

## Protein Engineering

## Computational Stabilization of a Non-Heme Iron Enzyme Enables Efficient Evolution of New Function

Brianne R. King, Kiera H. Sumida, Jessica L. Caruso, David Baker, and Jesse G. Zalatan\*

**Abstract:** Deep learning tools for enzyme design are rapidly emerging, and there is a critical need to evaluate their effectiveness in engineering workflows. Here we show that the deep learning-based tool ProteinMPNN can be used to redesign Fe(II)/ $\alpha$ KG superfamily enzymes for greater stability, solubility, and expression while retaining both native activity and industrially relevant non-native functions. This superfamily has diverse catalytic functions and could provide a rich new source of biocatalysts for synthesis and industrial processes. Through systematic comparisons of directed evolution trajectories for a non-native, remote C(sp<sup>3</sup>)-H hydroxylation reaction, we demonstrate that the stabilized redesign can be evolved more efficiently than the wild-type enzyme. After three rounds of directed evolution, we obtained a 6-fold activity increase from the wild-type parent and an 80-fold increase from the stabilized variant. To generate the initial stabilized variant, we identified multiple structural and sequence constraints to preserve catalytic function. We applied these criteria to produce stabilized, catalytically active variants of a second Fe(II)/ $\alpha$ KG enzyme, suggesting that the approach is generalizable to additional members of the Fe(II)/ $\alpha$ KG superfamily. ProteinMPNN is user-friendly and widely accessible, and our results provide a framework for the routine implementation of deep learning-based protein stabilization tools in directed evolution workflows for novel biocatalysts.

## Introduction

Directed evolution is a powerful method to generate enzymes for new chemical transformations.<sup>[1,2]</sup> However, catalytic functional groups often have destabilizing effects on protein structure, and altering active site groups for new reactions can lead to unstable, non-functional proteins.<sup>[3–11]</sup> Initiating a directed evolution campaign from a stabilized variant can be an effective way to overcome this problem.<sup>[12,13]</sup> Typically, a starting point for directed evolution is obtained by screening a library of candidate enzymes for a desired promiscuous activity. If a thermostable homolog with similar catalytic properties is identified, it can then be used as a starting point for evolution of the desired function.<sup>[14,15]</sup> Alternatively, there are a variety of strategies to produce stable variants using directed evolution,<sup>[16,17]</sup> ancestral reconstruction,<sup>[14]</sup> protein recombination,<sup>[18,19]</sup> or computational engineering.<sup>[20–22]</sup> These methods are often time- and resource-intensive, highlighting the need for simple and accessible alternatives.

An important class of enzymes where stabilization would be useful is the non-heme iron(II)  $\alpha$ -ketoglutarate-dependent oxygenase (Fe(II)/ $\alpha$ KG) superfamily.<sup>[23–26]</sup> These enzymes have emerged as a rich source of potential new biocatalysts. Fe(II)/ $\alpha$ KG enzymes can perform remote, asymmetric C(sp<sup>3</sup>)-H oxyfunctionalization reactions on small molecule substrates using a conserved, radical-mediated mechanism. These transformations are synthetically challenging,<sup>[27,28]</sup> and a biocatalytic alternative could allow expedient and sustainable diversification of simple building blocks to a range of complex polyfunctional compounds. The advantages offered by this enzyme family include a high degree of chemical flexibility in the iron-containing active site due to multiple open coordination sites, the utilization of benign molecular oxygen as an oxidant, and use of the inexpensive and readily available co-factor  $\alpha$ KG. However, Fe(II)/ $\alpha$ KGs can be relatively unstable,<sup>[29,30]</sup> which may limit their practical applications in organic synthesis and in industrial process applications.

The recent use of an Fe(II)/ $\alpha$ KG in an industrial-scale drug biosynthesis pathway highlights both the potential advantages and drawbacks of this family for biocatalysis. An engineered Fe(II)/ $\alpha$ KG was used to catalyze an enantioselective C(sp<sup>3</sup>)-H hydroxylation to produce a key intermediate for the anti-cancer drug belzutifan.<sup>[30]</sup> The reaction could be performed on kilogram-scale and bypassed five steps of the pre-existing chemical synthesis route. Notably, this effort required an extensive, large-scale directed evolution campaign. Furthermore, early rounds of screening yielded stabilizing mutations before significant improvements in turnover could be obtained in later rounds. These findings highlight the potential utility of Fe(II)/ $\alpha$ KG for practical,

[\*] B. R. King, K. H. Sumida, J. L. Caruso, J. G. Zalatan  
Department of Chemistry  
University of Washington  
Seattle, Washington 98195, USA  
E-mail: zalatan@uw.edu  
K. H. Sumida, D. Baker  
Institute of Protein Design  
University of Washington  
Seattle, Washington 98195, USA  
D. Baker  
Howard Hughes Medical Institute  
University of Washington  
Seattle, Washington 98195, USA

industrial-scale green chemistry but also the importance of enzyme stability in the evolution of new function.

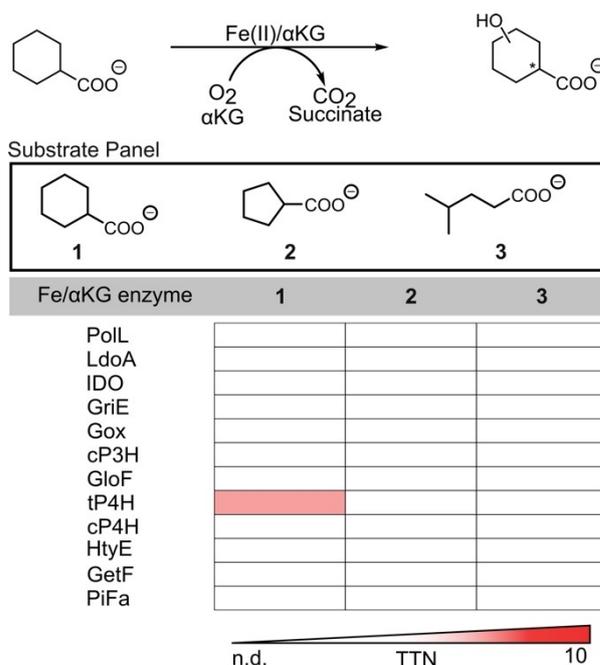
Recent applications of deep learning to protein design have provided new and relatively straightforward methods to stabilize protein scaffolds,<sup>[22,31,32]</sup> and there is broad interest in applying these approaches to directed evolution.<sup>[33]</sup> Here we demonstrate that the deep learning-based tool ProteinMPNN<sup>[31,32]</sup> enables more efficient optimization of a synthetically relevant, non-native C(*sp*<sup>3</sup>)-H hydroxylation reaction in an Fe(II)/ $\alpha$ KG family member. A critical step was identifying appropriate design criteria to prevent modification of residues important for catalytic function, which includes both active site and remote positions. With a stabilized starting point for site-saturation mutagenesis, we observed substantially larger increases in non-native activity compared to the same mutations in the wild-type parent enzyme. This systematic comparison of the wild-type parent and the stabilized redesign provides a critical benchmark for the field to evaluate the effectiveness of these tools. We suggest that this designed stabilization approach should be routinely used in future directed evolution campaigns with the Fe(II)/ $\alpha$ KG superfamily and will likely be effective in a broad range of other enzyme families. There have been many recent reports applying machine learning tools to design variant libraries or pick residues for functional optimization,<sup>[22,33–38]</sup> and stabilized scaffolds can readily be coupled to these approaches for tailored engineering.

