supervised convolutional neural network	Activity and stability	(Lu et al., 2022)
Data-driven	Carbonic anhydrase II and Δ^5 -3-ketosteroid isomerase	Constrained hallucination and inpainting	New enzyme	(Wang et al., 2022)
Data-driven	Luciferase	Family-wide hallucination	New enzyme	(Yeh et al., 2023)

for backbone atoms, meeting the requirements for analyzing catalytic mechanisms and facilitating rational enzyme design (AlQuraishi, 2021). Structure-based rational computational design can be further classified into two subclasses based on whether it involves the computational redesign of existing enzymes or creating new enzymes with novel functions (Dinmukhamed et al., 2021).

2.1. Computational engineering of existing enzymes

In 1993, Arnold introduced the concept of directed evolution for enzymes, promoting the idea of restructuring the structure and function of natural enzymes through the random replacement of amino acid residues (Chen and Arnold, 1993). Presently, computational strategies based on virtual screening have emerged as effective means to obtain target mutants. For instance, Zhang et al. utilized both molecular dynamics (MD) simulations and binding free energy calculation in their investigation, building upon the structure and function of phosphoserine aminotransferase (SerC). They incorporated site-directed engineering strategy for virtual screening. Finally, the mutant SerC R42W:R77W exhibiting a minimum $\Delta\Delta G$ ($\Delta\Delta G_{\text{binding}} = \Delta G_{\text{binding}}(\text{mutant}) - \Delta G_{\text{binding}}(\text{wild type})$) of -0.6 kcal/mol, was experimentally obtained. L-phosphoserine was replaced by L-homoserine as the preferred substrate for SerC, representing a 4.2-fold increase in activity compared to the wildtype enzyme (Zhang et al., 2019). Similarly, Mori et al. applied a comparable strategy to tailor ferulic acid decarboxylase (FDC), targeting the unnatural substrate *cis,cis*-muconic acid (ccMA) to efficiently produce 1,3-butadiene. The best mutant, AnFDC Y394H:T395Q, displayed a remarkable 1002-fold increase in 1,3-butadiene titer (Mori et al., 2021). Recently, focused rational iterative site-specific mutagenesis (FRISM) was developed, drawing inspiration from the iterative saturation mutagenesis and the combined-active site saturation test. Differing from previous methods, the FRISM introduces only 3–5 amino acids in the hotspots, each with notable characteristics, including distinct steric, electrical, hydrophilic, and hydrophobic properties. This approach eliminates the need for saturation mutagenesis libraries and large-scale screening, thereby reducing computational resources and time while ensuring positive results (Li et al., 2021a; Li et al., 2021c). In line with the FRISM strategy, Xu et al. achieved mutants of fatty acid photodecarboxylase (CvFAP) capable of efficiently converting fatty acids with diverse chain lengths into corresponding hydrocarbons (Xu et al., 2022). Taking into account steric hindrance, they substituted ten amino acids in hotspots with A/L/F/Y to adjust the steric hindrance of the substrate binding pocket. This led to CvFAP mutants (CvFAP I398L and CvFAP P460A/G462A) exhibiting a substantial 29 to 552 times increase in enzyme activity (k_{cat}/K_m) for decarboxylation of short- and medium-fatty acids (C3–C14). The FRISM strategy also demonstrated effectiveness in designing the stereoselectivity of the enzyme. Li et al. applied rational design using FRISM to stabilize elaidic acid binding in CvFAP by introducing π - π interactions with the substrate's double bond. The V453E mutant exhibited a remarkable 1000 times enhancement in trans-over-cis selectivity than the natural enzyme (Li et al., 2021a). Additionally, several high-throughput computational enzyme engineering approaches have been reported. CADEE (computer-aided directed evolution of enzymes) employs *in silico* enzyme directed evolution to reduce the need for extensive screening and enhance the efficiency of identifying desired mutants (Amrein et al., 2017). EnzyHTP software is able to automatically handle the simulation flow, covering tasks such as structural model construction, MD simulation, quantum mechanics (QM)/molecular mechanics (MM) calculation, and modeling data analysis. This automation streamlines the identification of successful enzyme mutants in a high-throughput manner (Shao et al., 2022, 2023; Yang et al., 2023).

The rational design methods discussed above, relies on the idea of directed evolution, often require virtual iterative mutation of hot spots. In contrast, it is more resource-efficient to directly modify the local chemical micro-environment, involving factors such as the hydrogen

bond network, steric hindrance of amino acid residues, electrostatic interactions, etc., which significantly impact enzyme performance (Fig. 1). Effective hydrogen bonding within the active site and protein structure plays a vital role in stabilizing enzyme-substrate complexes, transition states, and protein folding processes. Bulky groups such as large amino acid side chains can obstruct or restrict access to the enzyme's active site, impacting substrate binding and catalytic activity. Additionally, electrostatic interactions involving charged residues and substrates can enhance catalysis by stabilizing charged transition states or forming specific binding interactions. Lin and coworkers employed 4-hydroxyphenylpyruvate dioxygenase (HPPD) as an example, manipulated its hydrogen bond network in the amide-rich zone, resulting in 100% conversion of the substrate to the desired intermediate 4-hydroxyphenylacetic acid (HPA) (Lin et al., 2021). In addition, the size of the active pocket can influence substrate conversion, particularly when amino acid residues with steric effects in the active center may impede the movement of the substrate to the reactive position. For example, the catalytic mechanism of santalene synthases was explored through multiscale simulations. The results shown that the restricted space within the active pocket of the enzyme SanSyn, caused by the key residue F441, specifically led to the specific production of the α -santalene. To enhance the active pocket space, mutant F441V was created, resulting in increased conformational change of the reaction intermediate and the generation of both α - and β -santalene (Zha et al., 2022). Unlike the previously mentioned FRISM method, this approach provides a more targeted and efficient strategy by directly identifying the crucial amino acid. This reduces the need for extensive screening and experimental validation. Moreover, Su et al. introduced a strategy for enzyme redesign towards unnatural substrates, identifying key residues with distinct conformations through restricted MD simulations. This approach was employed to construct a dehydrogenase/reductase variant, showcasing apparent conformational differences in simulations upon modifying the substrate-binding pocket. The relative sizes of the C1 and C2 pockets were tuned to enhance stereoselectivity, resulting in mutants with efficient asymmetric reduction of various substrates (Su et al., 2020). In addition, modifying the electrostatic potential or electric field applied by the enzyme to essential atoms of the substrate plays an important role in catalytic activity. This impact occurs through various mechanisms such as stabilizing the transition state of the reaction, altering the activation energy required for the transition state, promoting or preventing proton transfers during catalysis and changing the flexibility and dynamics of the enzyme's active site (Warshel et al., 2006). A quantum QM/MM dynamics strategy was used to redesign a promiscuous esterase Bs2 from *Bacillus subtilis*, aiming to improve its amidase activity. They transferred a spatially relevant aspartate residue from *Candida antarctica* lipase B (CALB) to Bs2, improving the electrostatics of transition state formation. This mutation exhibited a 1.3-fold improvement in catalytic efficiency (Galmés et al., 2022). Furthermore, integrating electric field optimization into enzyme engineering processes can yield substantial catalytic enhancements. An electric field optimization scheme was developed to improve the activity of a synthetic Kemp elimination enzyme KE15. By systematically considering the impact of strategic mutations on the local electric field along the reaction axis, researchers created a new Kemp eliminase with a nearly 50-fold increase in the k_{cat} value (Bhowmick et al., 2016; Vaissier et al., 2018). Except the electric field, substrate positioning dynamics (SPD) also serves as a significant factor for enzyme catalysis *via* aligning the substrate in a reactive conformation. The noteworthy contribution of the non-electrostatic component of SPD to regulating enzyme kinetics was recently validated by employing Kemp eliminase as an illustrative example. A distal mutant R154W created by high-throughput enzyme modeling, displayed favorable SPD. This led to an increased proportion of reactive conformations and consequently yielded the lowest activation free energy (Jiang et al., 2023).

