

1 **Evolution-guided diffusion generative model enables**  
2 **large-step exploration of functional protein sequence**  
3 **space from single sequences**

4 Xing Zhang<sup>1§</sup>, Jinle Tang<sup>1§</sup>, Tingkai Zhang<sup>1,2§</sup>, Zhihang Chen<sup>3,4</sup>, Zhe Zhang<sup>1</sup>, Jian  
5 Zhan<sup>1,5,6\*</sup>, and Yaoqi Zhou<sup>1\*</sup>

6

7 <sup>1</sup>Institute of Systems and Physical Biology, Shenzhen Bay Laboratory, Shenzhen,  
8 518132, China

9 <sup>2</sup>School of Medicine, Southern University of Science and Technology, Shenzhen,  
10 518055, China

11 <sup>3</sup>Shenzhen Medical Academy of Research and Translation (SMART), Shenzhen, 518107,  
12 China

13 <sup>4</sup>Tsinghua University, Beijing 100084, China

14 <sup>5</sup>Ribopeutic (Shenzhen) Co., Ltd., Shenzhen, 518000, China

15 <sup>6</sup>Ribopeutic Inc., Hangzhou, 310018, China

16 <sup>§</sup>co-first authors. These authors are contributed to the manuscript equally.

17 <sup>\*</sup>Corresponding authors: Yaoqi Zhou, +86-(755) 2684 6275, [zhouyq@szbl.ac.cn](mailto:zhouyq@szbl.ac.cn); Jian

18 Zhan +86-(755) 2684 6275, [zhanjian@szbl.ac.cn](mailto:zhanjian@szbl.ac.cn)

19

## 20 **Abstract**

21 Protein evolution in nature and in the laboratory proceeds through incremental, largely  
22 undirected mutational steps, restricting exploration to local regions of sequence space and  
23 limiting access to remote yet potentially functional proteins. We present EvoGUD, a  
24 single-sequence-conditioned diffusion framework for large-step exploration of protein  
25 sequence space under learned evolutionary constraints. EvoGUD-generated sequences  
26 preserve natural-like co-evolutionary structure in representation space despite large  
27 sequence divergence. When assembled as virtual multiple sequence alignments, these  
28 sequences substantially improve AlphaFold3 single-sequence inference, restoring much  
29 of the backbone accuracy and atomic-level side-chain realism for recent deposited protein  
30 monomers as well as protein-nucleic-acid complexes, without evolutionary database  
31 search. Moreover, EvoGUD enables functional discovery in remote sequence space,  
32 yielding active variants of the adenine base-editing enzyme TadA in targeted validation  
33 experiments (80% success rate) and large numbers of functional variants of the  
34 intrinsically disordered antitoxin CcdA in high-throughput selection assays (19% success  
35 rate). Together, these results establish EvoGUD as a single-sequence, evolution-aware  
36 generative framework for large-step navigation of protein sequence space, with direct  
37 implications for structure modeling and functional protein discovery in previously  
38 unexplored sequence space.

## 39 **Introduction**

40 Protein evolution has generated an extraordinary diversity of molecular structures and  
41 biochemical functions, yet the mechanisms by which evolution explores protein sequence  
42 space are intrinsically constrained. Natural evolution proceeds through incremental,  
43 largely undirected mutations accumulated over long timescales, while laboratory directed  
44 evolution accelerates selection but still relies on random, local mutational steps<sup>1,2</sup>.  
45 Consequently, both processes tend to explore narrow neighborhoods around existing  
46 solutions, leading to uneven sampling of protein fitness landscapes and limited access to  
47 distant yet potentially functional regions of sequence space<sup>3</sup>. A central challenge in  
48 protein science is therefore to enable large-step exploration of protein sequence space—  
49 jumping to remote regions—while preserving the functional and structural constraints  
50 characteristic of a target protein family.

51 Computational approaches to this challenge broadly fall into structure-based and  
52 sequence-based strategies. Structure-centric methods have achieved remarkable success  
53 in *de novo* protein design, but their reliance on explicit backbone templates or structural  
54 hypotheses limits their applicability when reliable structures are unavailable, particularly  
55 for intrinsically disordered proteins<sup>4,5</sup>. Sequence-based protein language models provide  
56 an alternative by leveraging large-scale evolutionary data to generate functional  
57 sequences<sup>6</sup>; however, most operate in an unconditional or weakly conditioned regime and  
58 are not designed for targeted exploration of homolog families. Conditional sequence  
59 models and evolutionary augmentation approaches can expand sequence diversity but  
60 typically depend on pre-training with family-level labels, fine-tuning on known

61 homologs, or autoregressive sampling schemes that bias generation toward the statistical  
62 center of the training distribution<sup>7-9</sup>.

63 Recent diffusion-based generative models offer a complementary paradigm by generating  
64 all residues simultaneously through iterative denoising, enabling more effective capture  
65 of global context and long-range dependencies than token-by-token autoregressive  
66 models<sup>10,11</sup>. In protein modeling, however, most conditional diffusion frameworks  
67 function primarily as evolutionary inpainting methods: they rely on a multiple sequence  
68 alignment (MSA) as input and recover missing information for the masked query  
69 sequence by interpolating within evolutionary boundaries defined by known  
70 homologs<sup>12,13</sup>. This inward generation regime benefits from dense evolutionary  
71 coordinates supplied by the MSA but inherently limits the novelty of generated  
72 sequences. By contrast, single-sequence exploration requires outward extrapolation—  
73 inferring latent co-evolutionary constraints from a solitary query sequence to generate a  
74 coherent, diverged homolog family *de novo*, without access to an existing alignment.

75 Existing methods lack an explicit mechanism for this type of controlled, query-centric  
76 expansion.

77 Here we introduce EvoGUD (**E**volution-**g**uided **D**iffusion), a single-sequence–  
78 conditioned diffusion framework for large-step exploration of protein sequence space.

79 EvoGUD learns the statistical structure of natural homolog families by training on MSAs  
80 while conditioning exclusively on a single query sequence, enabling inference without  
81 homolog retrieval. A tunable conditioning strength ( $\gamma$ ) controls the balance between  
82 exploratory breadth and adherence to learned evolutionary constraints, allowing direct  
83 generation of substantially diverged yet evolutionarily consistent homolog families

84 within a single generative process. As a result, EvoGUD-generated sequences preserve  
85 natural-like co-evolutionary structure in representation space despite large sequence  
86 divergence, closely tracking natural homologs and substantially exceeding identity-  
87 matched random controls.

88 When assembled as virtual MSAs (vMSAs), EvoGUD-generated sequences substantially  
89 improve AlphaFold3<sup>14</sup> single-sequence inference, restoring both backbone accuracy and  
90 atomic-level side-chain realism for most monomers and protein-nucleic acid complexes,  
91 without evolutionary database search. More importantly, EvoGUD enables functional  
92 discovery in remote sequence space, yielding active TadA variants (an adenine base-  
93 editing enzyme)<sup>15</sup> and large numbers of functional CcdA antitoxin variants (an  
94 intrinsically disordered protein interacted with the toxic protein CcdB)<sup>16</sup> under purely  
95 sequence-level conditioning. Together, these results establish EvoGUD as a single-  
96 sequence-based, evolution-aware generative framework for large-step exploration of  
97 protein sequence space under evolutionary constraints, enabling downstream structure  
98 modeling and functional discovery.

99

## 100 **Results**

### 101 **Generating remote homologs from single protein sequences**

102 EvoGUD was designed to bridge the gap between raw sequence space and the higher-  
103 order evolutionary manifold captured by protein language models<sup>17</sup>. Conceptually (Fig.  
104 1a), both natural evolution and conventional laboratory directed evolution rely on small,  
105 largely undirected mutational steps, and therefore explore only a limited local  
106 neighborhood around a starting sequence<sup>1</sup>. EvoGUD instead enables direct generation of  
107 remote homologs by operating in a feature space defined by ESM-2, conditioning  
108 sequence generation on the query's per-residue embeddings and pairwise attention  
109 patterns<sup>17</sup>.

110 Concretely, EvoGUD was trained as a denoising diffusion model<sup>10,18</sup> to reconstruct  
111 natural MSAs (nMSAs) from noise while observing only query-derived ESM-2 features  
112 (Fig. 1b). During sampling, the model starts from a random amino-acid probability  
113 distribution and iteratively denoises it under conditioning of single query sequence (Fig.  
114 1c). A probability absorption step blends the evolving distribution with the model-  
115 predicted denoised distribution, and a single scalar parameter, the conditioning strength  $\gamma$ ,  
116 modulates the extent to which sampling remains anchored to the conditioning query  
117 sequence. As a result,  $\gamma$  provides an explicit and continuous control over the distance of  
118 generated sequences from the query sequence in sequence space.

119 We first quantified how the conditioning strength  $\gamma$  controls sampling distance by  
120 measuring the sequence identity between generated sequences and their corresponding  
121 query sequences using an independent test set of 159 proteins. This dataset was

122 constructed from recently released PDB monomers<sup>19</sup> and filtered to be non-redundant  
123 with respect to both the training data and within the set itself, applying a 40% sequence-  
124 identity cutoff, denoted as RecentPDB-monomer. (see Methods). Across a wide range of  
125  $\gamma$  values, EvoGUD produces smoothly tunable identity distributions, with increasing  $\gamma$   
126 leading to progressively higher similarity to the query sequence (Fig. 1d), along with  
127 improved attention similarity (Supplementary Fig. S1b), increased foldability according  
128 to pTM from ESMfold (Supplementary Fig. S1c), reduction of diversity according to  
129 intra-set sequence identity (Supplementary Fig. S1d) and decreased novelty according to  
130 the maximum sequence identity to nMSA sequences (Supplementary Fig. S1f). This  
131 demonstrates that EvoGUD does not rely on a fixed exploration regime, but instead  
132 enables controlled interpolation between aggressive exploration and conservative  
133 refinement.