## Results and Discussion

### Fe(II)/ $\alpha$ KGs with Promiscuous Activity for C–H Hydroxylation of Free Carboxylate Substrates

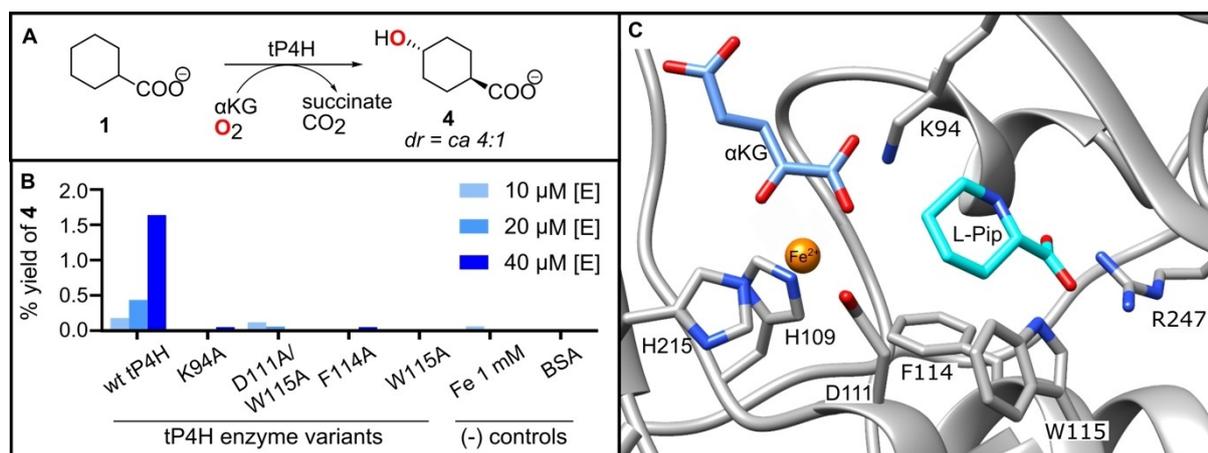
Within the Fe(II)/ $\alpha$ KG enzyme superfamily, free amino acid hydroxylases are attractive candidates for engineering new reactions.<sup>[23–26]</sup> Because Fe(II)/ $\alpha$ KG amino acid hydroxylases already have catalytic machinery to interact with amine and carboxylate functional groups in amino acids, we hypothesized that they might have promiscuous activity for substrates containing only an amine or only a carboxylate. These molecules are important feedstocks for early-stage oxyfunctionalization reactions in multi-step syntheses. Selectivity for remote C(*sp*<sup>3</sup>)-H hydroxylation reactions has been historically difficult to achieve with traditional transition metal catalysis, and a biocatalytic process could offer improved regio- and stereoselectivity.<sup>[27,28]</sup>

We initially screened a panel of 12 Fe(II)/ $\alpha$ KG amino acid hydroxylases for the ability to hydroxylate free carboxylates (Figure 1). The enzymes in this panel were chosen for their ability to hydroxylate free amino acids and their ease of expression, handling, and purification (Table S6). We chose a set of candidate carboxylate substrates (1–3) that are structurally analogous to the native amino acid substrates L-pipecolic acid (L-Pip), L-proline (L-Pro), and L-leucine (L-Leu). We used whole-cell biocatalysis and liquid chromatography-mass spectrometry (LC-MS) to detect products. We confirmed that the native amino acid reaction



**Figure 1.** Initial whole-cell reaction screen data with a panel of Fe(II)/ $\alpha$ KG amino acid hydroxylases and free acid substrate analogues. Each enzyme was screened against all three substrates. A white panel indicates that product was not detectable (n.d.), which is the case for all reactions except tP4H with substrate 1. Reactions were performed in whole cell from 50 mL expression cultures where whole cell volume was 1/20<sup>th</sup> the expression volume. Reactions were carried out in MOPS (pH 7.0, 50 mM) with 20 mM substrate, 60 mM  $\alpha$ KG (as disodium salt), 1 mM ferrous ammonium sulfate, and 1 mM L-ascorbic acid.

products are detectable with all 12 members of the enzyme panel (Table S6). We then screened for promiscuous activity with carboxylates, and observed that one enzyme, tP4H,<sup>[39]</sup> has detectable activity with substrate 1 (Figure 1). The total turnover number (TTN) with this substrate was approximately 5 after incubation for 24 h with 10  $\mu$ M enzyme, about 130-fold lower than the TTN for the corresponding native amino acid substrate. The reaction of tP4H with substrate 1 gives the *trans* product with a *d.r.* of 4:1 (Figure S4). tP4H produces exclusively *trans* product with its native substrate L-Pip,<sup>[39]</sup> suggesting that the free carboxylate substrate 1 and the native substrate are positioned similarly in the enzyme active site with respect to the iron center. To confirm that tP4H and not a contaminating enzyme was responsible for the observed non-native activity, we mutated active site residues that are involved in Fe(II), substrate, or  $\alpha$ KG binding. Because tP4H does not have an experimental structure, we identified these active site residues using an Alphafold2<sup>[40]</sup> model (Figure S9) and comparisons to structures of the highly homologous Fe(II)/ $\alpha$ KG enzyme GriE.<sup>[41,42]</sup> In all cases, active site mutations produced activity decreases for the non-native substrate 1 (Figure 2). We also observed increased product yield with increasing wild-type tP4H concentration (Figure 2). The overall yields remain relatively low, which is typical for promiscuous, non-native reactions.<sup>[1]</sup> Together these results confirm that tP4H is responsible for the non-native reaction.