The computational engineering of natural proteins with new functions holds the potential to revolutionize sustainable biological

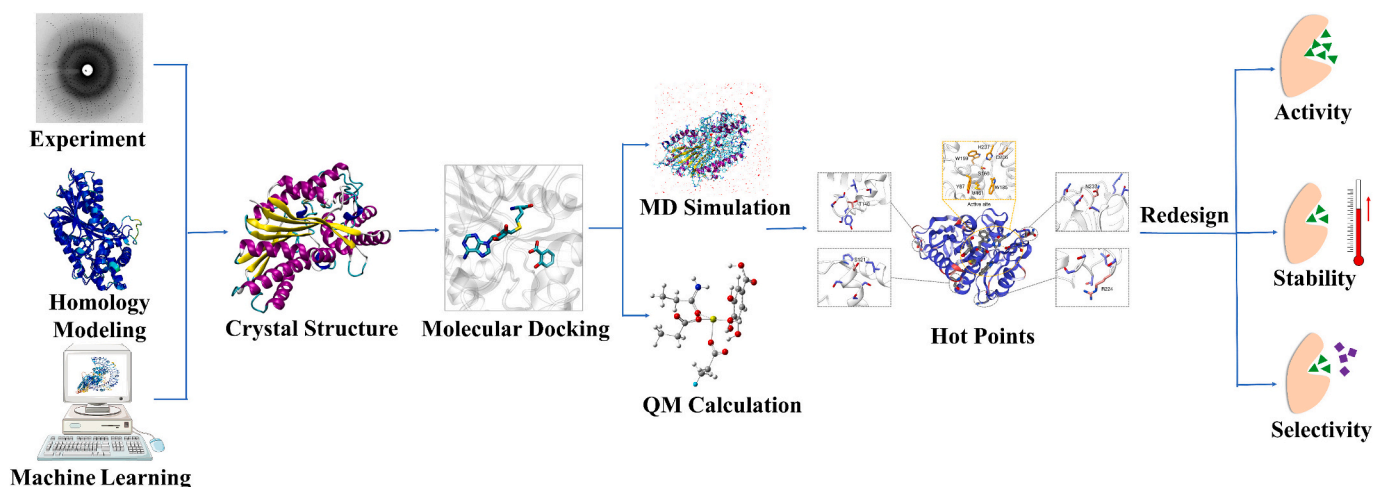


Fig. 1. Flowchart depicting the systematic process for enzyme redesign based on protein structure. Initially, protein crystal structures are acquired using experimental techniques, homologous modeling, or machine learning methodologies. Subsequently, the structural analysis is conducted to scrutinize the substrate-enzyme interaction by employing molecular docking, molecular dynamics (MD) simulations, and quantum mechanics (QM) calculations. This comprehensive examination aims to pinpoint key hot spot residues crucial for enzyme function. Finally, the target enzyme is achieved by strategically designing and incorporating modifications to the identified hot spot residues.

manufacturing. However, the strategies mentioned earlier are mainly utilized to boost or control the selectivity or activity of enzymatic reactions, proving less effective for engineered enzymes with aspirations to introduce entirely new functions. This limitation stems from the difficulty in transitioning the sequence space of the active site to entirely new protein fitness landscapes (Cui et al., 2022). Nonetheless, Rosetta redesign, which considers amino acid mechanical constraints to redesign suitable active pockets, has shown promise in creating efficient mutants with novel functions (Zanghellini et al., 2006). Rosetta redesign was first applied in the re-computational engineering of the benzaldehyde binding pocket of formolase (FLS) to enhance the catalytic efficiency for formaldehyde and the formose reaction. Following four cycles of computational design and subsequent experimental assessment, one of the 121 designed mutants, called Des1, exhibited a 26-fold increase in catalytic activity than the wildtype (Siegel et al., 2015). Despite these promising results, Rosetta redesign faces challenges that need to be addressed. First, the challenges associated with Rosetta's computational enzyme design include its random and often insufficient sampling of enzyme reaction conformations (Osuna, 2021; Tyka et al., 2012; Tyka et al., 2011; Wijma and Janssen, 2013). Additionally, the scoring function used by Rosetta compromises calculation accuracy to enhance speed, a trade-off coupled with limitations in calculation time and starting models (Bender et al., 2016). Therefore, the energy function cannot accurately describe the relevant biological state of the enzyme-catalyzed reaction. To address this, the Rosetta redesign strategy has integrated multi-state conformation design and high-throughput parallel MD simulations to enlarge the sample pool (Löffler et al., 2017). The second computational enzyme design challenge pertains to arranging catalytic residues optimally, especially in the near attack conformation (NAC), which is disrupted during dynamic catalytic processes (Ruscio et al., 2009). By relying on the NAC strategy, the catalytic selectivity by computational design (CASCO) framework addresses this challenge through the design of substrate binding sites in a predetermined orientation. This process includes generating steric hindrance to inhibit undesired substrate binding modes, and the subsequent ranking of results is carried out through high-throughput MD simulations. The application of this strategy resulted in the successful acquisition of a highly stereoselective limonene epoxide hydrolase, achieved with relatively minimal experimental work (Wijma et al., 2015). Recently, Wu et al. redesigned an aspartase AspB to enable the cross-addition of a diverse of nucleophilic amines to unsaturated acids. This process starting by docking ligands with the protein crystal structure through QM/MM

calculations. Subsequently, different docking results were simulated using MD simulations. The outcomes of the MD simulations were then utilized in the Rosetta enzyme design. The mutants produced in this process were sorted according to NAC frequency, penalty scores of constraints and total energy scores. Mutants ranked highest were chosen for experimental validation (Fig. 2). Finally, they achieved C–N lyases that exhibit cross-compatibility for unnatural nucleophiles and electrophiles, demonstrating the ability to produce diverse non-canonical amino acids with outstanding catalytic efficiency, regioselectivity, and enantioselectivity (Cui et al., 2021).

Except for enzyme activity, the enzyme stability that can maintain its functionality and structural integrity across diverse and challenging conditions is also paramount for ensuring its sustained effectiveness in industrial applications. To improve the thermostability of designed enzymes, manipulation of the electric charge distribution on the enzyme surface (Wang et al., 2020), disulfide bonds (Yang et al., 2019), rational truncation (Reich et al., 2014) and flexible loop replacement (Damjanović et al., 2014) are the most commonly used strategies. These strategies are primarily employed to decrease the entropy of unfolded proteins and improve the stability of protein conformation (Dagan et al., 2013; Kumar et al., 2004; Vieille and Zeikus, 1996). By analyzing the three-dimensional structure of *Rhizomucor miehei* lipase (RML), predictions for enhanced thermostability in RML variants can be made using tools such as Rosetta ddg_monomer, FoldX, and I-Mutant. Rosetta ddg_monomer utilizes energy-based calculations to predict mutation-induced changes in free energy, offering detailed physics-based modeling albeit requiring more expertise. FoldX employs an empirical force field to estimate protein stability changes, striking a balance between accuracy and user-friendliness. I-Mutant specializes in predicting mutations' effects on protein stability, particularly concerning protein unfolding and thermostability. If enhancing enzyme thermal stability is the goal, I-Mutant is a focused option to consider. However, for broader considerations such as structural or functional changes, or for more intricate physics-based modeling, Rosetta or FoldX may be more suitable choices. Additionally, Disulfide by Design 2 (DbD2), SSBOND, MODIP, and BridgeD were employed to identify potential disulfide regions and residues capable of forming disulfide bonds. DbD2 primarily serves protein design and engineering needs by optimizing both stability and function. SSBOND predicts and analyzes disulfide bonds crucial for stability and structure. MODIP focuses on protein-protein interactions and their impact on stability and function, while BridgeD identifies and analyzes bridging water molecules affecting