134 To further assess whether EvoGUD preserves evolutionary features beyond simple  
135 sequence identity shown in Fig. 1d, we compared nMSA sequences, EvoGUD-generated  
136 sequences, and random sequences to the query sequence in the representation space of  
137 ESM-2, according to cosine similarity between ESM-2 attention maps by using the  
138 RecentPDB-monomer set. As shown in Fig. 1e, cosine similarity of ESM-2 attention of  
139 EvoGUD-generated sequences to the query sequence closely tracks the similarity  
140 trajectory of natural homologs (natural MSAs) and remain substantially higher than the  
141 similarity of random controls across all identity regimes. Notably, even at low sequence  
142 identity ( $< 0.3$ ), EvoGUD maintains a median attention-space similarity of 0.874,  
143 compared to 0.578 for random sequences, indicating that EvoGUD samples realistic  
144 regions of the evolutionary manifold rather than merely matching identity statistics.

145 Comparable trends were observed using ESM-2 embedding similarity instead of attention  
146 (Supplementary Fig. S2). Together, these results show that EvoGUD enables controllable  
147 generation of structurally plausible remote homologs while preserving higher-order  
148 evolutionary constraints, with  $\gamma$  acting as an explicit knob to balance exploration and  
149 constraint.

150

### 151 **Monomeric structure prediction with vMSA from EvoGUD**

152 We next asked whether vMSAs generated by EvoGUD can replace nMSAs for  
153 monomeric protein structure prediction with AlphaFold3<sup>14</sup> under single-sequence  
154 conditions. We employed EvoGUD to generate vMSAs across a range of conditioning  
155 strengths  $\gamma$  and vMSA depths, and the resulting structures were predicted using  
156 AlphaFold3 in single-sequence mode (AF3-SS) using an ensemble of 10 EvoGUD  
157 parameter settings (See Methods).

158 EvoGUD + AF3-SS substantially outperformed AF3-SS on the RecentPDB-monomer  
159 test set (Fig. 2a). The mean TM-score<sup>20</sup> increased from 0.476 for AF3-SS to 0.795 for  
160 EvoGUD + AF3-SS, approaching the performance of ESMFold (0.827) and AF3 with  
161 nMSAs and PDB templates (0.884). The fraction of targets with TM-score  $\geq 0.5$ —a  
162 commonly used criterion for correct fold assignment—increased from 40.9% for AF3-SS  
163 to 89.3% for EvoGUD + AF3-SS, comparable to ESMFold (92.4%) and approaching  
164 AF3 with nMSAs and templates (96.2%). These results indicate that most monomers can  
165 be modeled at near-native resolution from a single input sequence using EvoGUD.

166 To assess generalization beyond RecentPDB-monomer, we evaluated EvoGUD on the  
167 CASP15<sup>21</sup> monomer benchmark (N = 71), which was not used during parameter



168 selection. EvoGUD + AF3-SS again markedly improved backbone accuracy relative to  
169 AF3-SS and yielded TM-score distributions comparable to those observed on  
170 RecentPDB-monomer (Fig. 2b), demonstrating that the selected vMSA settings  
171 generalize to an independent community benchmark without re-tuning. On CASP15, the  
172 mean TM-score increased from 0.403 for AF3-SS to 0.593 for EvoGUD + AF3-SS,  
173 approaching the performance of ESMFold (0.640) and AF3 with nMSAs and PDB  
174 templates (0.740). The fraction of targets with TM-score  $\geq 0.5$  increased from 25.4% for  
175 AF3-SS to 64.8% for EvoGUD + AF3-SS, comparable to ESMFold (66.2%) and  
176 approaching AF3 with nMSAs and PDB structure templates (77.5%). A representative  
177 example is shown in Fig. 2c, where EvoGUD + AF3-SS recovers the correct overall fold  
178 and domain arrangement that is missed by AF3-SS.

179 Because AlphaFold3 is an all-atom diffusion model, we further evaluated the physical  
180 plausibility of the predicted structures by quantifying steric clashes between heavy atoms.  
181 ESMFold, despite its strong backbone accuracy, exhibited substantially higher clash  
182 counts, with a broad tail of severely mispacked models on both RecentPDB-monomer  
183 and CASP15 (Fig. 2d and 2e). In contrast, EvoGUD + AF3-SS markedly reduced steric  
184 clashes relative to ESMFold, producing distributions much closer to those of AF3. A  
185 representative local side-chain environment is shown in Fig. 2f, where EvoGUD + AF3-  
186 SS yields well-packed, stereochemically reasonable side chains, whereas ESMFold  
187 displays strained rotamers and steric overlaps<sup>22</sup>.

188 Together, these results show that EvoGUD enables single-sequence AlphaFold3 to  
189 recover both high backbone accuracy and AF3-like all-atom realism for monomeric  
190 proteins, without requiring time-consuming searches for natural homologs. We next

191 asked whether similar gains extend beyond monomers to multimeric and protein–nucleic-  
192 acid complexes.

193

### 194 **Protein-NA complex prediction with vMSA from EvoGUD**

195 We next evaluated whether EvoGUD can extend single-sequence prediction beyond  
196 monomeric proteins to multi-chain assemblies, focusing on protein–nucleic-acid  
197 (protein–NA) complexes (Fig. 3). We curated a benchmark of 165 experimentally  
198 determined protein–DNA/RNA complexes and generated vMSAs for the protein chains  
199 only, while keeping the nucleic-acid sequences fixed. These sequences were assembled  
200 into 10 sets of vMSA and supplied, together with the original NA sequence, to AF3-SS.  
201 For each target we selected the final EvoGUD + AF3-SS model by the highest ipTM  
202 confidence score. As shown in Fig. 3a, EvoGUD + AF3-SS substantially improved the  
203 accuracy of the protein subunits within the protein–NA complexes (median TM-score =  
204 0.872) relative to AF3-SS (median TM-score = 0.436), close to the distribution of AF3  
205 (median TM-score = 0.928), as in the case of monomeric proteins.

206 For complex structures evaluated by interface local distance difference test (iLDDT<sup>23</sup>) (as  
207 in AF3<sup>14</sup>), AF3-SS frequently produced mis-docked subunits and distorted protein–  
208 nucleic-acid contacts, with most targets exhibiting low iLDDT values below 0.4 (Fig.  
209 3b). In contrast, EvoGUD + AF3-SS markedly shifted the distribution toward higher  
210 iLDDT values, correctly recovering 55 well-docked complexes out of 90 that are  
211 correctly predicted by AF3 using nMSAs for both protein and nucleic acid sequences  
212 together with PDB templates—substantially exceeding the 14 complexes recovered by  
213 AF3-SS. These results indicate that vMSAs generated from single sequences can

214 effectively guide accurate docking of protein chains onto DNA and RNA for most cases  
215 studied.  
216 To visualize these trends on a per-complex basis, we compared the change in TM-score  
217 and iLDDT relative to AF3-SS for each method (Fig. 3c). Most EvoGUD + AF3-SS  
218 predictions fall in the quadrant corresponding to simultaneous gains in subunit accuracy  
219 and interface quality, demonstrating that vMSAs derived from single sequences can  
220 improve both global folding and interfacial organization within the same model.  
221 One structural example illustrates these effects at the level of individual assemblies (Fig.  
222 3d). In the shown protein–NA complex (PDB: 7U7C)<sup>24</sup>, AF3-SS fails to correctly  
223 position protein subunits relative to the nucleic acid and produces poorly resolved  
224 interfaces. In contrast, EvoGUD + AF3-SS accurately recovers the overall architecture  
225 and correctly docks the protein chains onto the nucleic-acid scaffold, yielding an  
226 interface geometry that closely matches the experimental structure and approaches the  
227 AF3 baseline. Thus, vMSAs from EvoGUD can capture the evolution signals not only in  
228 monomeric structures but also in interfacial structures for docking.

229

### 230 **Locating remote functional enzymes by EvoGUD**

231 To test whether EvoGUD can recover functional enzymes from remote regions of  
232 sequence space, we selected TadA, a bacterial tRNA adenosine deaminase that has been  
233 repurposed through extensive directed evolution into the catalytic core of adenine base  
234 editors (ABEs), enabling programmable A•T→G•C DNA conversion<sup>15</sup>. Although highly  
235 effective, previously reported TadA variants remain closely related to their ancestral  
236 sequences, motivating exploration of more distant sequence solutions<sup>25,26</sup>.