**Figure 2.** Validation of tP4H activity with free acid **1**. A) Reaction of tP4H with substrate **1** to form *trans*-4-hydroxycyclohexane carboxylic acid **4**. B) Yield of **4** after reaction of **1** with tP4H variants, as well as negative control reaction with Fe(II) and bovine serum albumin (BSA). The Fe 1 mM control was run in the absence of added enzyme. Purified enzyme concentration was varied between 10–40  $\mu\text{M}$  with 20 mM **1**, 40 mM  $\alpha\text{KG}$ , 1 mM ferrous ammonium sulfate, and 1 mM L-ascorbic acid in MES buffer (50 mM, pH 6.8). Reactions were carried out for 24 h at 25 °C and quantified with analytical LC–MS. C) Structural model of the tP4H showing key active site residues. Fe(II),  $\alpha\text{KG}$ , and L-Pip were modeled in Chimera.<sup>[43]</sup>

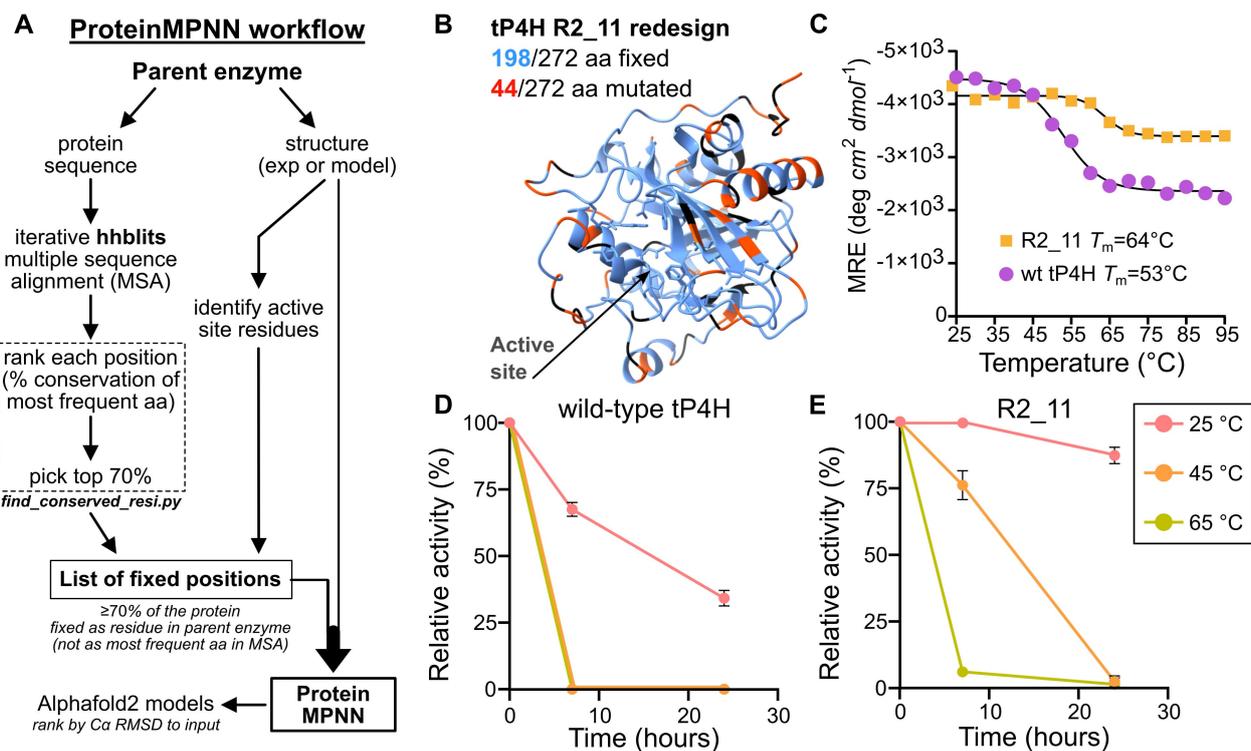
### Stabilization of tP4H with ProteinMPNN

To improve tP4H activity towards substrate **1**, we began a directed evolution campaign but quickly encountered limitations due to poor enzyme stability. First, we found that tP4H variants were difficult to express and purify due to enzyme insolubility. Additionally, we found that the parent wild-type tP4H enzyme loses activity with time (Figure 3). These observations are consistent with prior reports on tP4H behavior.<sup>[39]</sup> It is possible for enzyme stability to improve during directed evolution, whether by selecting more stable variants during each round or by chance. For example, the evolved Fe(II)/ $\alpha\text{KG}$  PSEFE had slight stability improvements compared to wild-type, despite the researchers not selecting for improved stability.<sup>[44]</sup> In another case, the Fe(II)/ $\alpha\text{KG}$  UbP4H was successfully screened for improved stability after initial screening rounds destabilized the enzyme.<sup>[29]</sup> However, the benefits of starting with a highly stable enzyme for directed evolution are well-established.<sup>[12–14]</sup> A stabilized protein scaffold could potentially increase the population of active, properly folded protein or provide access to other mutants that otherwise do not fold.

We used the deep learning-based tool ProteinMPNN<sup>[31,32]</sup> to generate stabilized variants of tP4H through sequence design. We first used ProteinMPNN to redesign the entire tP4H sequence. Unsurprisingly, we found that the predicted sequences eliminated key active site residues, which is likely to disrupt enzymatic activity (Figure S10). This behavior is consistent with the well-established propensity for catalytic active site residues to be destabilizing.<sup>[3–10]</sup> To preserve catalytic function, we fixed the active site residues in all subsequent design efforts. We defined the active site as any residues that contact the amino acid substrate, Fe(II), or the  $\alpha\text{KG}$  cofactor, based on our AlphaFold2 model (Figure S9) and comparisons to the structure of the closely related enzyme GriE bound to L-Leu.<sup>[42]</sup> Because other residues

throughout the protein could also be important, we also tested four additional strategies using either sequence conservation or distance metrics. To identify important conserved residues, we constructed a multiple sequence alignment (MSA) and selected tP4H residues conserved in at least 35 %, 70 %, or 95 % of sequences (Methods and Table S4). Alternatively, we fixed any residues with side chains within a 10 Å sphere from the substrate binding pocket. Using these five starting points (fix active site only, active site + 35 %/70 %/95 % conservation, active site + 10 Å sphere), we generated 48 ProteinMPNN sequences per method, generated AlphaFold2 models, and selected 4 each (20 total) for activity screens. Selection was based on calculated top-ranked C $\alpha$ -RMSD values matched to the input tP4H structure. We obtained only one variant that had any detectable activity, with catalytic efficiency ~35-fold lower than wild-type tP4H (Figure S11). This variant was designed from the sequence where >35 % conserved residues are fixed, which constrains more residues than the >70 % or >95 % cutoffs. This result suggests that even weakly conserved residues may need to be fixed to maintain activity. Further analysis of the variant with detectable activity revealed an approximately 2-fold increase in the  $K_M$  for the  $\alpha\text{KG}$  cofactor and an approximately 3-fold decrease in  $k_{\text{cat}}$  (Figure S11). We identified two residues in proximity to  $\alpha\text{KG}$ , L228 and V230, that were mutated in the redesigned sequence. These sequence changes may have contributed to improved stability at the expense of cofactor binding and positioning, leading to the decrease in activity. Notably, L228 and V230 were fixed in the designs generated from fixed active site + 10 Å sphere, but none of these designs had detectable activity. Taken together, these findings suggest that additional criteria will be needed to identify critical functional residues that should be fixed prior to sequence redesign.