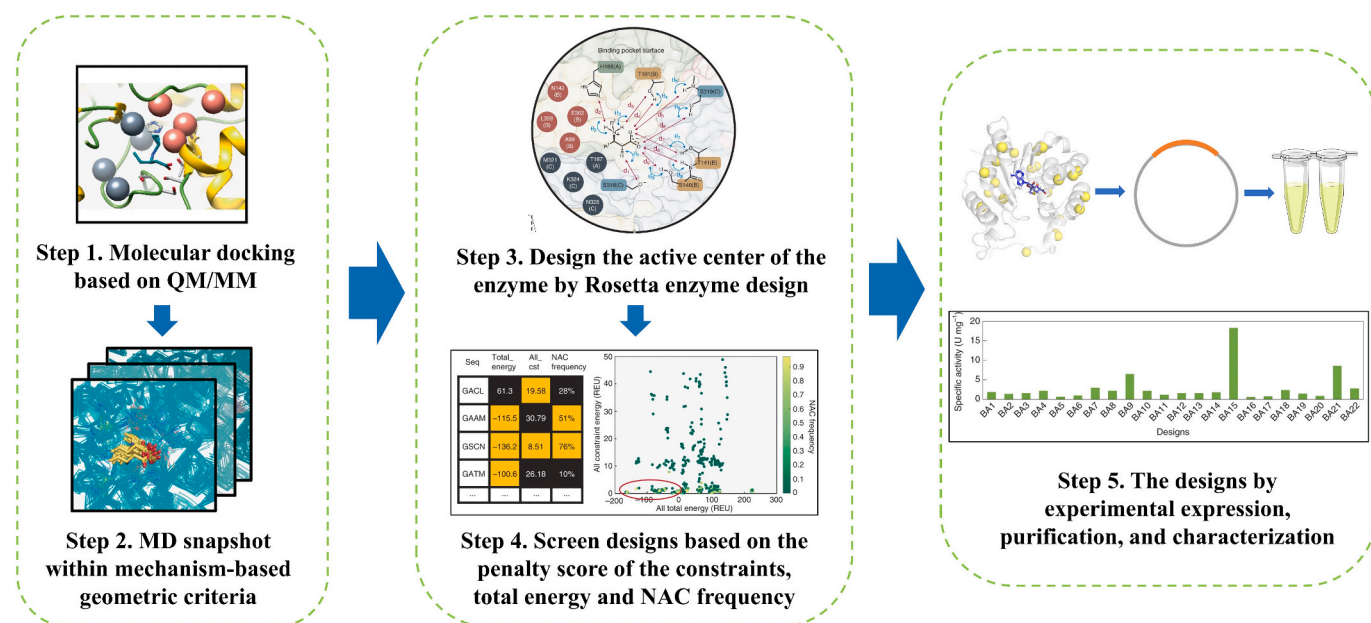


Fig. 2. Schematic representation of the enzyme redesign process employing Rosetta redesign methodology.

stability and dynamics. DbD2 and SSBOND target specific aspects like mutations or disulfide bonds, while MODIP and BridgeD offer broader insights into interactions and hydration. The most stable mutant identified in the study exhibited a remarkable 12.5 times enhancement in half-life at 70 °C, indicating a substantial enhancement in the thermal stability of RML (Li et al., 2018a). Truncation and replacement of flexible loops have proven to be effective strategies for significantly improving enzyme stability. However, identifying the specific loops that impact enzyme stability remains challenging, primarily because of the significant rearrangement of the tertiary structure. This complexity makes it difficult to pinpoint the exact loops that contribute to or influence the stability of the enzyme (Nagi and Regan, 1997; Yu et al., 2017). Han et al. employed computational-aided design method to examine the influence of the C-terminus on the functionality of GH2 fungal glucuronidase. Through structural analysis, they determined that the C-terminus of the enzyme was redundant for maintaining normal function. The C-terminus deletion induced an optimal conformation in the active site pocket, facilitating substrate binding and catalysis. Modifying the length of the C-terminal through truncation resulted in a PGUS D591:604 mutant that exhibited 3.8 times improvement in half-life at 65 °C, 2.4 times increase in enzyme activity, and 1.8 times increase in expression level. These enhancements signify significant improvements in both kinetic and thermodynamic stability of the enzyme (Han et al., 2018). Wijma and coworkers established a computational libraries framework, FRESKO, for rapid enzyme stabilization. This approach involves initially computing a large number of folding energies ($\Delta\Delta G_{\text{fold}} = \Delta G_{\text{fold}}(\text{mutant}) - \Delta G_{\text{fold}}(\text{wild type})$) to identify potential stable point mutations. Subsequently, high-throughput MD simulations and visual inspection was utilized to screen all potential point mutations. Finally, experimental validation is conducted by combining the identified mutations to significantly enhance the enzyme thermal stability. While computational design methods excel at precise mutation and structural modifications for targeted engineering purposes, FRESKO can utilize available data to inform design decisions, especially when detailed sequence and structural information are accessible. Successful applications of FRESKO in enzyme stability engineering have been demonstrated across various enzymes (Bu et al., 2018; Floor et al., 2014; Fürst et al., 2019; Meng et al., 2020; Wu et al., 2016). Especially, FRESKO was used to improve the thermostability of a homodimeric pyridoxal-5-phosphate (PLP) ω -transaminase. Numerous surface point

mutations were predicted by the computational tool. The beneficial mutations, which included P9A, E38Q, A60V, S87D/N, M128F, and I154V, resulted in improved thermostability, co-solvent resistance, isopropylamine compatibility, catalytic activity and the complete retention of enantioselectivity. These improvements were attributed to introduction of salt-bridge interactions or extra hydrogen bonding, increase of hydrophobic interactions, imposition of π -stacking, reduction of hydrophobic surface exposure, altering electrostatic surface charge distribution and the alleviation of steric strain (Meng et al., 2020).

2.2. Computational design of new enzymes with novel function

The redesign of enzymes relies on modifications to engineer existing enzymes that have functions similar to the target function or have the appropriate geometry and sufficient stability to tolerate variants needed for incorporating the target functions. However, natural enzymes contain a limited number and type of chemical reactions, which limits the potential of these enzymes to be used by the industry. Creating novel enzymes to catalyze non-natural reactions is a promising solution, which allows for breaking out from the locally optimal configurations found in natural protein sequences, installing amino acids in protein scaffolds with non-classical biological functions to enable new chemical reactions (Lu et al., 2009; Nanda and Koder, 2010).

The creating novel functional proteins with the ability to catalyze unnatural reactions is still in its early stages. Notable advancements were made by Houk and Baker in their cooperative effort to develop an “inside-out” protocol, which involves dividing the enzyme into two distinct components: a catalytic part and a scaffold part. The protocol then separately calculates and designs these two parts, following a series of general steps: 1) utilizing QM, an ideal transition state, known as theozyme, is generated within the active site. This theozyme represents the most stable configuration during the enzymatic reaction; 2) the generated transition state is transplanted onto the protein scaffold; 3) The design and optimization of the amino acid residues around the theozyme aim to guarantee favorable stacking and stability in folding (Kiss et al., 2013). Building upon this foundation, electronic and geometric properties of the transition state are incorporated to generate enzymes with novel catalytic functions. The earliest algorithms for grafting theozyme onto protein scaffolds primarily include RosettaMatch and SABER (Fig. 3). The RosettaMatch algorithm is designed to

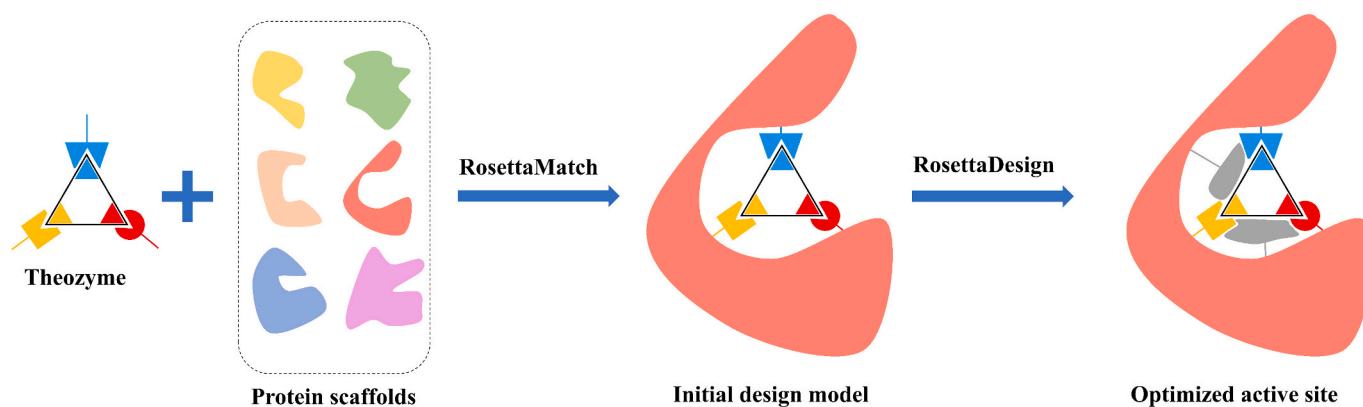


Fig. 3. Key steps in the “inside-out” protocol for enzyme *de novo* design.