237 Based on the joint behavior of sequence identity, foldability, diversity, and novelty  
238 (Supplementary Fig. S3), we selected a conditioning strength of  $\gamma = 2$  and generated  
239 1,024 sequences conditioned on the wild type TadA from *Staphylococcus aureus*<sup>27</sup>. Ten  
240 representative TadA variants were selected for experimental evaluation (see Methods).  
241 Structural evaluation using AlphaFold3 (AF3) showed that providing EvoGUD-generated  
242 sequences as vMSAs substantially improved structure prediction compared with single-  
243 sequence input, yielding well-folded TadA-like architectures with markedly higher pTM  
244 confidence score and TM-score (shown in Fig. 4a).  
245 The ten EvoGUD-generated variants were then evaluated with a trimethoprim (TMP)  
246 resistance reversion assay in *E. coli*, in which TadA-mediated A→G editing restores a  
247 functional *R67* dihydrofolate reductase (DHFR) reporter gene and as a result, a stronger  
248 active TadA variant will grow more colonies (Fig. 4b). Some examples are shown in Fig.  
249 4c. Eight of the ten variants restored TMP resistance, demonstrating robust A•T→G•C  
250 DNA-editing activity despite low sequence identity to both the query sequence (0.38–  
251 0.40) and any known TadA homologs (0.57–0.62). Quantification of editing activity was  
252 made according to the number of colonies before and after TMP selection (see Methods).  
253 It revealed a reproducible range of activities across variants, with the activity of the wild  
254 type TadA falling within the distribution of EvoGUD-generated sequences as shown in  
255 Fig. 4d. Thus, EvoGUD can recover functional TadA DNA-editing enzymes from  
256 previously unexplored, remote regions of sequence space.  
257

## 258 **Application to the intrinsically disordered antitoxin CcdA**

259 To further test the limits of EvoGUD, we examined whether it could generate functional  
260 variants of the intrinsically disordered antitoxin CcdA. CcdA lacks a stable structure in  
261 isolation and acquires its functional conformation only upon binding its cognate toxin  
262 CcdB, posing a stringent challenge for protein design<sup>16,28</sup>.

263 Monomeric CcdA and CcdB sequences were derived from the *E. coli* CcdA–CcdB  
264 complex (PDB: 3HPW)<sup>16</sup>. These sequences were then assembled into a symmetrized  
265 CcdB–G<sub>50</sub>–CcdA–G<sub>50</sub>–CcdB fusion construct. The long flexible link was specifically  
266 designed as the query for EvoGUD-based CcdA generation under CcdB conditioning.  
267 Based on the joint behavior of identity, foldability, diversity, and novelty (Supplementary  
268 Fig. S4), we selected a conditioning strength of  $\gamma = 2$  for downstream experiments and  
269 generated a pooled library of 5,623 unique CcdA variants for experimental screening (see  
270 Methods).

271 For structural verification, only a vMSA constructed from EvoGUD-generated CcdA  
272 sequences was inputted into AlphaFold3, without using any nMSAs or vMSA for CcdB.  
273 A representative subset of eight CcdA variants generated at  $\gamma = 2$  was selected to form the  
274 vMSA, as prediction accuracy decreased when larger numbers of low- $\gamma$  sequences ( $> 64$ )  
275 were included (Supplementary Fig. S5). As shown in Fig. 5a, AF3 predictions using this  
276 CcdA-only vMSA recover a coherent CcdA–CcdB complex. Relative to the  
277 experimentally determined *E. coli* CcdA–CcdB structure, the predicted model shows high  
278 agreement at both the subunit and interface levels, as reflected by elevated pTM, ipTM,  
279 TM-score, and iLDDT values. Notably, despite containing sequences for CcdA only, the  
280 vMSA improves the predicted structures of both the CcdA antitoxin and the flanking

281 CcdB toxin subunits, as well as their binding interface. In contrast, AF3 single-sequence  
282 inference fails to recover the CcdA structure or its docked configuration within CcdB,  
283 indicating that EvoGUD-generated CcdA sequence ensembles provide sufficient context  
284 for accurate prediction of the full toxin–antitoxin complex, without natural or virtual  
285 MSA for CcdB.

286 The variant library was subjected to an experimental CcdB toxin selection assay, and  
287 variant counts were obtained by high-throughput sequencing before and after selection in  
288 two independent biological replicates (Fig. 5b). Fitness was inferred from before/after  
289 enrichment using early-stop variants as internal negative controls, followed by  
290 Benjamini–Hochberg false discovery rate (BH-FDR)<sup>29</sup> filtering ( $q \leq 0.01$ ) and  
291 normalization (see Methods). Across the two replicates, 1,110 and 1,153 variants passed  
292 the survival test, respectively. Requiring consistent enrichment in both experiments  
293 identified 1,072 functional variants, corresponding to an overall success rate of 19%  
294 (1,072 of 5,623 variants). Both replicates show a substantial population of variants with  
295 normalized  $\log_2$  fitness exceeding that of wild-type CcdA (*E. coli* CcdA, or EcCcdA, Fig.  
296 5c), and inferred fitness values for the same variants are highly correlated between  
297 replicates (Fig. 5d), enabling robust ranking of functional variants.

298 To experimentally validate the statistical classification, we evaluated a subset of nine  
299 individual variants spanning the fitness range using spot survival assays (Supplementary  
300 Fig. S6 and S7). These variants (EvoGUD generated CcdA, denoted as egCcdA) were  
301 chosen according to their rank based on average normalized  $\log_2$  fitness across replicates  
302 between 3 and 18. (egCcdA-1 denotes the highest-ranked variant). Across both plate  
303 experiments, all nine tested variants exhibited detectable protection from CcdB toxicity.

304 No growth was observed in the CcdB-only control, demonstrating that the enrichment in  
305 the high-throughput pooled assay reflects genuine antitoxin function. As illustrative  
306 examples, a representative subset of five variants is shown in the main text (Fig. 5c and  
307 5e), chosen to illustrate the correspondence between inferred fitness and phenotypic  
308 strength across a wide dynamic range. High-ranked variants such as egCcdA-1 and  
309 egCcdA-14 displayed activity comparable to or exceeding that of wild-type CcdA at  
310 different dilutions (y-axis), whereas egCcdA-78 showed moderately reduced activity,  
311 consistent with its lower inferred fitness. As a variant positioned near the statistical  
312 decision boundary, egCcdA-933 still exhibited weak but detectable rescue relative to the  
313 negative control (Fig. 5f), validating the sensitivity of the fitness-based classification.  
314 Thus, EvoGUD can generate large numbers of functional CcdA variants under purely  
315 sequence-level conditioning, despite the absence of a stable ground-state fold. Unedited  
316 plate images for all tested variants are provided in Supplementary Figs. S6 and S7,  
317 ensuring full transparency of the experimental results.  
318

## 319 **Discussion**

320 EvoGUD was designed to enable large-step exploration of protein sequence space while  
321 preserving the higher-order evolutionary constraints that characterize natural protein  
322 families. A central finding of this work is that EvoGUD-generated sequences remain  
323 embedded within realistic evolutionary manifolds despite substantial sequence  
324 divergence. In ESM-2 embedding and attention spaces, generated sequences closely track  
325 natural homologs and remain far more consistent with the query in representation space  
326 than identity-matched random controls, indicating that EvoGUD captures co-evolutionary  
327 structure beyond residue-level conservation. The conditioning strength  $\gamma$  provides  
328 continuous control over the balance between exploratory breadth and evolutionary  
329 constraint.

330 These properties translate directly into improved structure prediction. By assembling  
331 EvoGUD-generated sequences as vMSAs, AlphaFold3 single-sequence inference  
332 recovers much of the accuracy and atomic detail typically associated with natural  
333 homolog searches. Across monomer benchmarks, EvoGUD-assisted predictions are near  
334 standard AlphaFold3 performance. For protein–DNA and protein–RNA complexes,  
335 EvoGUD-derived vMSAs further improve interface geometry, demonstrating that the  
336 generated sequence families encode actionable co-evolutionary signals for multimeric  
337 recognition.

338 Besides structure modeling, EvoGUD provides a general framework for functional  
339 protein discovery in remote sequence space. For the TadA enzyme, EvoGUD identified  
340 highly divergent yet functional DNA-editing variants, revealing functional solutions  
341 inaccessible to stepwise directed evolution. Notably, these active variants reside at



342 approximately 40% sequence identity to the query, corresponding—in natural  
343 evolutionary terms—to divergence accumulated over on the order of  $\sim 10^9$  years  
344 (Supplementary Table S4)<sup>30</sup>. Similarly, in the CcdA–CcdB toxin–antitoxin system,  
345 EvoGUD enabled large-scale discovery of functional antitoxin variants spanning 54–73%  
346 sequence identity, even among low-ranked candidates, indicating robust preservation of  
347 context-dependent functional constraints without explicit structural or biophysical  
348 scoring.

349 A key requirement for sequence generative models is generalizability across protein  
350 families and evolutionary distances. EvoGUD exhibits robust generalization across  
351 multiple validation regimes, including stringent non-redundant subsets, with generated  
352 sequences consistently following the same representation-space trajectories as natural  
353 homologs irrespective of training-set proximity (Supplementary Fig. S8). Controlled  
354 experiments on TadA further show that excluding or retaining close homologs during  
355 training produces only minor shifts in generative behavior without evidence of collapse  
356 (Supplementary Fig. S9). Notably, the CcdA system provides complementary insight:  
357 although CcdA homologs were present in the training data, CcdA-only conditioning  
358 yielded low identity and reduced foldability, whereas introducing an unseen fusion  
359 context with its cognate binding partner systematically shifted generation toward  
360 functionally coherent sequence space (Supplementary Fig. S10). Together, these results  
361 indicate that EvoGUD does not rely on memorization of training sequences but is  
362 primarily shaped by the evolutionary constraints supplied at inference time.