To generate stabilized variants that maintain catalytic activity, we performed another set of ProteinMPNN se-



**Figure 3.** Stability and activity of wild-type tP4H and ProteinMPNN design R2\_11. A) Flowchart of the ProteinMPNN workflow that successfully produced stabilized variants that retained catalytic activity. See Supporting Information for detailed guidelines for each computational step. B) tP4H structure (AlphaFold2 model, Figure S9) color-coded to show sites fixed in the design process (blue, Supplementary spreadsheet—ProteinMPNN sequences\_metrics) and sites mutated in the ProteinMPNN R2\_11 redesign (orange-red). Sites colored black were neither fixed nor redesigned in the R2\_11 variant. Side chains for first shell active site residues (Table S4) are shown in blue. C) Temperature-dependent CD spectroscopy of wild-type tP4H and R2\_11.  $T_m$  values were calculated using the Boltzmann sigmoid function in GraphPad Prism. D) Activity-stability analysis of wild-type tP4H. E) Activity-stability analysis of R2\_11. For (D) and (E), relative activity was determined using PBP assay described in Supporting Information. Values are mean  $\pm$  SD for three replicates.

quence redesigns with three new strategies to fix important residues. In each case, we fixed the active site as defined above plus residues L228 and V230. For the first approach, we fixed all residues at tP4H positions conserved in 35 % of the MSA. We chose this cutoff because it was the only one from our initial set that produced a stabilized variant with any detectable activity, and we expected that fixing L228 or V230 could further improve these designs. For the second and third approaches, we identified highly conserved positions regardless of whether the wild-type tP4H residue is the most highly conserved amino acid. These strategies were based on previous work suggesting that more stringent constraints are necessary to maintain activity in ProteinMPNN redesign.<sup>[32]</sup> Every tP4H amino acid position was ranked based on the % conservation of the most frequent amino acid present in the MSA, and the top 50 % or 70 % were fixed. These positions were fixed as the wild-type tP4H residue, even if they were different from the consensus most frequent amino acid in the MSA. Together with fixed active site residues, these criteria resulted in 148/272 (54 %) or 198/272 (73 %) fixed residues across the entire 272 amino acid protein. Using these three strategies, we selected 32 designs each of 48 generated for a total of 96 sequences (Supplementary spreadsheet—ProteinMPNN sequences\_metrics). Of these designs, 69 expressed detectable quanti-

ties of protein by SDS-PAGE and 11 had detectable activity above background for the native substrate. For the active enzymes we proceeded to measure thermostability and kinetic parameters for the native L-Pip substrate and the promiscuous carboxylate substrate **1**. The variant with the highest  $k_{cat}$  for L-Pip was R2\_11 (Table S7), with a  $k_{cat}$  of 0.10 s<sup>-1</sup> compared to 0.14 s<sup>-1</sup> for wild-type. R2\_11 was designed from the method where the top 70 % ranked conserved residues were fixed, and had 44 designed mutations compared to the wild-type sequence (Figure 3). R2\_11 has modestly slower (ca. 3-fold) non-native carboxylate hydroxylase activity compared to wild-type tP4H and exhibits an 11 °C increase in thermal melting temperature ( $T_m$ ) as measured by temperature-dependent circular dichroism (CD) spectroscopy (Figure 3). When activity is measured as a function of time, R2\_11 maintains activity over a timescale of days, which is a substantial improvement compared to wild-type tP4H (Figure 3). The modest decrease in promiscuous activity is unsurprising because ProteinMPNN does not consider catalytic activity, and there is no expectation that global protein stabilization would either maintain or optimize a non-native reaction.

### Stabilization of GriE with ProteinMPNN

After successfully identifying sequence constraints for ProteinMPNN-mediated stabilization of tP4H while maintaining catalytic function, we evaluated whether the same approach would be effective with a second Fe(II)/ $\alpha$ KG amino acid hydroxylase, GriE. This enzyme could benefit from stabilization because, although it expresses well and is soluble, it loses activity at room temperature over 24 h (Figure S12). As with tP4H, we fixed the top 70 % ranked conserved residues along with catalytic residues identified from the GriE crystal structure (Supplementary spreadsheet –ProteinMPNN sequences\_metrics).<sup>[42]</sup> We generated 48 redesigned sequences, selected the top 32, and found that 29 were expressed as soluble enzyme, and 27 showed activity with the GriE native substrate L-Leu. The top design based on stability and kinetic parameters, GM\_A9, showed a similar catalytic efficiency ( $k_{\text{cat}}/K_M$ ) and an approximately 4-fold decrease in  $k_{\text{cat}}$  compared to wild-type GriE (Figure S12&S13, Table S8). One design, GM\_A11, had a 2-fold faster initial rate with L-Leu compared to GM\_A9 but this design was unstable by temperature-dependent CD and thus was not chosen for further analysis. A decrease in  $k_{\text{cat}}$  is not surprising given that increased stability could reduce conformational flexibility and negatively impact catalytic function.<sup>[4,6,7]</sup>

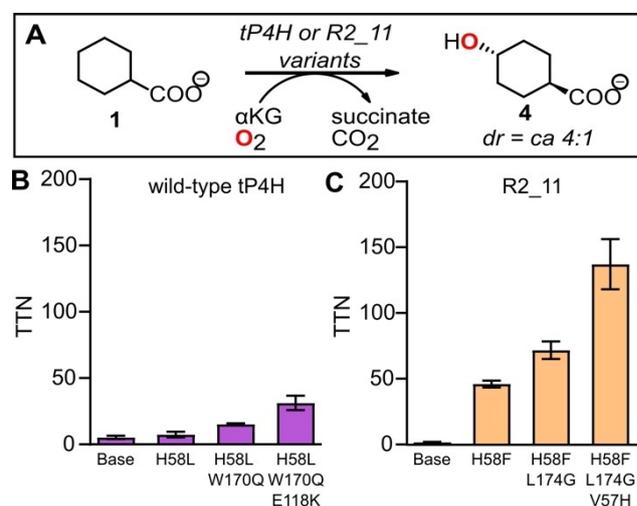
We next screened the stabilized GriE redesign GM\_A9 for substrate promiscuity. Previously, wild-type GriE has been shown to accept substrates with increased substrate chain lengths but has weaker activity towards substrates with substitution at C3.<sup>[45]</sup> We chose two previously identified non-native substrates to test: L-norleucine (L-Nle) and L-allo-isoleucine (L-allo-Ile). L-Nle was chosen as a representative substrate with increased chain length compared to L-Leu. L-allo-Ile was chosen because it has a methyl group substitution at C3. L-isoleucine also has a C3 methyl group, but it is not detectably hydroxylated by wild-type GriE<sup>[45]</sup> and was therefore not included in this analysis. We observed detectable activity with L-Nle but not for L-allo-Ile (Table S9). Similar to wild-type GriE, the GM\_A9 variant maintained a preference for the extended chain L-Nle substrate over the C3-substituted L-allo-Ile substrate. The GM\_A9 reactions with L-Leu and L-Nle were 11- and 4-fold slower than wild-type GriE reactions, respectively (Table S9). These results suggest that our ProteinMPNN protocol can be readily applied to other Fe(II)/ $\alpha$ KG enzymes to stabilize proteins while maintaining synthetically relevant catalytic function that can be a foothold for further optimization by directed evolution.

