

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

4 Xing Zhang^{1§}, Jinle Tang^{1§}, Tingkai Zhang^{1,2§}, Zhihang Chen^{3,4}, Zhe Zhang¹, Jian
5 Zhan^{1,5,6*}, and Yaoqi Zhou^{1*}

6

7 ¹Institute of Systems and Physical Biology, Shenzhen Bay Laboratory, Shenzhen,
8 518132, China

9 ²School of Medicine, Southern University of Science and Technology, Shenzhen,
10 518055, China

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

13 ⁴Tsinghua University, Beijing 100084, China

14 ⁵Ribopeutic (Shenzhen) Co., Ltd., Shenzhen, 518000, China

15 ⁶Ribopeutic Inc., Hangzhou, 310018, China

16 [§]co-first authors. These authors contributed equally.

17 ^{*}Corresponding authors: Yaoqi Zhou, +86-(755) 2684 6275, zhouyq@szbl.ac.cn; Jian

18 Zhan +86-(755) 2684 6275, 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 with an 80% success
33 rate in targeted validation experiments, and 1,072 functional variants of the intrinsically
34 disordered antitoxin CcdA from 5,623 unique sequences in high-throughput selection
35 assays. 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^{1,2}.
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 space³. 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^{4,5}. Sequence-based protein language models provide
56 an alternative by leveraging large-scale evolutionary data to generate functional
57 sequences⁶; 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 distribution⁷⁻⁹.

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^{10,11}. 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^{12,13}. 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 (γ) 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 AlphaFold3¹⁴ 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)¹⁵ and large numbers of functional CcdA antitoxin variants (an
94 intrinsically disordered protein interacted with the toxic protein CcdB)¹⁶ 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 models¹⁷. 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 sequence¹. 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 patterns¹⁷.

110 Concretely, EvoGUD was trained as a denoising diffusion model^{10,18} 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 γ ,
116 modulates the extent to which sampling remains anchored to the conditioning query
117 sequence. As a result, γ 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 γ 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 monomers¹⁹ 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 γ values, EvoGUD produces smoothly tunable identity distributions, with increasing γ
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 γ 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 AlphaFold3¹⁴ under single-sequence
154 conditions. To this end, we used EvoGUD to generate candidate vMSAs across a range of
155 conditioning strengths γ and vMSA depths, and then applied a genetic algorithm to
156 adaptively select a target-specific subset of EvoGUD-generated sequences for
157 AlphaFold3 single-sequence inference (AF3-SS) (see Methods).

158 EvoGUD + AF3-SS substantially outperformed AF3-SS on the RecentPDB-monomer
159 test set (Fig. 2a). The mean TM-score²⁰ increased from 0.476 for AF3-SS to 0.820 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 ≥ 0.5 —a
162 commonly used criterion for correct fold assignment—increased from 40.9% for AF3-SS
163 to 91.8% 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 CASP15²¹ monomer benchmark (N = 71). EvoGUD + AF3-SS again markedly improved

168 backbone accuracy relative to AF3-SS and yielded TM-score distributions comparable to
169 those observed on RecentPDB-monomer (Fig. 2b). On CASP15, the mean TM-score
170 increased from 0.403 for AF3-SS to 0.650 for EvoGUD + AF3-SS, slightly exceeding the
171 performance of ESMFold (0.640) and approaching AF3 with nMSAs and PDB templates
172 (0.740). The fraction of targets with TM-score ≥ 0.5 increased from 25.4% for AF3-SS to
173 69.0% for EvoGUD + AF3-SS, again slightly exceeding ESMFold (66.2%) and
174 approaching AF3 with nMSAs and PDB structure templates (77.5%). A representative
175 example is shown in Fig. 2c, where EvoGUD + AF3-SS recovers the correct overall fold
176 and domain arrangement that is missed by AF3-SS.

177 Because AlphaFold3 is an all-atom diffusion model, we further evaluated the physical
178 plausibility of the predicted structures by quantifying steric clashes between heavy atoms.
179 Despite its strong backbone accuracy, ESMFold exhibited substantially higher clash
180 counts, with a broad tail of severely mispicked models on both RecentPDB-monomer
181 and CASP15 (Fig. 2d and 2e). By contrast, EvoGUD + AF3-SS, while achieving
182 backbone accuracy comparable to ESMFold, markedly reduced steric clashes relative to
183 ESMFold and producing distributions much closer to those of AF3. A representative
184 local side-chain environment is shown in Fig. 2f, where EvoGUD + AF3-SS yields well-
185 packed, stereochemically reasonable side chains, whereas ESMFold displays strained
186 rotamers and steric overlaps²².

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

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

192

193 **Protein-NA complex prediction with vMSA from EvoGUD**

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

206 For complex structures evaluated by interface local distance difference test (iLDDT²³) (as
207 in AF3¹⁴), 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)²⁴, 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 conversion¹⁵. Although highly

235 effective, previously reported TadA variants remain closely related to their ancestral
236 sequences, motivating exploration of more distant sequence solutions^{25,26}.
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*²⁷. 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^{16,28}.

263 Monomeric CcdA and CcdB sequences were derived from the *E. coli* CcdA–CcdB
264 complex (PDB: 3HPW)¹⁶. These sequences were then assembled into a symmetrized
265 CcdB–G₅₀–CcdA–G₅₀–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- γ 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)²⁹ 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 γ provides
328 continuous control over the balance between exploratory breadth and evolutionary
329 constraint. Together, these results suggest that single protein sequences contain richer
330 exploitable evolutionary structure than is apparent from sequence identity alone, and that
331 this structure can be computationally read out to support outward exploration in sequence
332 space.