363 EvoGUD adopts a modular alternative to end-to-end single-sequence structure prediction  
364 pipelines. By decoupling protein sequence feature extraction from structure inference

365 with intermediate sequence generation, improvements in protein language models or  
366 structure predictors can be incorporated by retraining only a lightweight generative  
367 adapter, requiring on the order of a single GPU-day. This contrasts sharply with the  
368 hundreds of GPU-weeks typically required to train or adapt full-scale structure prediction  
369 models<sup>17</sup>, enabling rapid iteration and reuse of advances in representation learning.  
370 Although EvoGUD's performance is bounded by the fidelity of underlying sequence  
371 representations, its modular design provides a scalable and extensible framework for  
372 integrating advances in representation learning, enabling large-scale evolutionary  
373 exploration that can be naturally combined with local refinement strategies such as  
374 directed evolution. In this way, EvoGUD bridges global sequence-space exploration and  
375 functional protein engineering while remaining grounded in evolutionary realism.  
376

## 377 **Methods**

### 378 **Training Dataset and Data Preprocessing**

379 EvoGUD was trained on OpenProteinSet-PDB<sup>31</sup>, a curated reconstruction of the  
380 AlphaFold2 training dataset as implemented in OpenFold<sup>32</sup>, comprising 131,487 protein  
381 chains with precomputed MSAs. To ensure sequence integrity and computational  
382 consistency, chains containing unknown amino acids (“X”) were excluded. Sequences  
383 were further filtered by length, retaining chains with  $30 \leq L \leq 1000$  amino acids. After  
384 filtering, the final dataset comprised 117,556 entries, which were partitioned into a  
385 training set of 116,756 entries and a validation set of 800 entries.

386

### 387 **Model architecture and conditioning**

388 EvoGUD is built on a Diffusion Transformer (DiT) backbone<sup>33</sup> with a modified adaLN-  
389 Zero conditioning mechanism (Supplementary Fig. S11). Whereas the original DiT  
390 conditions on global features, EvoGUD incorporates sequence-specific evolutionary  
391 context derived from a protein language model.  
392 During both training and inference, EvoGUD conditions the denoising process on  
393 representations extracted from the ESM-2 3B model, including per-residue embeddings  
394 (2,560 dimensions), and pairwise attention maps (36 layers  $\times$  40 heads; 1,440  
395 dimensions)<sup>17</sup>. These features are linearly projected into a 128-dimensional latent space.  
396 In each DiT block, ESM-2 embeddings are injected via cross-attention as key–value  
397 pairs, while projected attention maps are added to the attention logits as a pairwise bias.  
398 The resulting representations generate the six modulation parameters (shift, scale, and

399 gate for attention and feed-forward sublayers) used in the adaLN-Zero operation,  
400 enabling position-wise, evolution-aware modulation of the network.  
401 The model comprises 6 DiT blocks with a hidden dimension of 128 (3.06 M parameters  
402 total). A global dropout rate of 0.1 was applied. To stabilize early training, the linear  
403 layers producing adaLN modulation parameters were zero-initialized, such that each DiT  
404 block initially behaves as an identity mapping. Model outputs are projected to a  
405 categorical distribution over a 21-token alphabet (20 amino acids plus a gap token).  
406

#### 407 **Diffusion formulation**

408 EvoGUD adopts a continuous-time Gaussian diffusion framework with a cosine noise  
409 schedule<sup>34</sup>. The forward process gradually perturbs one-hot encoded amino-acid  
410 sequences with Gaussian noise according to a cumulative signal retention coefficient  $\bar{\alpha}_t$ ,  
411 defined over normalized time  $u \in [0,1]$  as:

$$412 \quad \bar{\alpha}_u = \frac{f(u)}{f(0)}, \quad f(u) = \cos^2 \left( \frac{u}{1+s} + \frac{1}{1+s} \cdot \frac{\pi}{2} \right)$$

413 where  $s$  is a small offset to prevent the noise level from becoming too small at  $t = 0$ . The  
414 schedule was discretized into  $T = 100$  steps for training. The model is trained to predict  
415 the original categorical distribution of clean sequences from noisy inputs, conditioned on  
416 ESM-2 features.

417

#### 418 **Model Training**

419 EvoGUD was trained to reconstruct natural homolog sequences drawn from MSAs using  
420 a query-centric batching strategy. For each optimization step, a single query sequence

421 was paired with 64 target sequences sampled from its associated MSA. For shallow  
422 MSAs, target sequences were oversampled; for deep MSAs, 64 members were randomly  
423 subsampled. The query sequence was always included in the batch to anchor  
424 reconstruction.

425 Each target sequence was one-hot encoded and independently corrupted by Gaussian  
426 noise at a randomly sampled timestep  $t_i \in [1, T]$ . The model predicted the denoised  
427 categorical distribution  $p(x_0 | x_{t,i}, cond)$  over the 21-token alphabet.

428 To emphasize evolutionary diversity rather than trivial sequence conservation, a per-  
429 residue weighted cross-entropy loss was applied: positions differing from the query were  
430 assigned weight = 1.0, whereas positions identical to the query or corresponding to gaps  
431 were down-weighted (weight = 0.1).

432 Training was performed for 100 epochs using the AdamW optimizer<sup>35</sup> with a learning  
433 rate of  $1 \times 10^{-3}$  and automatic mixed precision. Each epoch comprised 10,000 unique  
434 queries sampled from the whole training set, totaling  $\sim 10^6$  optimization steps. Training  
435 required approximately 25 h on a single NVIDIA A100 GPU and was implemented in  
436 PyTorch 2.2<sup>36</sup>.

437

### 438 **Single-sequence conditional sampling via probability absorption**

439 At inference time, EvoGUD generates homolog families from a single query sequence  
440 using a probability absorption sampling scheme that bridges continuous diffusion  
441 dynamics with discrete sequence space.

442 Sampling begins from isotropic Gaussian noise  $x_T \sim \mathcal{N}(0, I)$  and proceeds over  $T = 100$   
443 discrete reverse-diffusion steps. At each step, the model predicts a position-wise amino-

444 acid probability distribution conditioned on ESM-2 features of the query. A discrete  
445 sequence is obtained by deterministic argmax decoding, excluding the gap token to  
446 generate gap-free sequences.

447 The sampled sequence is projected back into latent space as a centered one-hot  
448 representation, scaled by a conditioning strength parameter  $\gamma$ :

$$449 \quad \hat{x}_0 = \gamma \cdot (\text{one\_hot}(S^0) - 0.5)$$

450 The latent state is then updated via the Gaussian reverse transition, blending the absorbed  
451 identity signal with stochastic noise. The parameter  $\gamma$  controls the trade-off between  
452 evolutionary adherence and exploratory breadth: higher values promote conservative,  
453 high-confidence homologs, whereas lower values allow broader exploration of remote  
454 sequence space.

455

#### 456 **Benchmark test sets**

457 **RecentPDB-monomer.** To evaluate structure prediction performance, an independent  
458 test set was curated from PDB entries released between January 1 and July 1, 2024,  
459 following AlphaFold3 benchmarking principles<sup>14</sup>. Only protein-only monomers solved by  
460 X-ray crystallography at  $\leq 2.0$  Å resolution were retained. Chains were filtered to lengths  
461 of 30–500 residues and subjected to 40% sequence-identity filtering both within the set  
462 and against the training/validation data. The final set comprised 159 non-redundant  
463 monomers.

464 **RecentPDB-multimer.** Protein–nucleic-acid complexes were derived from the  
465 AlphaFold3 benchmark dataset<sup>14</sup>. After excluding entries lacking protein chains or

466 containing non-canonical residues, the final set comprised 165 complexes. vMSAs were  
467 generated only for protein components, while nucleic-acid sequences were held fixed.

468

#### 469 **Ensemble settings for EvoGUD + AF3-SS**

470 A grid search over conditioning strength  $\gamma$  and vMSA depth on the RecentPDB-monomer  
471 benchmark set revealed a broad operating regime in which EvoGUD-derived vMSAs  
472 substantially improved structure prediction accuracy relative to AF3-SS (Supplementary  
473 Fig. S5). Within this regime, increasing  $\gamma$  and vMSA depth (# of generated vMSA  
474 sequences) led to substantial gains in both predicted confidence (pTM<sup>14</sup>) and backbone  
475 accuracy (TM-score), indicating that vMSAs capture much of the evolutionary  
476 information normally supplied by natural homologs. Based on this analysis, we selected a  
477 small ensemble of 10 EvoGUD parameter settings spanning this robust regime ( $\gamma = 1$   
478 with #vMSA  $\in \{2,4,8,16\}$  and  $\gamma = 2$  with #vMSA  $\in \{2,4,8,16,32,64\}$ ) and fixed these  
479 settings for all subsequent analyses.

480

#### 481 **Generation and validation of TadA variants**

482 Wild-type *Staphylococcus aureus* TadA (PDB: 2B3J)<sup>27</sup> was used as the query for  
483 EvoGUD sequence generation. Conditioning strength  $\gamma$  was selected based on predicted  
484 foldability, novelty, and diversity ( $\gamma = 2$ ; Supplementary Fig. S3). A total of 1,024 TadA  
485 variants were generated under a co-generation identity constraint of  $\geq 35\%$  relative to  
486 wild-type TadA, subsequently clustered at 70% sequence identity, and ten representatives  
487 from the largest clusters were selected for experimental validation.