### Directed Evolution of Wild-Type tP4H for Carboxylate C–H Hydroxylation Activity

We next sought to improve the non-native carboxylate C( $sp^3$ )–H hydroxylase activity through directed evolution. We prioritized tP4H because carboxylate hydroxylase activity was detectable in both the wild-type and ProteinMPNN-stabilized variant, which allows for direct com-

parisons. We conducted three rounds of directed evolution for both enzymes by varying first- and second-shell substrate binding residues identified in the active site from our AlphaFold2 structural model. We defined the first shell as any residues that contact the amino acid substrate, based on comparisons to the structure of ligand-bound GriE.<sup>[42]</sup> We defined the second shell as any residues that make contacts with first shell residues. We used the 22c-trick method for single site-saturation mutagenesis at each target position, and we screened 70 colonies for each position to ensure >95 % library coverage (Table S1).<sup>[46]</sup>

We first performed directed evolution with wild-type tP4H. For the first screening round, we chose three tP4H active site residues based on their potential role in substrate specificity: H58, F114, and L174. Based on our tP4H structural model (Figure S9B), H58 likely contacts the amine of native amino acid substrates and is presumably not needed or detrimental for carboxylate substrates that lack an amine. F114 likely provides a substrate hydrophobic contact, and L174 is part of a loop that could affect substrate binding. We screened whole cell biocatalysis reactions in 96 well plates for improved TTN and 80 % *trans* selectivity in reactions with substrate **1**. Based on production of hydroxylated product **4**, the top 5 % of mutants were chosen for validation with purified enzymes. We obtained several variants with modest activity improvements, and the best performer was mutant H58L with a TTN of 7 (Figure 4A). In a second round starting from H58L, we rescreened mutants at F114 and L174 and screened an additional 14 first and second shell residues (Table S1). We identified the improved variant H58 L/W170Q with a TTN of 15. In a third round starting from H58 L/W170Q, we screened 9 residues



**Figure 4.** A) C–H hydroxylation of substrate **1** with tP4H, R2\_11 and associated variants. B) Directed evolution of wild-type tP4H. C) Directed evolution of stabilized variant R2\_11. Reactions were carried out for 24 h at 25 °C using purified enzyme (10–20  $\mu$ M) in MES buffer (50 mM, pH 6.8), with 20 mM cyclohexane carboxylic acid **1**, 40 mM  $\alpha$ KG, 1 mM ferrous ammonium sulfate, and 1 mM ascorbic acid. Concentration of **4** in quenched reaction samples was quantified by analytical LC–MS analysis. For (B) and (C), values are mean  $\pm$  SD for three replicates.

that showed activity increases in previous rounds and identified the improved variant H58 L/W170Q/E118 K with a TTN of 31 (Figure 4A & Table S1). Overall, after three rounds of directed evolution for improved carboxylate hydroxylase activity with wild-type tP4H we obtained a 6-fold improvement in TTN and maintained >80 % selectivity for the *trans* reaction product (Figure S4).

#### Directed Evolution of ProteinMPNN-Stabilized tP4H for Carboxylate C–H Hydroxylation Activity

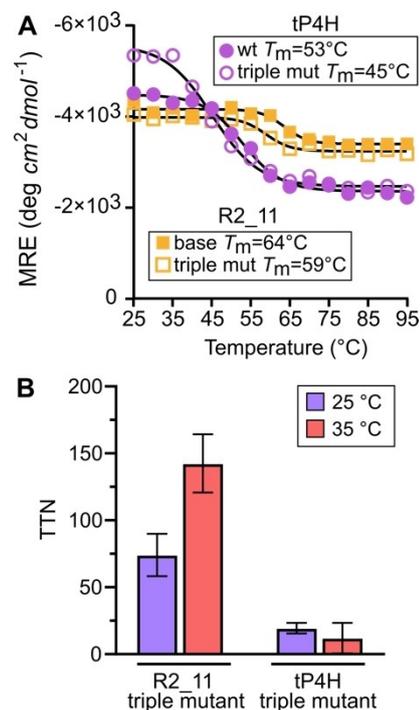
To optimize carboxylate hydroxylase activity in ProteinMPNN-stabilized tP4H, we conducted directed evolution using a similar strategy to our approach with wild-type tP4H, with minor modifications. In the first round of site saturation mutagenesis, we started with a larger pool of 19 first- and second-shell residues including the three sites from prior round one (H58, F114, and L174) and the 16 additional sites from prior round 2 (Table S1). As before, we screened whole cell biocatalysis reactions for improved TTN and >80 % *trans* selectivity with carboxylic acid substrate **1**. The top hit was H58F, which is the same position but a different mutant than the previous round one winner, H58L. R2\_11\_H58F displayed a 27-fold increase in TTN relative to the parent R2\_11 (Figure 4B). This effect is substantially bigger than the <2-fold improvement obtained with H58L relative to wild-type tP4H. Notably, the 27-fold increase in a single round from stabilized R2\_11 was already larger than the total 6-fold improvement from three rounds of directed evolution from wild-type tP4H. Given the strong performance of the H58F mutant, we also evaluated its effect in the wild-type tP4H background and observed a small, <2-fold increase in TTN, similar to the effect of H58 L on wild-type tP4H (Figure S14). Thus, the strong, 27-fold improvement with the H58F mutant depends on the context of the stabilized R2\_11 backbone. Context-dependent activity increases have previously been observed in stabilized variants, supporting the general idea that stability can promote evolvability.<sup>[12]</sup>

In a second round of screening from R2\_11\_H58F, we selected the 18 residues that were screened in prior round two from wild-type tP4H (Table S1). We retained this large pool of residues to ensure a direct comparison to the tP4H directed evolution workflow. This round identified improved variant L174G. R2\_11\_H58F/L174G has a TTN of 72. This TTN is a 1.6-fold improvement from parent R2\_11\_H58F and outstrips any variant obtained from the wild-type tP4H backbone (Figure 4B). In a third round of screening from R2\_11\_H58F/L174G, we selected 9 residues that were screened in prior round 3 from wild-type tP4H (Table S1). This round identified the improved variant V57H with a TTN of 138, a 1.7-fold improvement from the previous round.