identify backbone positions on a given protein scaffold and graft essential catalytic residues onto the scaffold. His process is focused on maintaining and stabilizing the reaction transition state by preserving a corresponding geometry (Richter et al., 2011; Zanghellini et al., 2006). The resulting protein, equipped with grafted catalytic residues, may demonstrate the ability to catalyze the enzymatic reaction. It's noteworthy that the RosettaMatch approach may require a higher number of mutations to adapt to new catalytic functions and substrates, potentially causing significant disruption to enzyme stability and catalytic performance. On the other hand, SABER, developed by Houk and Nosrati, offers an alternative strategy. Instead of inserting theozyme into predefined active sites, SABER searches the Protein Data Bank database to obtain native protein scaffolds with suitable catalytic functions. Once a match is found, SABER selectively mutates amino acid residues that could interfere with the correct binding of the substrate. This mutation process aims to preserve the integrity of the ideal theozyme (Kiss et al., 2013). Importantly, the success of these computational designs is ultimately validated through experimental verification (Richter et al., 2011). Siegel and coworkers employed the “inside-out” enzyme engineering strategy for asymmetric catalysis of bimolecular Diels-Alder reactions. First, the mechanism of the Diels-Alder reaction was elucidated through QM calculations. A transition state model comprising essential amino acid residues and substrates was generated. Second, the transition state model was matched to a group of protein scaffolds with high stability using RosettaMatch. The use of these scaffolds facilitated the positioning of the theozyme without inducing notable steric clashes with the protein backbone. Third, a total of 106 potential active sites were matched on a stable protein scaffold. RosettaDesign was employed to refine each match, maximizing transition state binding without conflicting with bound substrates or products. Finally, a total of 84 variants were chosen, taking into account the satisfaction of shape complementarity, catalytic geometry and binding energy. The experimental verification provided evidence that two enzymes, employing different matched protein scaffolds, were effective in catalyzing the Diels-Alder reaction. Subsequently, 6 site mutations were introduced to the candidate DA_20_00 to improve its activity. This mutational strategy successfully increased its catalytic activity. The experimental results aligned with the expected design, as DA_20_10 demonstrated the ability to catalyze the formation of 3R,4S products with >97% enantiomeric excess (Siegel et al., 2010). In addition to Diels-Alder reactions, *de novo* enzyme design has been developed for other reactions, including Kemp elimination reactions (Röthlisberger et al., 2008), Retro-Aldol reactions (Jiang et al., 2008), ester hydrolases (Richter et al., 2012), and Morita-Baylis-Hillman reactions (Bjelic et al., 2013; Burke et al., 2019; Crawshaw et al., 2022).

In recent years, advances in creating enzymes with novel functions has faced challenges, and several factors contribute to this trend. First, many unnatural chemical reactions involve intricate transition states,

and the difficulty of *de novo* design is exacerbated by the complexity of these transition states, constrained by both technical and theoretical limitations. Secondly, based on current observations, the catalytic performance of *de novo*-designed enzymes is significantly lower than that of natural enzymes, rendering them less suitable for practical applications. To address this, there is a necessity to enhance the enzyme efficiency of *de novo* enzymes through iterative rounds of directed evolution or re-rational design (Basler et al., 2021; Broom et al., 2020; Khersonsky et al., 2012; Khersonsky et al., 2011; O'Reilly, 2022; Röthlisberger et al., 2008; Siegel et al., 2010; Völler, 2020). For instance, Weitzner et al. identified a comprehensive network of hydrogen bonds involving catalytic residues in an aldolase designed after directed evolution. Recognizing the limitation of the RosettaMatch method, which primarily considers interactions between protein side chains and the transition state while overlooking hydrogen bonding between residues, they developed a side-chain optimization method called HBNNetGen. This method enhanced *de novo* enzyme design by incorporating a network of hydrogen bonds between residues in the transition state model (Obexer et al., 2017; Weitzner et al., 2019). Efficient search for protein scaffolds is crucial for successful *de novo* enzyme design. Addressing this, Zhang et al. developed ProdaMatch, an algorithm that rapidly and accurately matches protein scaffolds throughout the entire protein database (Zhang et al., 2020). To overcome the obstacle of a limited number and variety of protein scaffolds in the current database, Baker and coworkers developed an enumeration algorithm for creating novel proteins with different pocket structures and arbitrary functions, which can generate an almost infinite number of novel proteins. After two rounds of extensive experimental testing and refinement, the enhanced algorithm can produce proteins that maintain their folded structure at elevated temperatures. These proteins also showcase greater pocket diversity compared to natural NTF2-like proteins. This method enables the creation of active site geometries better suited for specific designs and promises to change the status quo of accomplishing *de novo* design of enzymes by reusing limited amounts of naturally occurring NTF2-like proteins (Basanta et al., 2020). Expanding on this foundation, Baker et al. advanced their research by integrating the catalytic site of luciferase into a computationally designed NTF2 protein scaffold, leading to the creation of a luciferase variant with specific targeting capabilities for 2-deoxycoelenterazine (Yeh et al., 2023). Except for structural enhancements in the *de novo* enzyme design process, researchers have proposed novel energy functions to assist in the design. Traditional energy functions often assume independent and additive effects for statistical energy terms, such as backbone dihedral angle, solvent exposure, and secondary structure type. Liu et al. introduced an energy function to design amino acid sequences based on a given backbone structure, named A Backbone-based Amino Acid-Usage-Survey. Derived mainly from statistics in the PDB, the ABACUS model utilizes kernel density estimation and neural network learning to represent multidimensional,

high-order correlations of energy between different structures in known protein structures. This approach introduces neural network-based statistical energy terms that capture multidimensional features not easily described by traditional statistical methods (Huang et al., 2022).

3. Sequence-based enzyme computational engineering and design

Both homology modeling and deep learning-based protein structure prediction approaches are capable of providing accurate and reliable three-dimensional structures for proteins. However, the effectiveness of these methods relies on the presence of resolved crystal structures for homologous proteins in the database. If an appropriate template for homology modeling is unavailable or of poor quality, the prediction may be compromised. Additionally, obtaining the crystal structure experimentally is a time-consuming and labor-intensive process. Furthermore, once the protein structure is acquired, structure-based computational design methods demand a significant understanding of theoretical chemistry and computational techniques to explore chemical mechanisms and transition states. This could pose challenges for researchers lacking expertise in these areas. In contrast, sequence-based computational design methods offer a partial solution to these challenges. These methods allow the direct design of enzymes from protein sequences, bypassing the need for explicit structural information. This approach enables researchers to glean insights from the natural evolution of enzymes over billions of years. With the wealth of sequence data made available by next-generation sequencing methods, researchers can employ phylogenetic analysis to uncover the underlying principles of enzyme evolution (Stephens et al., 2015).

Consensus Design (CD) and Ancestral Sequence Reconstruction (ASR) are two of the most widespread sequence-based computational design methods (Fig. 4) (Musil et al., 2018). CD relies on homologous protein sequences, typically utilizing tens to hundreds of such sequences to generate a Multiple Sequence Alignment (MSA). The distribution frequency of each amino acid at a specific position is then calculated, and a user-defined conservation threshold is applied to identify “consensus” amino acids during evolution. The underlying assumption is that the most common amino acid at a specific position is more likely to be stable and evolutionarily conserved, potentially affecting enzyme catalysis (Magliery, 2015; Porebski and Buckle, 2016; Steipe et al., 1994). Teze et al. introduced a sequence-based method for creating highly efficient glycosyltransferases (GTs) from retaining glycoside hydrolases (GHs) without relying on enzyme structure information. The method involves collecting a substantial number of sequences within the GH family, clustering to reduce redundancy, and iteratively performing

MSA to choose sequences with >5% identity. The conserved residues are identified through iterative increases in the sequence identity threshold, and the most conserved residues in GH are replaced with structural analogs. This method has proven to be effective in engineering of enzymes from various GH families (Teze et al., 2021). To address the challenge that conserved amino acid residues may not always be ideal mutation targets, Sullivan et al. enhanced the consensus design by removing nearly invariant positions and sites exhibit high statistical correlations with other positions. This modification improved the identification rate of stable mutations in *Saccharomyces cerevisiae* triosephosphate isomerase from 50% to 90% (Sullivan et al., 2012). Except for consensus design, structure-guided consensus design methods leverage structural information to identify potentially damaging mutations. Vázquez-Figueroa et al. proposed a method incorporating distance between potential mutation sites and active sites, secondary structure information, and the overall count of intramolecular contacts. This method led to a substantial improvement in the stability of mutant glucose-1-dehydrogenase, with a 10^6 -fold improvement when combining mutation sites (Vázquez-Figueroa et al., 2007). Another efficient structure-guided consensus design method involves analyzing molecular fluctuations by considering crystallographic B-factors (Parthasarathy and Murthy, 2000). Despite its strengths, consensus design has limitations, such as an inability to explain epistasis between mutation results and noticeable phylogenetic bias when certain subfamilies dominate the MSA (Aerts et al., 2013; Cole and Gaucher, 2011; Hochberg and Thornton, 2017; Lehmann et al., 2000).