488 TadA activity was quantified using a trimethoprim resistance reversion assay in *E. coli*  
489 based on a premature stop-codon reporter. Editing activity was defined as per-base  
490 mutation rates estimated from the observed frequency of TMP-resistant colonies, using a  
491 Luria–Delbrück approximation<sup>37</sup>. Full experimental protocols are provided in the  
492 Supplementary Information.

493

#### 494 **Generation and validation of CcdA variants**

495 CcdA variants were generated using a conditional single-chain strategy in which the  
496 antitoxin sequence was embedded within a symmetrized CcdB–G<sub>50</sub>–CcdA–G<sub>50</sub>–CcdB  
497 fusion, with wild type CcdA and CcdB sequences from *E. coli* (PDB: 3HPW)<sup>16</sup>. Here, G<sub>50</sub>  
498 denotes a 50-amino-acid poly-glycine linker that spatially separates the CcdA and CcdB  
499 domains while preserving sequence-level context<sup>38</sup>. Residues outside the CcdA region  
500 were fixed during sampling. Conditioning strength  $\gamma = 2$  was selected according to  
501 predicted foldability, novelty, and diversity (Supplementary Fig. S4).

502 A pooled library of 10,000 CcdA variants was generated under a co-generation identity  
503 constraint of  $\leq 75\%$  relative to wild-type CcdA and subsequently clustered at 90%  
504 sequence identity, yielding 5,623 unique sequences. Functional selection was performed  
505 using a toxin-rescue assay in *E. coli*, followed by deep sequencing.

506 Variant fitness was estimated using a Poisson-based  $\log_2$  enrichment model and false-  
507 discovery-rate (FDR) correction<sup>29</sup>. Variants passing  $\text{FDR} \leq 0.01$  in two independent  
508 experiments were considered functional. Full experimental protocols are provided in the  
509 Supplementary Information.

510



## 511 **Reference**

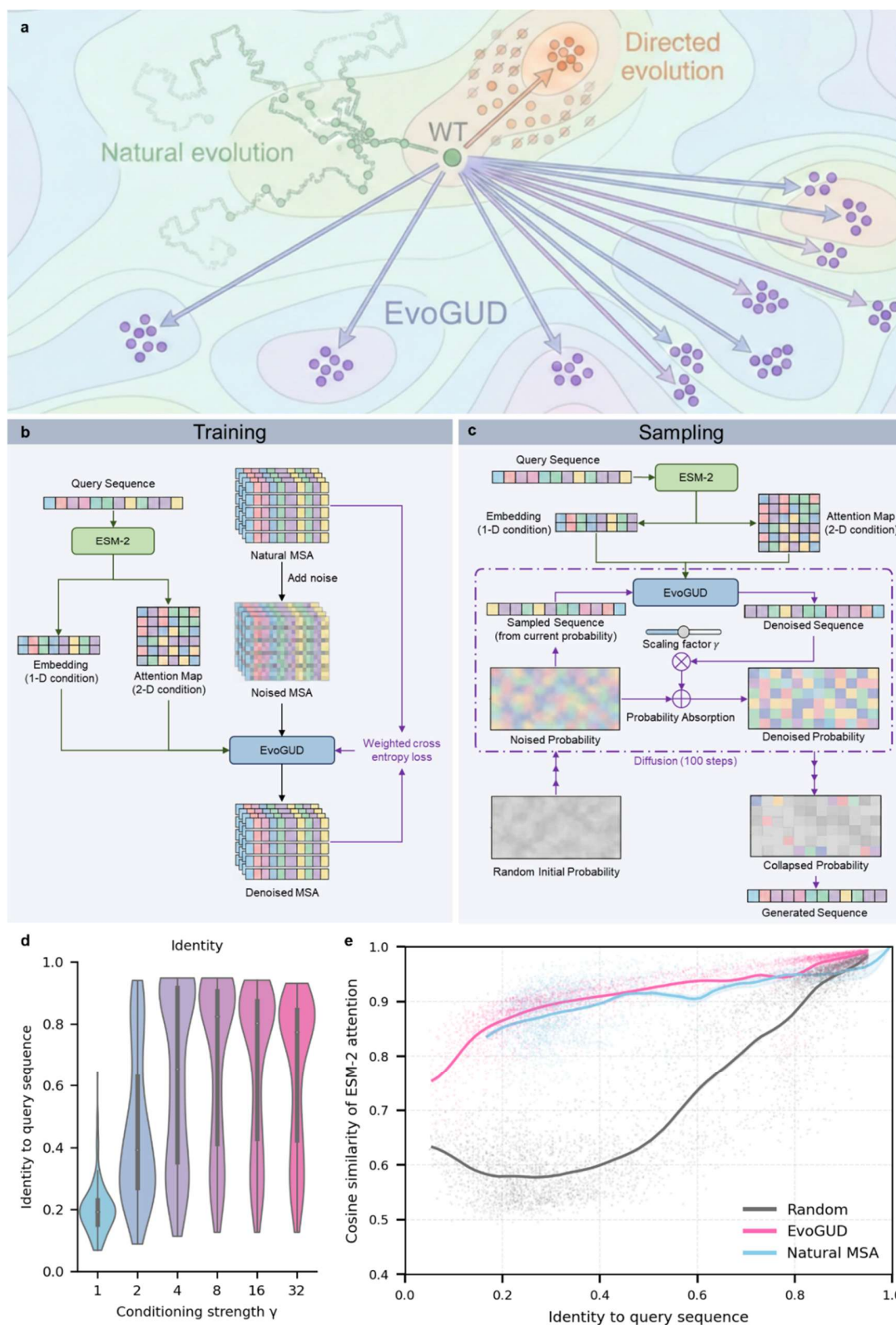
- 512 1. Arnold, F. H. Design by directed evolution. *Accounts of chemical research* **31**, 125–131  
513 (1998).
- 514 2. Romero, P. A. & Arnold, F. H. Exploring protein fitness landscapes by directed  
515 evolution. *Nature reviews Molecular cell biology* **10**, 866–876 (2009).
- 516 3. Poelwijk, F. J., Tănase-Nicola, S., Kiviet, D. J. & Tans, S. J. Reciprocal sign epistasis is a  
517 necessary condition for multi-peaked fitness landscapes. *Journal of theoretical biology*  
518 **272**, 141–144 (2011).
- 519 4. Huang, P.-S., Boyken, S. E. & Baker, D. The coming of age of de novo protein design.  
520 *Nature* **537**, 320–327 (2016).
- 521 5. Dauparas, J. *et al.* Robust deep learning–based protein sequence design using  
522 ProteinMPNN. *Science* **378**, 49–56 (2022).
- 523 6. Rives, A. *et al.* Biological structure and function emerge from scaling unsupervised  
524 learning to 250 million protein sequences. *Proceedings of the National Academy of*  
525 *Sciences* **118**, e2016239118 (2021).
- 526 7. Madani, A. *et al.* Large language models generate functional protein sequences across  
527 diverse families. *Nature biotechnology* **41**, 1099–1106 (2023).
- 528 8. Nijkamp, E., Ruffolo, J. A., Weinstein, E. N., Naik, N. & Madani, A. Progen2: exploring  
529 the boundaries of protein language models. *Cell systems* **14**, 968–978. e3 (2023).
- 530 9. Zhang, J. *et al.* Unsupervisedly prompting AlphaFold2 for accurate few-shot protein  
531 structure prediction. *Journal of Chemical Theory and Computation* **19**, 8460–8471  
532 (2023).

- 533 10. Ho, J., Jain, A. & Abbeel, P. Denoising diffusion probabilistic models. *Advances in*  
534 *neural information processing systems* **33**, 6840–6851 (2020).
- 535 11. Trippe, B. L. *et al.* Diffusion probabilistic modeling of protein backbones in 3d for the  
536 motif-scaffolding problem. *arXiv preprint arXiv:2206.04119* (2022).
- 537 12. Alamdari, S. *et al.* Protein generation with evolutionary diffusion: sequence is all you  
538 need. *BioRxiv* 2023.09. 11.556673 (2023).
- 539 13. Sgarbossa, D., Lupo, U. & Bitbol, A.-F. Generative power of a protein language model  
540 trained on multiple sequence alignments. *Elife* **12**, e79854 (2023).
- 541 14. Abramson, J. *et al.* Accurate structure prediction of biomolecular interactions with  
542 AlphaFold 3. *Nature* **630**, 493–500 (2024).
- 543 15. Gaudelli, N. M. *et al.* Programmable base editing of A•T to G•C in genomic DNA  
544 without DNA cleavage. *Nature* **551**, 464–471 (2017).
- 545 16. De Jonge, N. *et al.* Rejuvenation of CcdB-poisoned gyrase by an intrinsically  
546 disordered protein domain. *Molecular cell* **35**, 154–163 (2009).
- 547 17. Lin, Z. *et al.* Evolutionary-scale prediction of atomic-level protein structure with a  
548 language model. *Science* **379**, 1123–1130 (2023).
- 549 18. Austin, J., Johnson, D. D., Ho, J., Tarlow, D. & Van Den Berg, R. Structured  
550 denoising diffusion models in discrete state-spaces. *Advances in neural information*  
551 *processing systems* **34**, 17981–17993 (2021).
- 552 19. Berman, H. M. *et al.* The protein data bank. *Nucleic acids research* **28**, 235–242  
553 (2000).