Overall, the ProteinMPNN-stabilized tP4H directed evolution campaign produced an 80-fold improvement in TTN from the base R2\_11 redesign, compared to a modest 6-fold improvement in the wild-type tP4H evolutionary trajectory. Although the R2\_11 parent starts about 3-fold

slower than wild-type tP4H, the much larger improvement over three rounds of directed evolution produced an R2\_11 triple mutant with a 4.5-fold higher TTN than the triple mutant obtained from wild-type tP4H (Figure 4).

In addition to a more efficient directed evolution trajectory, the R2\_11 triple mutant maintains high stability relative to the triple mutant derived from wild-type tP4H (Figure 5A, Figure S15), with only a modest decrease in thermal stability compared to the R2\_11 parent. Higher stability allows reactions to be run more efficiently, both at higher temperatures and for less time. For example, after 6 h at 35 °C, the R2\_11 triple mutant reaches a mean TTN of 142 for the non-native reaction with carboxylate **1** to form product **4** with 4:1 selectivity for the *trans* reaction product (Figure 5B). The TTN after 6 h at 35 °C is comparable to the TTN after 24 h at 25 °C. In contrast, the tP4H triple mutant shows a slight decrease in TTN at 35 °C, likely due to enzyme instability at higher temperatures (Figure 5B). The stability profile of the R2\_11 triple mutant suggests that this enzyme will be more robust towards further engineering compared to the tP4H triple mutant. Future engineering efforts with the R2\_11 mutant could include improvements to key reaction metrics like turnover, selectivity, and increased substrate scope.



**Figure 5.** A) Temperature-dependent CD of the R2\_11 and tP4H parent enzymes and triple mutants.  $T_m$  values were calculated using the Boltzmann sigmoid function in GraphPad Prism. B) TTN for formation of **4** (*d.r.* 4:1, Figure S4) with the R2\_11 and tP4H triple mutants at two different temperatures. Reactions were carried out for 6 h at 25 °C and 35 °C using purified enzyme (15  $\mu$ M) in MES buffer (50 mM, pH 6.8), with 20 mM cyclohexane carboxylic acid **1**, 40 mM  $\alpha$ KG, 1 mM ferrous ammonium sulfate, and 1 mM ascorbic acid. Values are mean  $\pm$  SD for three replicates.

## Conclusion

Directed evolution is a powerful tool to engineer enzymes for new-to-nature reactions. However, many enzyme starting points for evolution may lack the stability required to reach user-defined optimum fitness after multiple rounds of mutagenesis. Here we show that the deep learning-based tool ProteinMPNN can be used to stabilize the Fe(II)/ $\alpha$ KG enzyme superfamily members tP4H and GriE with straightforward sequence constraints to maintain catalytic activity. Consistent with previous results using ProteinMPNN,<sup>[32]</sup> the top tP4H design was identified by using the most conservative of our chosen methods for fixing residues during sequence redesign. Applying the same method to the related enzyme GriE readily produced stabilized variants with catalytic activity.

Wild-type and redesigned tP4H both exhibit novel reactivity towards remote C(sp<sup>3</sup>)-H hydroxylation of a free carboxylic acid substrate. We directly compared evolutionary trajectories of wild-type tP4H with the stabilized variant R2\_11 and demonstrated superior performance of the stabilized redesign variant. Future work will determine if this design method is generalizable to optimize directed evolution for other enzymes and enzyme families. Further improvements to deep learning models, or an improved understanding of the underlying mechanisms for enzyme stabilization, could be necessary for broad generalizability. Additional systematic comparisons will also be necessary to evaluate ProteinMPNN relative to other methods for enzyme stabilization. For example, PROSS<sup>[20,47]</sup> uses physics-based energy calculations to generate stabilized sequences, and MutCompute<sup>[22,48]</sup> uses a deep learning approach to identify individual point mutations. Both approaches are distinct from the complete sequence redesign produced from ProteinMPNN. Directly comparing the ability and efficiency for each approach to generate stable variants for directed evolution could identify tradeoffs and potential advantages for each method. User-friendly computational tools are rapidly emerging, and our work suggests that these tools should be routinely incorporated into enzyme engineering workflows to efficiently optimize catalytic fitness for new biocatalysts.<sup>[33]</sup>

## Supporting Information

Supporting Information includes materials, experimental and analytical methods, compound characterization data (Figure S4–S7), and enzyme characterization data (Figure S3, Figure S8, Figure S11–S15 and Tables S6–S9) (PDF). A supplementary spreadsheet includes accession numbers for all Fe(II)/ $\alpha$ KG enzymes used in this work, full nucleotide and amino acid sequences for all reported wild-type and enzyme variants, oligonucleotide sequences for cloning and for all ProteinMPNN designs, ProteinMPNN design screen criteria results (XLSX). The supplementary file `find_conserved_resi.txt` contains the python script to parse multiple sequence alignments. The authors have cited additional references within the Supporting Information.<sup>[49–57]</sup>

## Acknowledgements

We thank Dr. Wolfgang Hüttel at the University of Freiburg and Hans Renata at Rice University for the donation of the wild-type tP4H and GriE expression vectors, respectively, and for their advice on tP4H and GriE reactions in various formats. We also thank Jonathan Zhang and Susanna Vazquez Torres for their early contributions to Fe(II)/ $\alpha$ KG reaction screening, Jue Wang for assistance with the python script to parse sequence alignments, and Dr. Martin Sadilek in University of Washington Mass Spectrometry Facility for his continued support and helpful advice in analytical method development. This work was supported by U.S. National Institutes of Health grants T32GM008268 (B.R.K., J.L.C.) and R35GM124773 (J.G.Z.), and by the Open Philanthropy Project Improving Protein Design Fund (K.H.S., D.B.).

## Conflict of Interest

The authors declare no conflict of interest.

## Data Availability Statement

The data that support the findings of this study are available in the supplementary material of this article.