Proteins exhibit diverse catalytic functions, often evolving from pre-existing functions. However, the mechanisms underlying the emergence of new enzyme functions and their evolutionary processes remain inadequately understood. ASR employs the Bayesian inference or maximum likelihood method to explore the evolutionary history of homologous sequences from known MSA and appropriate phylogenetic models. This method infers the replacement probability of a defined amino acid at a specific site in the target enzyme and reconstructs the phylogenetic tree containing putative ancestral sequences (Selberg et al., 2021; Spence et al., 2021). Despite the inherent ambiguity in ASR, where even the most likely sequence obtained may not accurately represent the true ancestral sequence in evolutionary history, the method remains valuable. A cytochrome P450 was engineered by the ASR and the optimal mutant exhibited improved solvent stability and thermostability than a human cytochrome P450, despite sharing similar reactivity profiles. Additionally, this work also created a mutant ketol-acid reductoisomerase that showed an 8-fold increase in catalytic efficiency (Gumulya et al., 2018). ASR has also been instrumental in exploring the sequence-activity landscape of a self-sufficient P450

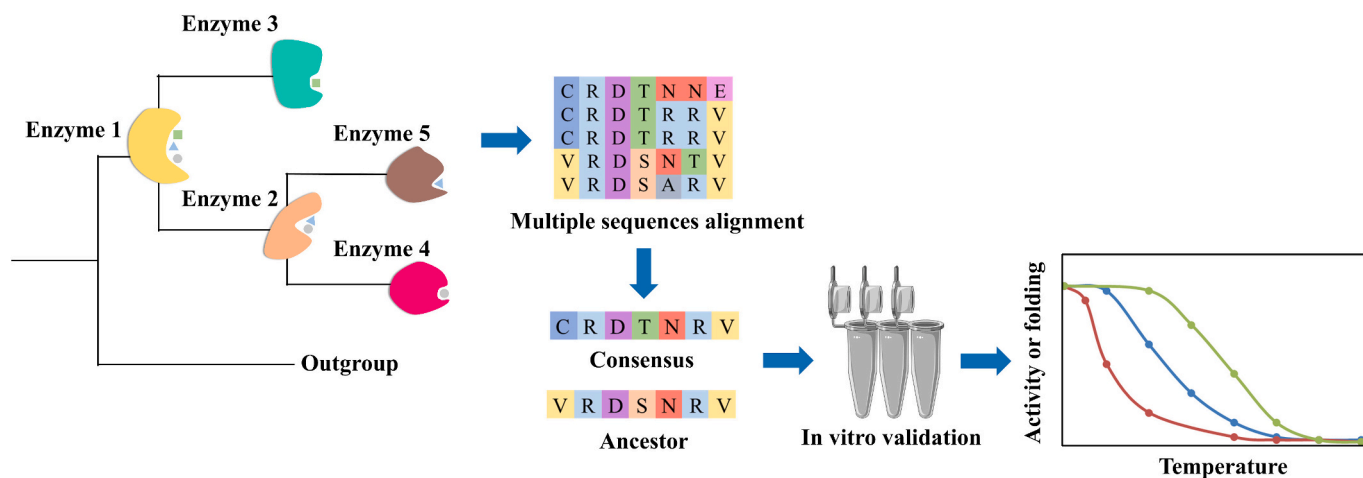


Fig. 4. Overview of consensus design (CD) and ancestral sequence reconstruction (ASR), illustrating the general workflow that encompasses the collection and alignment of homologous sequences, phylogenetic analysis, and experimental characterization.

monooxygenase CYP116B. It identified crucial mutational paths leading to changes in enzyme selectivity, resulting in ancestor variants with shifted regioselectivity from terminal to mid-chain hydroxylation of fatty acids (Jones et al., 2024). However, ASR faces challenges due to the high differentiation among extant proteins, leading to uncertainties in the predicted probabilities of many sites, often falling below 50%. This limitation underscores the difficulty in confidently determining the real ancestral sequence through ASR (Copley, 2021). However, despite these challenges, the protein scaffold obtained through ASR serves as an important starting point for protein design efforts owing to the inherent high thermostability and evolvability. ASR-derived protein scaffolds play a crucial role in studying protein conformational dynamics related to thermostability (Babkova et al., 2020), ligand binding specificity, allosteric regulation and catalytic activity (Jemth et al., 2018; Kaczmarek et al., 2020).

Recently, an escalating number of researchers have recognized the significant implications of conformational dynamics for the catalytic promiscuity and evolution of enzymes (Babkova et al., 2020; Campitelli et al., 2020; Crean et al., 2020; Gardner et al., 2020; James and Tawfik, 2003; Maria-Solano et al., 2018; Spence et al., 2021; St-Jacques et al., 2023; Zamora et al., 2020). James and Tawfik were among the first to propose that conformational plasticity contributes in enzyme evolution, enabling enzymes to bind new ligands and adopt conformations that catalyze novel chemical reactions (James and Tawfik, 2003). These conformations were generally rare in ancestral enzymes, but natural evolutionary pressures favored the prevalence of such rare conformations, leading to the creation of new catalytic reactions for previously unexplored substrates. Laboratory-based directed evolution can facilitate the generation of these new conformations more readily than in natural evolution, as the lower level of catalytic promiscuity can be detected and amplified in a controlled laboratory environment. Currently, two methods have been developed to fine-tune conformational dynamics for enhancing enzyme activity: 1) expanding the sampling of states that can impart new catalytic activity; 2) decreasing the sampling of non-productive conformations (Campbell et al., 2018). Recent research has demonstrated that ASR can provide a starting scaffold for modifying catalytic properties, with results that can be explained by conformational dynamics. Risso et al. reconstructed Precambrian β -lactamase using ASR as a scaffold for protein engineering. They designed an artificial Kemp eliminase with catalytic activity up to 10^7 higher than the uncatalyzed reaction, along with significant ester hydrolysis activity. Both experimental and computational analyses revealed that despite substantial differences in amino acid sequences and protein scaffolds across evolutionary time, the overall tertiary structure remained relatively unchanged (Risso et al., 2017). The “Shortest Path Map (SPM)” method is a notable advancement in computational enzyme engineering. This method utilizes a deep understanding of conformational dynamics and distal mutations, complemented by long-time-scale molecular dynamics (MD) simulations, to explore a wide range of conformations and predict mutations that can induce an allosteric impact on protein function. For instance, the direct evolution of a retro-aldolase enzyme led to a variant exhibiting a remarkable $>10^9$ rate enhancement. Through the analysis of enzyme active site conformational dynamics, it was discovered that long-range mutations played a crucial role in shifting populations of conformational states towards active states. The key amino acids identified through the SPM method were consistent with those mutated in the direct evolution process (Romero-Rivera et al., 2017). The SPM method also has proven successful in understanding of the catalytic mechanism of enzyme reactions. An engineered cytochrome P450 monooxygenase exhibited regio- and stereoselective hydroxylation activities towards steroid. The mechanism was investigated by SPM and the results revealed that the epistatic effects and conformational dynamics are influenced by distal interactions in loops, β -strands and helices, which control the substrate access tunnel. This regulation ultimately facilitates optimal catalysis (Acevedo-Rocha et al., 2021). Similarly, the molecular

mechanism of tryptophan synthase for stand-alone activity (Maria-Solano et al., 2019; Maria-Solano et al., 2021) and monoamine oxidase for broadening substrate scope (Curado-Carballada et al., 2019) were also elucidated using this approach. Studying the evolutionary history of native proteins enhances our understanding of protein engineering, emphasizing the importance of considering protein dynamics, neo-functionalization, and epistatic interactions in the design and engineering of proteins. Recently, several user-friendly web-based tools have been developed for improving enzyme performance without requiring extensive expertise or installation. For instance, Caver is a specialized tool used for the detailed analysis and visualization of substrate access tunnels and channels in protein structures (Stourac et al., 2019). Hot-Spot Wizard is designed to automatically identify “hotspots” in proteins, facilitating the engineering of substrate specificity, enzyme activity, or enantioselectivity (Pavelka et al., 2009). The Rosetta-based PROSS is specifically used for designing enzymes to ensure stability and high expression levels (Gomez de Santos et al., 2023). FuncLib is designed to generate multipoint mutations, specifically honing in on the active sites of enzymes (Kheronsky et al., 2018).