- 554 20. Zhang, Y. & Skolnick, J. Scoring function for automated assessment of protein  
555 structure template quality. *Proteins: Structure, Function, and Bioinformatics* **57**, 702–710  
556 (2004).
- 557 21. Kryshtafovych, A., Schwede, T., Topf, M., Fidelis, K. & Moult, J. Critical assessment  
558 of methods of protein structure prediction (CASP)—Round XV. *Proteins: Structure,*  
559 *Function, and Bioinformatics* **91**, 1539–1549 (2023).
- 560 22. Williams, C. J. *et al.* MolProbity: more and better reference data for improved all-atom  
561 structure validation. *Protein Science* **27**, 293–315 (2018).
- 562 23. Mariani, V., Biasini, M., Barbato, A. & Schwede, T. IDDT: a local superposition-free  
563 score for comparing protein structures and models using distance difference tests.  
564 *Bioinformatics* **29**, 2722–2728 (2013).
- 565 24. Chang, C., Lee Luo, C. & Gao, Y. In crystallo observation of three metal ion promoted  
566 DNA polymerase misincorporation. *Nature Communications* **13**, 2346 (2022).
- 567 25. Richter, M. F. *et al.* Phage-assisted evolution of an adenine base editor with improved  
568 Cas domain compatibility and activity. *Nature biotechnology* **38**, 883–891 (2020).
- 569 26. Lapinaite, A. *et al.* DNA capture by a CRISPR-Cas9–guided adenine base editor.  
570 *Science* **369**, 566–571 (2020).
- 571 27. Losey, H. C., Ruthenburg, A. J. & Verdine, G. L. Crystal structure of *Staphylococcus*  
572 *aureus* tRNA adenosine deaminase TadA in complex with RNA. *Nature structural &*  
573 *molecular biology* **13**, 153–159 (2006).
- 574 28. Loris, R. *et al.* Crystal structure of CcdB, a topoisomerase poison from *E. coli*. *Journal*  
575 *of molecular biology* **285**, 1667–1677 (1999).

- 576 29. Benjamini, Y. & Hochberg, Y. Controlling the false discovery rate: a practical and  
577 powerful approach to multiple testing. *Journal of the Royal statistical society: series B*  
578 *(Methodological)* **57**, 289–300 (1995).
- 579 30. Kumar, S., Stecher, G., Suleski, M. & Hedges, S. B. TimeTree: a resource for  
580 timelines, timetrees, and divergence times. *Molecular biology and evolution* **34**, 1812–  
581 1819 (2017).
- 582 31. Ahdritz, G. *et al.* OpenProteinSet: Training data for structural biology at scale.  
583 *Advances in Neural Information Processing Systems* **36**, 4597–4609 (2023).
- 584 32. Ahdritz, G. *et al.* OpenFold: Retraining AlphaFold2 yields new insights into its  
585 learning mechanisms and capacity for generalization. *Nature methods* **21**, 1514–1524  
586 (2024).
- 587 33. Peebles, W. & Xie, S. Scalable diffusion models with transformers. in *Proceedings of*  
588 *the IEEE/CVF international conference on computer vision* 4195–4205 (2023).
- 589 34. Nichol, A. Q. & Dhariwal, P. Improved denoising diffusion probabilistic models. in  
590 *International conference on machine learning* 8162–8171 (PMLR, 2021).
- 591 35. Loshchilov, I. & Hutter, F. Decoupled weight decay regularization. *arXiv preprint*  
592 *arXiv:1711.05101* (2017).
- 593 36. Paszke, A. *et al.* Pytorch: An imperative style, high-performance deep learning library.  
594 *Advances in neural information processing systems* **32**, (2019).
- 595 37. Luria, S. E. & Delbrück, M. Mutations of bacteria from virus sensitivity to virus  
596 resistance. *Genetics* **28**, 491 (1943).
- 597 38. Chen, X., Zaro, J. L. & Shen, W.-C. Fusion protein linkers: property, design and  
598 functionality. *Advanced drug delivery reviews* **65**, 1357–1369 (2013).

599 **Figures**



600

601 **Figure 1.** EvoGUD enables controllable, large-step exploration of protein sequence space  
602 from a single query.

603 a, Conceptual landscape illustrating the limitations of natural evolution and laboratory  
604 directed evolution, and how EvoGUD enables large-step sampling toward remote  
605 functional sequence regions.

606 b, Training of EvoGUD using nMSAs: the model is trained to denoise corrupted MSAs  
607 conditioned on query-derived ESM-2 embeddings and pairwise attention maps.

608 c, Sampling procedure: starting from a random probability distribution, EvoGUD  
609 iteratively denoises sequence probabilities under query conditioning, with a scaling factor  
610  $\gamma$  controlling the strength of probability absorption before collapsing to a discrete  
611 sequence.

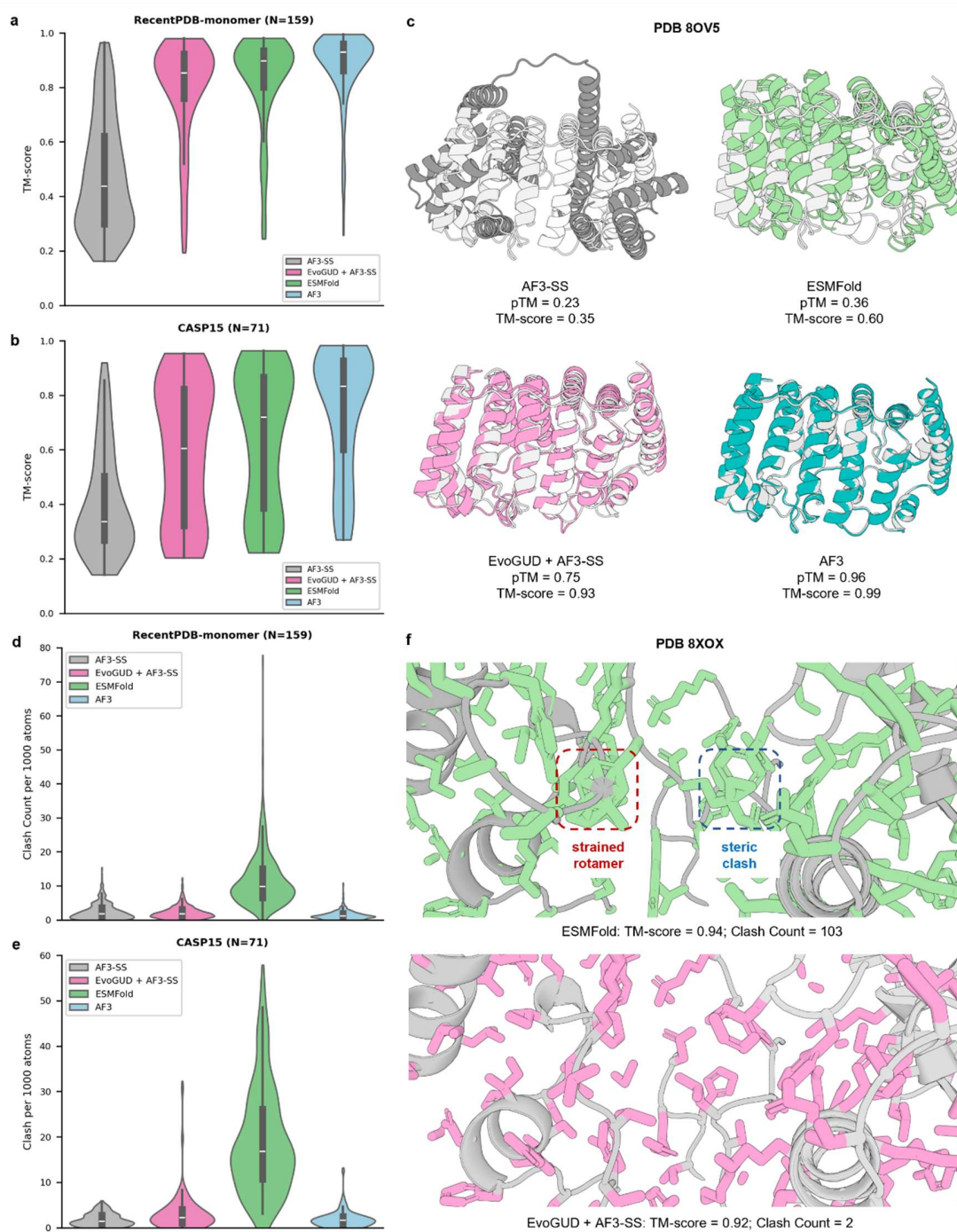
612 d, Distribution of sequence identity to the conditioning query as a function of  $\gamma$ ,  
613 demonstrating tunable control over sampling distance in sequence space. Additional  
614 validation metrics are shown in Supplementary Figures S1.

615 e, Evolutionary consistency beyond identity: cosine similarity of ESM-2 attention maps  
616 versus sequence identity for EvoGUD-generated sequences, compared with identity-  
617 matched random controls and natural MSA (nMSA) homologs, demonstrating  
618 preservation of natural-like co-evolutionary geometry across divergent regimes. EvoGUD  
619 data points correspond to 1,024 generated sequences per target across 159 RecentPDB-  
620 monomer proteins and six conditioning strengths ( $\gamma \in \{1, 2, 4, 8, 16, 32\}$ ), using the same  
621 generated sequences as in Fig. 1d and Supplementary Fig. S1b. Natural MSA sequences  
622 were included only when covering at least 80% of query positions to ensure comparable

623 alignment context. Solid lines denote mean trends and shaded regions indicate 95%

624 confidence intervals.