**Keywords:** Biocatalysis · Directed evolution · Hydroxylation · non-heme iron(II)  $\alpha$ -ketoglutarate-dependent oxygenase · Protein design

- [1] H. Renata, Z. J. Wang, F. H. Arnold, *Angew. Chem. Int. Ed.* **2015**, *54*, 3351–3367.
- [2] C. Zeymer, D. Hilvert, *Annu. Rev. Biochem.* **2018**, *87*, 1–27.
- [3] E. M. Meiering, L. Serrano, A. R. Fersht, *J. Mol. Biol.* **1992**, *225*, 585–589.
- [4] B. K. Shoichet, W. A. Baase, R. Kuroki, B. W. Matthews, *Proc. Nat. Acad. Sci.* **1995**, *92*, 452–456.
- [5] L. Giver, A. Gershenson, P.-O. Freskgard, F. H. Arnold, *Proc. Nat. Acad. Sci.* **1998**, *95*, 12809–12813.
- [6] B. M. Beadle, B. K. Shoichet, *J. Mol. Biol.* **2002**, *321*, 285–296.
- [7] R. A. Nagatani, A. Gonzalez, B. K. Shoichet, L. S. Brinen, P. C. Babbitt, *Biochemistry* **2007**, *46*, 6688–6695.
- [8] N. Tokuriki, F. Stricher, L. Serrano, D. S. Tawfik, *PLoS Comput. Biol.* **2008**, *4*, e1000002.
- [9] N. Tokuriki, C. J. Jackson, L. Afriat-Jurnou, K. T. Wyganowski, R. Tang, D. S. Tawfik, *Nat. Commun.* **2012**, *3*, 1257.
- [10] M. Goldsmith, D. S. Tawfik, *Curr. Opin. Struct. Biol.* **2017**, *47*, 140–150.
- [11] K. Tsuboyama, J. Dauparas, J. Chen, E. Laine, Y. M. Behbahani, J. J. Weinstein, N. M. Mangan, S. Ovchinnikov, G. J. Rocklin, *Nature* **2023**, *620*, 434–444.
- [12] J. D. Bloom, S. T. Labthavikul, C. R. Otey, F. H. Arnold, *Proc. Nat. Acad. Sci.* **2006**, *103*, 5869–5874.
- [13] Y. Gumulya, J.-M. Baek, S.-J. Wun, R. E. S. Thomson, K. L. Harris, D. J. B. Hunter, J. B. Y. H. Behrendorff, J. Kulig, S. Zheng, X. Wu, B. Wu, J. E. Stok, J. J. D. Voss, G. Schenk, U. Jurva, S. Andersson, E. M. Isin, M. Bodén, L. Guddat, E. M. J. Gillam, *Nat. Catal.* **2018**, *1*, 878–888.