4. Data-driven enzyme engineering and design

Indeed, the structure-based and sequence-based computational design approaches can offer powerful and complementary strategies for engineering of enzymes with desired properties. Despite the promising results shown by both methods, our knowledge of the sequence and structure of native enzymes is still in its infancy. Consequently, the computational design of enzymes, based on limited sequence or structural information, necessitates a substantial computational or experimental investment, thereby amplifying the overall cost of enzyme engineering. In the vast realm of nature, there exists substantial, untapped sequence and structural space. Efficient exploration of these uncharted territories within protein sequence space holds the potential for the discovery of novel and valuable enzymes. Integrating machine learning (ML) into enzyme engineering becomes pivotal, as it empowers the modeling of intricate sequence-function and structure-function relationships using existing data. ML facilitates the prediction of new, valuable enzymes that may prove challenging to attain through conventional methods. This approach not only enhances our understanding of these relationships but also guides the design of highly efficient enzymes, closely mirroring the efficacy observed in natural systems. Various ML algorithms have been employed for enzyme engineering, showcasing their versatility. For instance, random forests have been employed for predicting protein solubility (Yang et al., 2016; Yang et al., 2021). Support vector machines and decision trees have found utility in predicting changes in enzyme stability after mutation (Folkman et al., 2016; Teng et al., 2010) (Huang et al., 2007). K-nearest neighbor methods, such as the K-Nearest Neighbor Classifier, have been utilized for predicting enzyme functions and catalytic mechanisms (De Ferrari and Mitchell, 2014; Koskinen et al., 2015). Moreover, diverse scoring and clustering algorithms have been utilized for the swift annotation of functional sequences (Cozzetto et al., 2013; Falda et al., 2012). The main allure of ML in protein engineering lies in its ability to provide rapid predictions once trained on a suitable dataset. In contrast, rational design necessitates the construction of new models and extended periods of intensive computation, while directed evolution methods often require months of experimental works. The efficacy of ML relies on the quality of the training dataset and the efficiency of the underlying algorithms. Challenges such as the need for rigorous control over data collection and reporting, difficulties in standardizing data formats, the lack of large homogeneous datasets for training, slow establishment of new datasets for model testing, and the diversity of reactions, catalytic mechanisms and experimental conditions, significantly limit the widespread application of ML in designing biocatalysts (Fig. 5).

Supervised ML are commonly employed in enzyme engineering, where researchers represent enzyme fitness with a given label. The

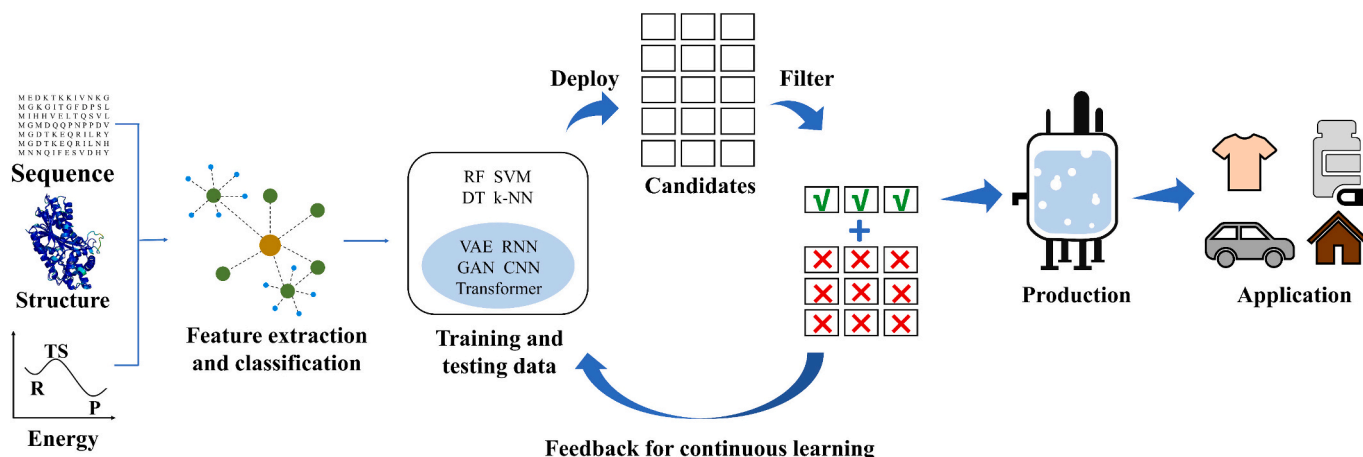


Fig. 5. Employ diverse machine learning techniques to categorize and train the various crucial parameters of enzymes, thereby obtaining corresponding models. Subsequently, apply these models to the systematic and rational design of enzymes, ultimately resulting in the delineation of the process for obtaining target enzymes.

choice of features and models distinguishes various ML methods. Generally, selected features in enzyme engineering can include protein sequences (Foroozandeh Shahraki et al., 2021; Greenhalgh et al., 2021; Vasina et al., 2022), structures (Feehan et al., 2021; Krivák and Hoksza, 2018; Unger et al., 2020), or energy data (Li et al., 2022; Pan et al., 2021; Xie et al., 2022). Sequence-based features have been predominantly chosen for enzyme engineering, allowing for the rapid engineering of sequence features and the generation of predictable feature models *in silico*. Despite the advantages of sequence-based ML in acquiring a large amount of data for building regression models, numerical transformation of sequence information is still required for accurate enzyme engineering predictions (Bonetta and Valentino, 2020). Fox et al. advanced directed evolution by introducing a statistical analysis strategy called ProSAR (Protein Sequence-Activity Relationship), which assigns specific weights (regression coefficients) to positions based on a simple transformation of the absence (0) or presence (1) of substitution. These weights are determined through statistical analysis of the initial dataset and the impact of each substitution. The ProSAR method effectively identifies beneficial mutations even in mutants with reduced enzyme function. By employing ProSAR, the volumetric productivity of the cyanation process increased approximately 4000-fold in a library of aldol dehalogenase mutants (Fox et al., 2007). Liao et al. proposed a similar method considering amino acid mutation weighting and recombination effect summarization. They tested eight different linear regression algorithms on the dataset and, after two rounds of design, achieved a mutant with enzyme activity 20 times higher than the wildtype among 95 proteinase K mutants (Liao et al., 2007). However, linear models may not be as effective in describing higher-order epistasis. Zaugg et al. constructed a model using support vector regression (SVR) on 136 mutant data of epoxide hydrolase enantioselectivity. Exploring various SVR kernels and utilizing kernel tricks, they found that nonlinear models outperformed linear models in predicting protein function (Zaugg et al., 2017). Li et al. efficiently identified highly robust limonene epoxide hydrolase variants by identifying epimutation interactions through the innov'SAR (innovative sequence-activity relationship) algorithm. Using Fourier transform (FT) to capture nonlinear interactions between the order and position of protein sequences, they obtained variants with higher unfolding stability and tolerance to aggregation (Li et al., 2021b). Russ et al. described a method for generating new artificial sequences with protein family properties solely from evolutionary sequence data. Considering the conservation of amino acid positions and the relatedness of amino acid pairs in evolution, they experimentally demonstrated that the predicted new artificial sequences exhibit similar catalytic functions to chorismate mutase enzymes. This work suggests that evolution-based statistical models are sufficient to

describe the extensive functional sequence space of a specific enzyme, laying the foundation for evolution-based artificial protease design (Russ et al., 2020).