625



626

627 **Figure 2.** EvoGUD restores MSA-level monomer performance and improves all-atom  
628 quality from a single sequence.



629 a, TM-score distributions on the RecentPDB-monomer test set (N = 159), evaluated using  
630 the selected EvoGUD ensemble settings (Supplementary Fig. S5). Predictions from AF3-  
631 SS (single sequence), EvoGUD + AF3-SS, ESMFold, and AF3 with nMSAs are  
632 compared. EvoGUD substantially improves backbone accuracy relative to AF3-SS and  
633 approaches the performance of ESMFold and AF3.

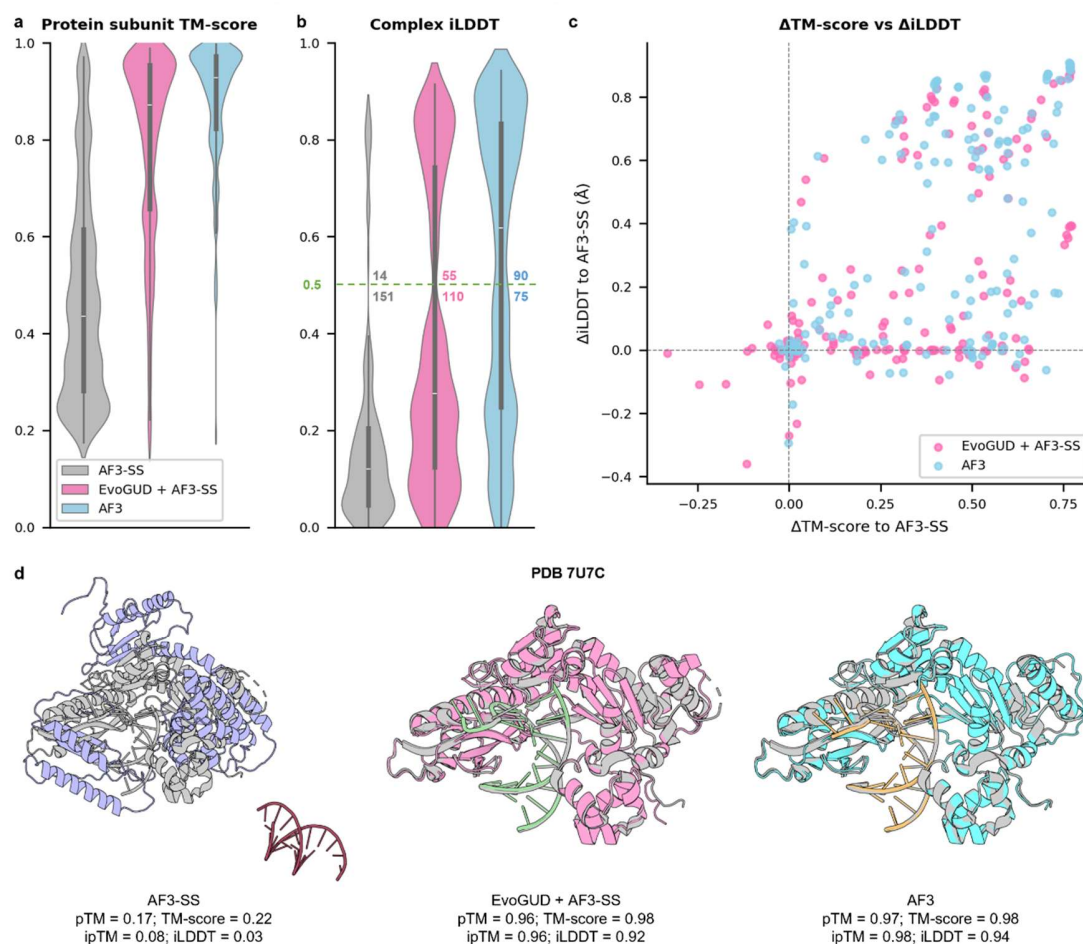
634 b, TM-score distributions on the CASP15 monomer benchmark (N = 71), evaluated using  
635 the same EvoGUD ensemble settings selected on RecentPDB-monomer. EvoGUD +  
636 AF3-SS generalizes to this independent benchmark without retuning.

637 c, Representative monomer example (PDB: 8OV5) illustrating global backbone accuracy.  
638 The experimental structure (white) is compared with predictions from AF3-SS (gray),  
639 ESMFold (green), EvoGUD + AF3-SS (pink), and AF3 (blue). EvoGUD + AF3-SS  
640 recovers the correct overall topology and domain arrangement.

641 d, All-atom steric clash counts per 1,000 atoms on the RecentPDB-monomer benchmark.  
642 Clash counts are computed as heavy-atom contacts closer than the sum of van der Waals  
643 radii with a 0.6 Å tolerance. EvoGUD + AF3-SS significantly reduces steric clashes  
644 relative to ESMFold.

645 e, All-atom steric clash counts on the CASP15 monomer benchmark. EvoGUD + AF3-SS  
646 maintains low clash rates comparable to AF3, indicating improved side-chain packing  
647 without nMSAs search.

648 f, Local side-chain environment example (PDB: 8XOX). ESMFold (green sticks, top)  
649 shows strained rotamers and steric clashes despite a high-quality backbone, whereas  
650 EvoGUD + AF3-SS (pink sticks, bottom) produces well-packed, stereochemically  
651 reasonable side chains.



652

653 **Figure 3.** EvoGUD enables single-sequence AlphaFold3 to model protein–nucleic-acid  
654 complexes.

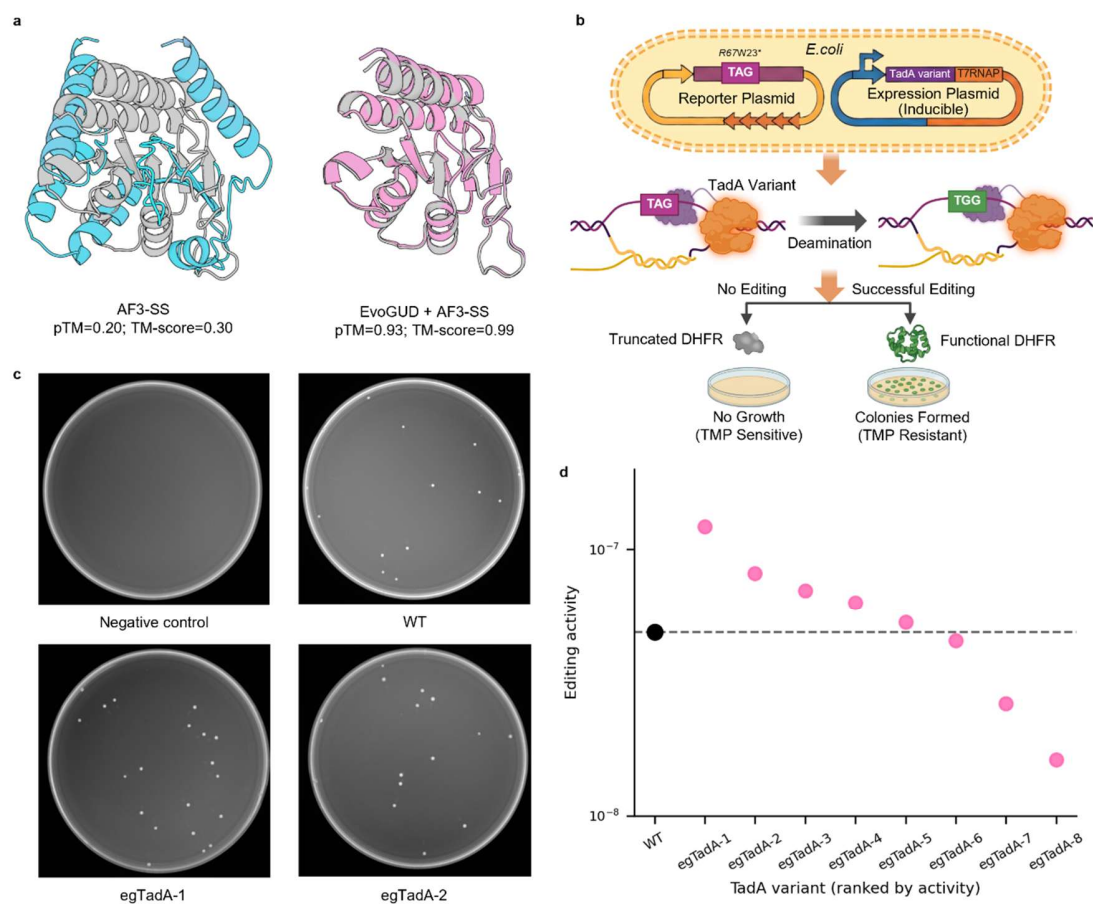
655 a, TM-score distributions of protein subunits extracted from predicted complexes for a  
656 benchmark of 165 protein–DNA/RNA assemblies, comparing AF3-SS (single sequence),  
657 EvoGUD + AF3-SS (vMSAs generated from single sequences), and AF3 with nMSAs  
658 and PDB template search (“AF3”). TM-scores are computed on protein subunits only,  
659 quantifying the correctness of individual protein folds within the predicted complexes.  
660 b, Interface LDDT (iLDDT) distributions for the same complexes, computed over  
661 protein–protein and protein–nucleic-acid interfaces following the AlphaFold3 evaluation

662 protocol. The horizontal dashed line at iLDDT = 0.5 is shown as a visual reference to  
663 illustrate the separation between lower- and higher-quality interface predictions observed  
664 in this benchmark. Numbers above and below the line report the counts of complexes on  
665 either side of this reference.