- [14] D. L. Trudeau, D. S. Tawfik, *Curr. Opin. Biotechnol.* **2019**, *60*, 46–52.
- [15] W. Besenmatter, P. Kast, D. Hilvert, *Proteins Struct. Funct. Bioinf.* **2007**, *66*, 500–506.
- [16] R. D. Socha, N. Tokuriki, *FEBS J.* **2013**, *280*, 5582–5595.
- [17] S. D. Stimple, M. D. Smith, P. M. Tessier, *AIChE J.* **2020**, *66*, DOI 10.1002/aic.16814.
- [18] Y. Li, D. A. Drummond, A. M. Sawayama, C. D. Snow, J. D. Bloom, F. H. Arnold, *Nat. Biotechnol.* **2007**, *25*, 1051–1056.
- [19] P. Heinzelman, C. D. Snow, I. Wu, C. Nguyen, A. Villalobos, S. Govindarajan, J. Minshull, F. H. Arnold, *Proc. Nat. Acad. Sci.* **2009**, *106*, 5610–5615.
- [20] A. Goldenzweig, M. Goldsmith, S. E. Hill, O. Gertman, P. Laurino, Y. Ashani, O. Dym, T. Unger, S. Albeck, J. Prilusky, R. L. Lieberman, A. Aharoni, I. Silman, J. L. Sussman, D. S. Tawfik, S. J. Fleishman, *Mol. Cell* **2016**, *63*, 337–346.
- [21] M. Musil, H. Konegger, J. Hon, D. Bednar, J. Damborsky, *ACS Catal.* **2019**, *9*, 1033–1054.
- [22] H. Lu, D. J. Diaz, N. J. Czarnecki, C. Zhu, W. Kim, R. Shroff, D. J. Acosta, B. R. Alexander, H. O. Cole, Y. Zhang, N. A. Lynd, A. D. Ellington, H. S. Alper, *Nature* **2021**, *604*, 662–667.
- [23] M. S. Islam, T. M. Leissing, R. Chowdhury, R. J. Hopkinson, C. J. Schofield, *Annu. Rev. Biochem.* **2018**, *87*, 585–620.
- [24] C. O. Herr, R. P. Hausinger, *Trends Biochem. Sci.* **2018**, *43*, 517–532.
- [25] A. Papadopoulou, F. Meyer, R. M. Buller, *Biochemistry* **2023**, *62*, 229–240.
- [26] C. R. Zwick, H. Renata, *Nat. Prod. Rep.* **2020**, *37*, 1065–1079.
- [27] E. Roduner, W. Kaim, B. Sarkar, V. B. Urlacher, J. Pleiss, R. Gläser, W. Einicke, G. A. Sprenger, U. Beifuß, E. Klemm, C. Liebner, H. Hieronymus, S. Hsu, B. Plietker, S. Laschat, *ChemCatChem* **2013**, *5*, 82–112.
- [28] C. He, W. G. Whitehurst, M. J. Gaunt, *Chem* **2019**, *5*, 1031–1058.
- [29] C. Liu, J. Zhao, J. Liu, X. Guo, D. Rao, H. Liu, P. Zheng, J. Sun, Y. Ma, *Appl. Microbiol. Biotechnol.* **2019**, *103*, 265–277.
- [30] W. L. Cheung-Lee, J. N. Kolev, J. A. McIntosh, A. A. Gil, W. Pan, L. Xiao, J. E. Velásquez, R. Gangam, M. S. Winston, S. Li, K. Abe, E. Alwedi, Z. E. X. Dance, H. Fan, K. Hiraga, J. Kim, B. Kosjek, D. N. Le, N. S. Marzijarani, K. Mattern, J. P. McMullen, K. Narsimhan, A. Vikram, W. Wang, J. Yan, R. Yang, V. Zhang, W. Zhong, D. A. DiRocco, W. J. Morris, G. S. Murphy, K. M. Maloney, *Angew. Chem. Int. Ed.* **2024**, *63*, e202316133.
- [31] J. Dauparas, I. Anishchenko, N. Bennett, H. Bai, R. J. Ragotte, L. F. Milles, B. I. M. Wicky, A. Courbet, R. J. de Haas, N. Bethel, P. J. Y. Leung, T. F. Huddy, S. Pellock, D. Tischer, F. Chan, B. Koepnick, H. Nguyen, A. Kang, B. Sankaran, A. K. Bera, N. P. King, D. Baker, *Science* **2022**, *378*, 49–56.
- [32] K. H. Sumida, R. Núñez-Franco, I. Kalvet, S. J. Pellock, B. I. M. Wicky, L. F. Milles, J. Dauparas, J. Wang, Y. Kipnis, N. Jameson, A. Kang, J. D. L. Cruz, B. Sankaran, A. K. Bera, G. Jiménez-Osés, D. Baker, *J. Am. Chem. Soc.* **2024**, *146*, 2054–2061.
- [33] J. Yang, F.-Z. Li, F. H. Arnold, *ACS Cent. Sci.* **2024**, *10*, 226–241.
- [34] B. J. Wittmann, Y. Yue, F. H. Arnold, *Cell Syst.* **2021**, *12*, 1026–1045.e7.
- [35] S. d'Oelsnitz, D. J. Diaz, W. Kim, D. J. Acosta, T. L. Dangerfield, M. W. Schechter, M. B. Minus, J. R. Howard, H. Do, J. M. Loy, H. S. Alper, Y. J. Zhang, A. D. Ellington, *Nat. Commun.* **2024**, *15*, 2084.
- [36] K. Ding, M. Chin, Y. Zhao, W. Huang, B. K. Mai, H. Wang, P. Liu, Y. Yang, Y. Luo, *Nat. Commun.* **2024**, *15*, 6392.
- [37] J. Yang, R. G. Lal, J. C. Bowden, R. Astudillo, M. A. Hameedi, S. Kaur, M. Hill, Y. Yue, F. H. Arnold, *bioRxiv* **2024**, 2024.07.27.605457.
- [38] N. Thomas, D. Belanger, C. Xu, H. Lee, K. Hirano, K. Iwai, V. Polic, K. D. Nyberg, K. G. Hoff, L. Frenz, C. A. Emrich, J. W. Kim, M. Chavarha, A. Ramanan, J. J. Agresti, L. J. Colwell, *bioRxiv* **2024**, 2024.03.21.585615.
- [39] C. Klein, W. Hüttel, *Adv. Synth. Catal.* **2011**, *353*, 1375–1383.
- [40] M. Mirdita, K. Schütze, Y. Moriwaki, L. Heo, S. Ovchinnikov, M. Steinegger, *Nat. Methods* **2022**, *19*, 679–682.
- [41] X. Hu, X. Huang, J. Liu, P. Zheng, W. Gong, L. Yang, *Acta Crystallogr. Sect. D* **2023**, *79*, 318–325.
- [42] P. Lukat, Y. Katsuyama, S. Wenzel, T. Binz, C. König, W. Blankenfeldt, M. Brönstrup, R. Müller, *Chem. Sci.* **2017**, *8*, 7521–7527.
- [43] E. F. Pettersen, T. D. Goddard, C. C. Huang, G. S. Couch, D. M. Greenblatt, E. C. Meng, T. E. Ferrin, *J. Comput. Chem.* **2004**, *25*, 1605–1612.
- [44] N. W. Goldberg, A. M. Knight, R. K. Zhang, F. H. Arnold, *J. Am. Chem. Soc.* **2019**, *141*, 19585–19588.
- [45] C. R. Zwick, H. Renata, *J. Am. Chem. Soc.* **2018**, *140*, 1165–1169.
- [46] S. Kille, C. G. Acevedo-Rocha, L. P. Parra, Z.-G. Zhang, D. J. Opperman, M. T. Reetz, J. P. Acevedo, *ACS Synth. Biol.* **2013**, *2*, 83–92.
- [47] Y. Peleg, R. Vincentelli, B. M. Collins, K.-E. Chen, E. K. Livingstone, S. Weeratunga, N. Leneva, Q. Guo, K. Remans, K. Perez, G. E. K. Bjerga, Ø. Larsen, O. Vaněk, O. Skořepa, S. Jacquemin, A. Poterszman, S. Kjær, E. Christodoulou, S. Albeck, O. Dym, E. Ainbinder, T. Unger, A. Schuetz, S. Matthes, M. Bader, A. de Marco, P. Storici, M. S. Semrau, P. Stolt-Bergner, C. Aigner, S. Suppmann, A. Goldenzweig, S. J. Fleishman, *J. Mol. Biol.* **2021**, *433*, 166964.
- [48] R. Shroff, A. W. Cole, D. J. Diaz, B. R. Morrow, I. Donnell, A. Annapareddy, J. Gollihar, A. D. Ellington, R. Thyer, *ACS Synth. Biol.* **2020**, *9*, 2927–2935.
- [49] M. P. Okoh, J. L. Hunter, J. E. T. Corrie, M. R. Webb, *Biochemistry* **2006**, *45*, 14764–14771.
- [50] L. Luo, M. B. Pappalardi, P. J. Tummino, R. A. Copeland, M. E. Fraser, P. K. Grzyska, R. P. Hausinger, *Anal. Biochem.* **2006**, *353*, 69–74.
- [51] J. C. Nolte, M. Schürmann, C.-L. Schepers, E. Vogel, J. H. Wübbeler, A. Steinbüchel, *Appl. Environ. Microbiol.* **2014**, *80*, 166–176.
- [52] J. Mattay, S. Houwaart, W. Hüttel, *Appl. Environ. Microbiol.* **2018**, *84*, e02370–17.
- [53] M. Remmert, A. Biegert, A. Hauser, J. Söding, *Nat. Methods* **2012**, *9*, 173–175.
- [54] J. Jumper, R. Evans, A. Pritzel, T. Green, M. Figurnov, O. Ronneberger, K. Tunyasuvunakool, R. Bates, A. Židek, A. Potapenko, A. Bridgland, C. Meyer, S. A. A. Kohl, A. J. Ballard, A. Cowie, B. Romera-Paredes, S. Nikolov, R. Jain, J. Adler, T. Back, S. Petersen, D. Reiman, E. Clancy, M. Zielinski, M. Steinegger, M. Pacholska, T. Berghammer, S. Bodenstein, D. Silver, O. Vinyals, A. W. Senior, K. Kavukcuoglu, P. Kohli, D. Hassabis, *Nature* **2021**, *596*, 583–589.
- [55] T. Wang, X. Jin, X. Lu, X. Min, S. Ge, S. Li, *Front. Genet.* **2024**, *14*, 1347667.
- [56] M. Davidson, M. McNamee, R. Fan, Y. Guo, W.-C. Chang, *J. Am. Chem. Soc.* **2019**, *141*, 3419–3423.
- [57] J. Mattay, W. Hüttel, *ChemBioChem* **2017**, *18*, 1523–1528.

Manuscript received: August 2, 2024

Accepted manuscript online: October 12, 2024

Version of record online: November 11, 2024