Deep learning simulates the human brain for analysis and learning *via* establishing neural networks. It offers several advantages over traditional methods. Firstly, deep learning stands out for its capacity to process vast amounts of data, yielding more precise predictions. Secondly, deep learning eliminates the need for intricate feature engineering, requiring only sequence, structure, or energy data to be directly input into the neural network for effective predictive performance. This eradicates the challenges associated with extensive feature engineering in the modeling process. Lastly, adaptive deep learning can seamlessly transition across different types of datasets, and transfer learning enables pre-trained deep neural networks to be applied to various applications within the same domain. However, deep learning has its limitations. It necessitates very large datasets, and obtaining enzyme data for enzyme engineering is both expensive and time-consuming. In cases where the dataset is small, classical machine learning tends to outperform deep learning. As biological information databases rapidly expand, an increasing number of deep learning algorithms are being applied to the computational design of enzymes, including Variational Autoencoders (VAE) (Doersch, 2016; Hawkins-Hooker et al., 2021; Lobzaev et al., 2022), Recurrent Neural Network (RNN) (Alley et al., 2019; Hawkins-Hooker et al., 2021; Lipton et al., 2015), Generative Adversarial Network (GAN) (Repecka et al., 2021), Convolutional Neural Networks (CNN) (Lu et al., 2022) and Transformer (Dubourg-Felonneau et al., 2021; Vaswani et al., 2017). Alley et al. employed a RNN to acquire statistical characterizations of proteins using 24 million UniRef50 sequences. This approach condenses arbitrary protein sequences into fixed-length backbones, estimating essential protein properties independently of structural or evolutionary data (Alley et al., 2019). By leveraging the original protein sequence, this approach addressed data scarcity issues in protein informatics and showcased exceptional performance in crucial engineering tasks, encompassing stability, function, and design. UniRep, based on a relatively small sequence dataset, demonstrated effectiveness in constructing precise models representing a protein's fitness landscape. Even with a limited number of functionally characterized mutants, this unsupervised deep learning model facilitated large-scale exploration of sequence space and achieved protein design optimization comparable to mutants obtained in previous high-throughput studies (Biswas et al., 2021). Lu and co-workers employed a self-supervised convolutional neural network operating in three dimensions for training nearly twenty thousand protein structures. The model acquired knowledge of the local chemical microenvironment surrounding amino acids, enabling the prediction of

locations in proteins where wild-type amino acids were not present. Through a computer-synthesized mutation scan and a stepwise combination strategy, this approach generated mutants with significant catalytic improvement (Lu et al., 2022). A generative maximum entropy model was used to optimize the Renilla luciferase enzyme. Their work illustrates that leveraging natural evolutionary information enables predictive improvement of enzyme stability and activity through the design of active sites and protein scaffolds (Xie et al., 2023). A deep learning model EnzyKR was developed to predict biocatalysts with high enantioselectivity. In contrast to existing k_{cat} predictors, EnzyKR integrates substrate chirality by utilizing dihedral angles, geometric features, and atomic distance maps derived from hydrolase-substrate pairs. The capability of EnzyKR to efficiently identify enzymes with high stereoselectivity is attributed to its consideration of variations in k_{cat} values between different enantiomers, leading to distinct stereoselectivities (Ran et al., 2023). For creating novel proteins, Norn et al. reversed the supervised model used for protein structure prediction into a model based on structural design sequences (Norn et al., 2020). Anishchenko et al. utilized deep neural networks to optimize randomly given amino acid sequences, generating new proteins with diverse sequences and predicted structures (Anishchenko et al., 2021). These methods, based on inverting deep networks trained to predict the natural structures of proteins, offer insights into *de novo* protein design, complementing traditional physics-based models. Recently, Wang et al. introduced two deep learning approaches for creating proteins with pre-specified functional sites. “Constrained hallucination” refers to a method that begins with a desired or predicted structure or functional site and optimizes sequences to match or contain these features. This approach imposes specific constraints or requirements during sequence optimization to ensure that certain features or properties are maintained or achieved. “Inpainting” starts from a known or predicted functional site and expands sequences or structures to fill in missing information or enhance specific elements using a RoseTTAFold network trained on the PDB database. It focuses on complementing or enhancing existing features, often without strict constraints on other parts of the sequence or structure. The authors successfully designed proteins with specific functions, including immunogens, receptor traps, metalloproteins, and enzymes, showcasing the capability to obtain desired protein scaffolds with functional sites (Wang et al., 2022). Using data-driven methods for enzyme engineering presents challenges and limitations. One primary issue is that the accuracy and robustness of these methods are limited by the quality and quantity of available data. To tackle this challenge, developing robust data-cleaning and data augmentation techniques can be beneficial. These strategies aim to remove errors, duplicates, and inconsistencies from the data, while also enhancing data diversity and representativeness, thus improving the overall reliability of data-driven approaches. Moreover, processing large-scale datasets and training complex models can be a time-consuming and resource-intensive task, which makes data-driven methods challenging for broad application. Developing parallel processing techniques and distributed computing infrastructure can improve the efficiency and scalability of data-driven approaches in enzyme engineering, making them more practical for real-world applications.

5. Conclusions and perspectives

In the review, we have delved into three major approaches to rational computational enzyme design, offering insights into their recent advancements. While each method provides a robust strategy for enzyme engineering, it is crucial to acknowledge the challenges that persist in the field. The foundation of structure-based computational design, which posits that protein structure dictates function, has led to significant progress. Through meticulous simulations and algorithmic analyses, this approach guides the targeted modification of natural enzymes for desired outcomes. Despite its success, challenges arise from its limited ability to explore the vast mutation sample space

comprehensively. The nascent field of *de novo* design faces hurdles such as constrained conformational sampling, insufficient protein scaffold diversity, and the need for improved accuracy in transition state descriptions. Sequence-based rational computational design starts from protein amino acid sequences, employing methods like consensus design and ancestral sequence reconstruction to leverage nature’s evolutionary skills for enzyme transformation. However, this approach may struggle to explain epistasis between mutations, especially those occurring at a distance. The bias towards natural substrates in evolutionary optimization poses limitations when aiming to improve the catalytic efficiency of non-natural substrates. To overcome these challenges, combining sequence and structural information is essential. The emergence of data-driven artificial intelligence and machine learning leverages vast amounts of sequence, structural, and energy data for efficient exploration of the combinatorial space of diverse enzyme sequences. Notably, deep learning methods, which can predict protein structure directly from sequence information, enhance retrofitting accuracy. However, challenges persist, particularly in *de novo* designing proteins with functional sites. Overcoming hurdles, such as embedding functional sites into designed protein scaffolds and predicting correctly folded amino acid sequences with functional sites, remains a priority. While data-driven machine learning methods hold promise for the future of enzyme design, current structure- and sequence-based methods remain crucial for elucidating catalytic mechanisms and rational enzyme design. Other challenges encountered in computational enzyme engineering have been thoroughly reviewed in recent articles (Kouba et al., 2023; Yang et al., 2023). In the future, the integration of machine learning with these traditional methods is vital to address these issues. This interdisciplinary approach aims to design biocatalysts with entirely new catalytic reactions, breaking through the bottleneck of limited reaction types and fostering advancements in synthetic biology. In conclusion, the synergistic application of these rational computational design approaches will likely lead to transformative breakthroughs in enzyme engineering, ultimately contributing to the flourishing fields of synthetic biology and metabolic engineering.

Declaration of competing interest

None.

Data availability

No data was used for the research described in the article.

Acknowledgments

This work was funded by National Key Research and Development Program of China (2022YFA2106100) and National Natural Science Foundation of China (22378016, 22078011 and 22238001).

References

- Acevedo-Rocha, C.G., Li, A., D’Amore, L., Hoebenreich, S., Sanchis, J., Lubrano, P., Ferla, M.P., Garcia-Borràs, M., Osuna, S., Reetz, M.T., 2021. Pervasive cooperative mutational effects on multiple catalytic enzyme traits emerge via long-range conformational dynamics. *Nat. Commun.* 12 (1), 1–13. <https://doi.org/10.1038/s41467-021-21833-w>.
- Aerts, D., Verhaeghe, T., Joosten, H.J., Vriend, G., Soetaert, W., Desmet, T., 2013. Consensus engineering of sucrose phosphorylase: the outcome reflects the sequence input. *Biotechnol. Bioeng.* 110 (10), 2563–2572. <https://doi.org/10.1002/bit.24940>.
- Alley, E.C., Khimulya, G., Biswas, S., AlQuraishi, M., Church, G.M., 2019. Unified rational protein engineering with sequence-based deep representation learning. *Nat. Methods* 16 (12), 1315–1322. <https://doi.org/10.1038/s41592-019-0598-1>.
- Almeida, T., Silvestre, A.J., Vilela, C., Freire, C.S., 2021. Bacterial nanocellulose toward green cosmetics: recent progresses and challenges. *Int. J. Mol. Sci.* 22 (6), 2836. <https://doi.org/10.3390/ijms22062836>.
- AlQuraishi, M., 2021. Protein-structure prediction revolutionized. *Nature* 596 (7873), 487–488. <https://doi.org/10.1038/d41586-021-02265-4>.