666 c, Per-complex changes in protein-subunit TM-score and interface iLDDT relative to  
667 AF3-SS. Each point corresponds to one complex, with  $\Delta$ TM-score on the x-axis and  
668  $\Delta$ iLDDT on the y-axis (positive values indicate improved interfaces).

669 d, Representative protein–nucleic-acid complex example (PDB:7U7C).

670



671

672 **Figure 4.** EvoGUD generates remotely homologous yet functional TadA variants.

673 a, AlphaFold3 (AF3) structure predictions for wild type *Staphylococcus aureus* TadA  
674 with EvoGUD-generated variants as vMSA, compared with AF3 predictions obtained  
675 from single-sequence input without MSA (AF3-SS). Aligned on reference PDB (2B3J)  
676 structure (gray).

677 b, Schematic of the trimethoprim (TMP) resistance reversion assay used to evaluate

678 TadA DNA-editing activity in *E. coli*. TadA-mediated A•T→G•C editing reverts a

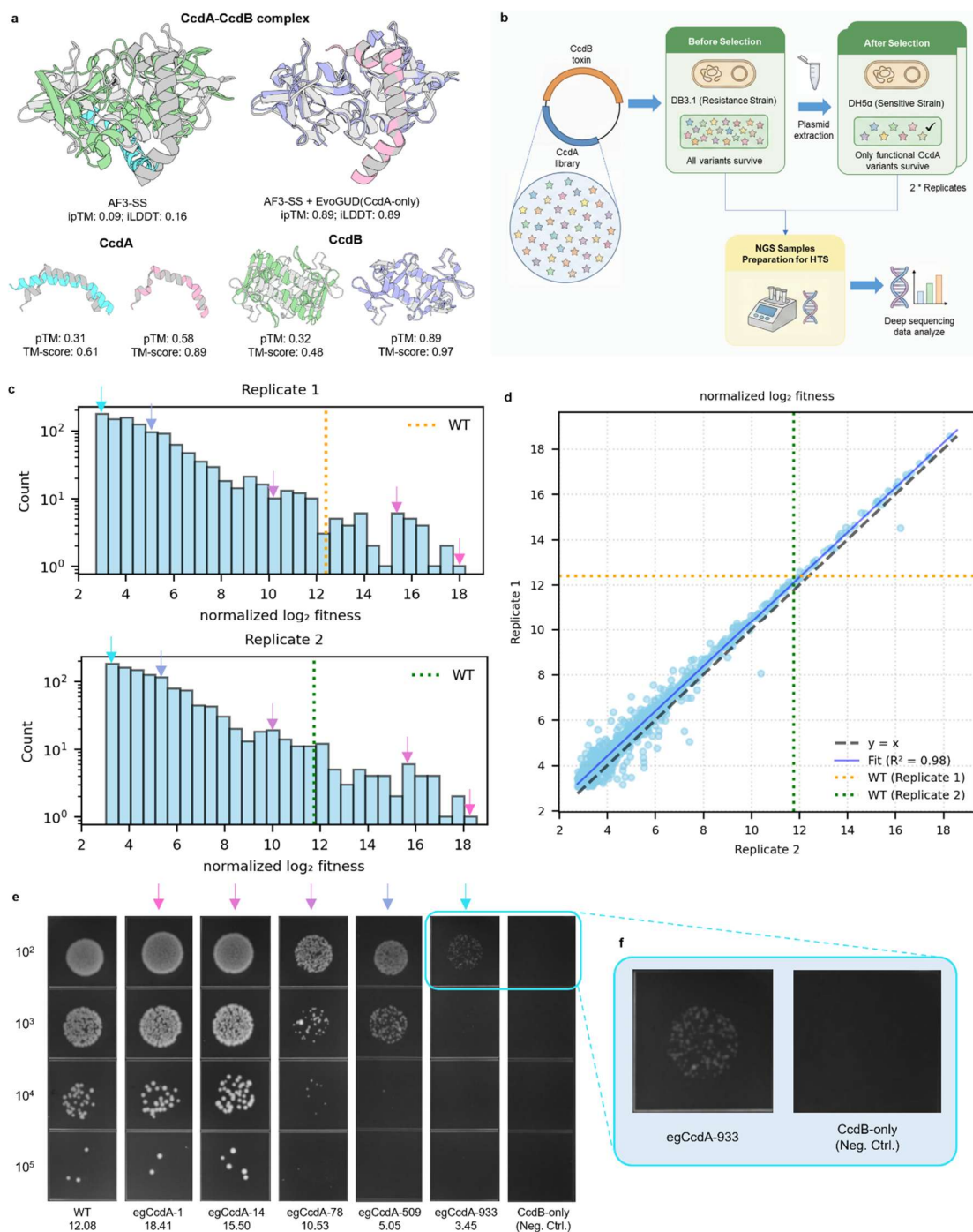
679 premature TAG stop codon in an *R67* dihydrofolate reductase (DHFR) reporter gene,

680 restoring the functional TGG codon and conferring TMP resistance.

681 c, Representative agar plate images showing TMP-resistant colony growth. A reporter-  
682 only strain lacking TadA expression serves as the negative control, wild type (WT) TadA  
683 is shown as a positive control, and two EvoGUD-generated variants (egTadA-1 and  
684 egTadA-2) illustrate functional recovery.

685 d, Quantification of DNA-editing activity for EvoGUD-generated TadA variants  
686 measured by TMP-resistance reversion. Each point represents one variant; WT is shown  
687 in black and EvoGUD-generated variants in pink.

688



689

690 **Figure 5.** Sequence-only conditional generation of the intrinsically disordered antitoxin

691 CcdA and high-throughput functional screening.

692 a, AlphaFold3 (AF3) predictions of the CcdA–CcdB complex for EvoGUD-generated  
693 CcdA variants. Left: AF3 single-sequence inference (AF3-SS) using the wild-type CcdA  
694 sequence. Right: AF3-SS using EvoGUD-generated CcdA sequences provided as vMSA  
695 ( $\gamma = 2$ , #MSA = 8). CcdA is shown in cyan (AF3-SS) or pink (AF3-SS + EvoGUD), and  
696 CcdB is shown in green (AF3-SS) or blue (AF3-SS + EvoGUD). Below, predicted  
697 structures of the CcdA and CcdB subunits are shown separately using the same color  
698 scheme. Reported ipTM, iLDDT, pTM, and TM-score values are indicated beneath each  
699 model.

700 b, Pooled selection workflow. A pooled oligo library of CcdA variants was constructed  
701 by cloning into a ccdB expression vector, pUC57-Kan-2BspQI-ccdB. The library was  
702 first propagated in the CcdB-resistant strain DB3.1 (Before selection), then subjected to  
703 selection in the CcdB-sensitive strain DH5 $\alpha$  (After selection; two biological replicates).  
704 Plasmids were extracted and deep-sequenced to infer variant fitness.

705 c, Distributions of normalized  $\log_2$  fitness for two biological replicates. Fitness is  
706 computed from before/after sequencing using early-stop variants as negative controls to  
707 estimate a baseline distribution (median and MAD), followed by BH-FDR filtering ( $q \leq$   
708 0.01) and normalization to obtain normalized  $\log_2$  fitness. Dotted lines mark WT; arrows  
709 indicate variants chosen for plate assays.

710 d, Cross-validation of normalized  $\log_2$  fitness between replicates ( $N = 1072$ ), with  $y = x$   
711 reference, fitted trend, and WT reference lines.

712 e, Plate-based validation of selected variants. Serial dilution spot assays for WT, selected  
713 variants, and a CcdB-only negative control; numbers under each label denote normalized  
714  $\log_2$  fitness used for selection.

715 f, Zoomed comparison highlighting egCcdA-933 versus the CcdB-only negative control.

716



717 **Data availability**

718 The source code, sampling script, and model weight are soon publicly available at  
719 <https://github.com/EricZhangSCUT/EvoGUD>.

720

721 **Author contributions**

722 XZ developed the computational methods, performed computational evaluation, and  
723 managed the overall project. JT and TZ designed experiments and performed  
724 experimental validations. ZC and ZZ performed NGS data processing together with XZ.  
725 JZ and YZ initiated and supervised the project, and YZ provided funding support. YZ and  
726 XZ drafted the initial manuscript. All authors contributed to manuscript revision and  
727 approved the final version.

728

729 **Acknowledgements**

730 This work was supported by the National Natural Science Foundation of China (Grant  
731 No. 92370202) and the National Key R&D Program of China (Grant No.  
732 2021YFF1200400). We acknowledge the High-Performance Computing Cluster at  
733 Shenzhen Bay Laboratory (SZBL) and the high-performance computing resources of the  
734 Shenzhen Medical Academy of Research and Translation (SMART) for providing  
735 computational support.

736

737 **Conflict of Interest**

738 All authors declare no financial interest. Jian Zhan is the founder and CEO of Ribopeptic,  
739 and Yaoqi Zhou is the scientific founder of Ribopeptic.