- Vázquez-Figueroa, E., Chaparro-Riggers, J., Bommarius, A.S., 2007. Development of a thermostable glucose dehydrogenase by a structure-guided consensus concept. *ChemBioChem* 8 (18), 2295–2301. <https://doi.org/10.1002/cbic.200700500>.
- Vieille, C., Zeikus, J.G., 1996. Thermozymes: identifying molecular determinants of protein structural and functional stability. *Trends Biotechnol.* 14 (6), 183–190. [https://doi.org/10.1016/0167-7799\(96\)10026-3](https://doi.org/10.1016/0167-7799(96)10026-3).
- Völler, J.-S., 2020. Ensemble-based enzyme design. *Nat. Catal.* 3 (10), 774. <https://doi.org/10.1038/s41929-020-00529-2>.
- Wang, Y., Xu, M., Yang, T., Zhang, X., Rao, Z., 2020. Surface charge-based rational design of aspartase modifies the optimal pH for efficient β -aminobutyric acid production. *Int. J. Biol. Macromol.* 164, 4165–4172. <https://doi.org/10.1016/j.ijbiomac.2020.08.229>.
- Wang, J., Lisanza, S., Juergens, D., Tischer, D., Watson, J.L., Castro, K.M., Ragotte, R., Saragovi, A., Milles, L.F., Baek, M., Anishchenko, I., Yang, W., Hicks, D.R., Expósito, M., Schlichthaerle, T., Chun, J.-H., Dauparas, J., Bennett, N., Wicky, B.I.M., Muenks, A., DiMaio, F., Correia, B., Ovchinnikov, S., Baker, D., 2022. Scaffolding protein functional sites using deep learning. *Science* 377 (6604), 387–394. <https://doi.org/10.1126/science.abn2100>.
- Warshel, A., Sharma, P.K., Kato, M., Xiang, Y., Liu, H., Olsson, M.H., 2006. Electrostatic basis for enzyme catalysis. *Chem. Rev.* 106 (8), 3210–3235. <https://doi.org/10.1021/cr0503106>.
- Weitzner, B.D., Kipnis, Y., Daniel, A.G., Hilvert, D., Baker, D., 2019. A computational method for design of connected catalytic networks in proteins. *Protein Sci.* 28 (12), 2036–2041. <https://doi.org/10.1002/pro.3757>.
- Wijma, H.J., Janssen, D.B., 2013. Computational design gains momentum in enzyme catalysis engineering. *FEBS J.* 280 (13), 2948–2960. <https://doi.org/10.1111/febs.12324>.
- Wijma, H.J., Floor, R.J., Bjelic, S., Marrink, S.J., Baker, D., Janssen, D.B., 2015. Enantioselective enzymes by computational design and *in silico* screening. *Angew. Chem. Int. Ed.* 54 (12), 3726–3730. <https://doi.org/10.1002/anie.201411415>.
- Wu, B., Wijma, H.J., Song, L., Rozeboom, H.J., Poloni, C., Tian, Y., Arif, M.I., Nuijens, T., Quaedflieg, P.J.L.M., Szymanski, W., Feringa, B.L., Janssen, D.B., 2016. Versatile peptide C-terminal functionalization via a computationally engineered peptide amidase. *ACS Catal.* 6 (8), 5405–5414. <https://doi.org/10.1021/acscatal.6b01062>.
- Xie, W.J., Asadi, M., Warshel, A., 2022. Enhancing computational enzyme design by a maximum entropy strategy. *Proc. Natl. Acad. Sci. USA* 119 (7), e2122355119. <https://doi.org/10.1073/pnas.2122355119>.
- Xie, W.J., Liu, D., Wang, X., Zhang, A., Wei, Q., Nandi, A., Dong, S., Warshel, A., 2023. Enhancing luciferase activity and stability through generative modeling of natural enzyme sequences. *Proc. Natl. Acad. Sci.* 120 (48), e2312848120 <https://doi.org/10.1073/pnas.2312848120>.
- Xu, E., Campanella, O.H., Ye, X., Jin, Z., Liu, D., BeMiller, J.N., 2020. Advances in conversion of natural biopolymers: A reactive extrusion (REX)-enzyme-combined strategy for starch/protein-based food processing. *Trends Food Sci. Technol.* 99, 167–180. <https://doi.org/10.1016/j.tifs.2020.02.018>.
- Xu, W., Chen, Y., Li, D., Wang, Z., Xu, J., Wu, Q., 2022. Rational design of fatty acid photodecarboxylase enables the efficient decarboxylation of medium-and short-chain fatty acids for the production of gasoline bio-alkanes. *Mol. Catal.* 524, 112261 <https://doi.org/10.1016/j.mcat.2022.112261>.
- Yang, Y., Niroula, A., Shen, B., Vihinen, M., 2016. PON-sol: prediction of effects of amino acid substitutions on protein solubility. *Bioinformatics* 32 (13), 2032–2034. <https://doi.org/10.1093/bioinformatics/btw066>.
- Yang, M., Yang, S.-X., Liu, Z.-M., Li, N.-N., Li, L., Mou, H.-J., 2019. Rational design of alginate lyase from *Microbulbifer* sp. Q7 to improve thermal stability. *Mar. Drugs* 17 (6), 378. <https://doi.org/10.3390/md17060378>.
- Yang, Y., Zeng, L., Vihinen, M., 2021. PON-Sol2: prediction of effects of variants on protein solubility. *Int. J. Mol. Sci.* 22 (15), 8027. <https://doi.org/10.3390/ijms22158027>.
- Yang, Z.J., Shao, Q., Jiang, Y., Jurich, C., Ran, X., Juarez, R.J., Yan, B., Stull, S.L., Gollu, A., Ding, N., 2023. Mutexa: a computational ecosystem for intelligent protein engineering. *J. Chem. Theory Comput.* 19 (21), 7459–7477. <https://doi.org/10.1021/acs.jctc.3c00602>.
- 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., 2023. *De novo* design of luciferases using deep learning. *Nature* 614 (7949), 774–780. <https://doi.org/10.1038/s41586-023-05696-3>.
- Yu, H., Yan, Y., Zhang, C., Dalby, P.A., 2017. Two strategies to engineer flexible loops for improved enzyme thermostability. *Sci. Rep.* 7, 1–15. <https://doi.org/10.1038/srep41212>.
- Zamora, R.A., Ramirez-Sarmiento, C.A., Castro-Fernández, V., Villalobos, P., Maturana, P., Herrera-Morande, A., Komives, E.A., Guixé, V., 2020. Tuning of conformational dynamics through evolution-based design modulates the catalytic adaptability of an extremophilic kinase. *ACS Catal.* 10 (19), 10847–10857. <https://doi.org/10.1021/acscatal.0c01300>.
- Zanghellini, A., Jiang, L., Wollacott, A.M., Cheng, G., Meiler, J., Althoff, E.A., Röthlisberger, D., Baker, D., 2006. New algorithms and an *in silico* benchmark for computational enzyme design. *Protein Sci.* 15 (12), 2785–2794. <https://doi.org/10.1110/ps.062353106>.
- Zaugg, J., Gumulya, Y., Malde, A.K., Bodén, M., 2017. Learning epistatic interactions from sequence-activity data to predict enantioselectivity. *J. Comput. Aided Mol. Des.* 31 (12), 1085–1096. <https://doi.org/10.1007/s10822-017-0090-x>.
- Zha, W., Zhang, F., Shao, J., Ma, X., Zhu, J., Sun, P., Wu, R., Zi, J., 2022. Rationally engineering santalene synthase to readjust the component ratio of sandalwood oil. *Nat. Commun.* 13 (1), 2508. <https://doi.org/10.1038/s41467-022-30294-8>.
- Zhang, Y., Ma, C., Dischert, W., Soucaille, P., Zeng, A.P., 2019. Engineering of phosphoserine aminotransferase increases the conversion of L-homoserine to 4-hydroxy-2-ketobutyrate in a glycerol-independent pathway of 1, 3-propanediol production from glucose. *Biotechnol. J.* 14 (9), 1900003. <https://doi.org/10.1002/biot.201900003>.
- Zhang, S., Zhang, J., Zhu, Y., 2020. ProdaMatch: A fast and accurate active site matching algorithm for *de novo* enzyme design. *Comput. Chem. Eng.* 140, 106921 <https://doi.org/10.1016/j.compchemeng.2020.106921>.