333 These properties translate directly into improved structure prediction. By assembling
334 EvoGUD-generated sequences as vMSAs, AlphaFold3 single-sequence inference
335 recovers much of the accuracy and atomic detail typically associated with natural
336 homolog searches. Across monomer benchmarks, EvoGUD-assisted predictions are near
337 standard AlphaFold3 performance. For protein–DNA and protein–RNA complexes,
338 EvoGUD-derived vMSAs further improve interface geometry, demonstrating that the
339 generated sequence families encode actionable co-evolutionary signals for multimeric
340 recognition. Rather than forcing downstream models to consume opaque latent PLM
341 features, EvoGUD exposes PLM-derived evolutionary knowledge in the explicit

342 language of MSAs—a biologically meaningful and directly reusable interface for existing
343 structure predictors and other downstream tools.

344 Besides structure modeling, EvoGUD provides a general framework for functional
345 protein discovery in remote sequence space. For the TadA enzyme, EvoGUD identified
346 highly divergent yet functional DNA-editing variants, revealing functional solutions
347 inaccessible to stepwise directed evolution. Notably, these active variants reside at
348 approximately 40% sequence identity to the query, corresponding—in natural
349 evolutionary terms—to divergence accumulated over on the order of $\sim 10^9$ years
350 (Supplementary Table S4)³⁰. Similarly, in the CcdA–CcdB toxin–antitoxin system,
351 EvoGUD enabled large-scale discovery of functional antitoxin variants spanning 54–73%
352 sequence identity, even among low-ranked candidates, indicating robust preservation of
353 context-dependent functional constraints without explicit structural or biophysical
354 scoring. These results show that EvoGUD can jump directly from a single functional seed
355 sequence to experimentally testable remote variants in a single generative step, without
356 iterative computational triage. In this sense, EvoGUD complements rather than replaces
357 directed evolution: it expands the accessible parent space globally, whereas directed
358 evolution can subsequently refine promising regions locally.

359 A key requirement for sequence generative models is generalizability across protein
360 families and evolutionary distances. EvoGUD exhibits robust generalization across
361 multiple validation regimes, including stringent non-redundant subsets, with generated
362 sequences consistently following the same representation-space trajectories as natural
363 homologs irrespective of training-set proximity (Supplementary Fig. S8). Controlled
364 experiments on TadA further show that excluding or retaining close homologs during

365 training produces only minor shifts in generative behavior without evidence of collapse
366 (Supplementary Fig. S9). Notably, the CcdA system provides complementary insight:
367 although CcdA homologs were present in the training data, CcdA-only conditioning
368 yielded low identity and reduced foldability, whereas introducing an unseen fusion
369 context with its cognate binding partner systematically shifted generation toward
370 functionally coherent sequence space (Supplementary Fig. S10). Together, these results
371 indicate that EvoGUD does not rely on memorization of training sequences but is
372 primarily shaped by the evolutionary constraints supplied at inference time. This behavior
373 is important for broad use, because it suggests that EvoGUD is not restricted to
374 memorized families but instead provides a reusable mechanism for extracting structural
375 and functional priors from new single-sequence inputs.

376 EvoGUD adopts a modular alternative to end-to-end single-sequence structure prediction
377 pipelines. By decoupling protein sequence feature extraction from structure inference
378 with intermediate sequence generation, improvements in protein language models or
379 structure predictors can be incorporated by retraining only a lightweight generative
380 adapter, requiring on the order of a single GPU-day. This contrasts sharply with the
381 hundreds of GPU-weeks typically required to train or adapt full-scale structure prediction
382 models¹⁷, enabling rapid iteration and reuse of advances in representation learning. More
383 broadly, EvoGUD points to a general design principle for protein AI: explicit
384 evolutionary representations may provide a more interpretable and transferable interface
385 between foundation sequence models and downstream biological tasks than hidden
386 feature spaces alone. In this view, single-sequence expansion could serve as a general
387 foundation for biotechnology—linking PLMs to downstream models in computational

388 biology, while providing protein engineering with globally expanded starting points for

389 subsequent local optimization.

390

391 **Methods**

392 **Training Dataset and Data Preprocessing**

393 EvoGUD was trained on OpenProteinSet-PDB³¹, a curated reconstruction of the
394 AlphaFold2 training dataset as implemented in OpenFold³², comprising 131,487 protein
395 chains with precomputed MSAs. To ensure sequence integrity and computational
396 consistency, chains containing unknown amino acids (“X”) were excluded. Sequences
397 were further filtered by length, retaining chains with $30 \leq L \leq 1000$ amino acids. After
398 filtering, the final dataset comprised 117,556 entries, which were partitioned into a
399 training set of 116,756 entries and a validation set of 800 entries.

400

401 **Model architecture and conditioning**

402 EvoGUD is built on a Diffusion Transformer (DiT) backbone³³ with a modified adaLN-
403 Zero conditioning mechanism (Supplementary Fig. S11). Whereas the original DiT
404 conditions on global features, EvoGUD incorporates sequence-specific evolutionary
405 context derived from a protein language model.
406 During both training and inference, EvoGUD conditions the denoising process on
407 representations extracted from the ESM-2 3B model, including per-residue embeddings
408 (2,560 dimensions), and pairwise attention maps (36 layers \times 40 heads; 1,440
409 dimensions)¹⁷. These features are linearly projected into a 128-dimensional latent space.
410 In each DiT block, ESM-2 embeddings are injected via cross-attention as key–value
411 pairs, while projected attention maps are added to the attention logits as a pairwise bias.
412 The resulting representations generate the six modulation parameters (shift, scale, and

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

420

421 **Diffusion formulation**

422 EvoGUD adopts a continuous-time Gaussian diffusion framework with a cosine noise
423 schedule³⁴. The forward process gradually perturbs one-hot encoded amino-acid
424 sequences with Gaussian noise according to a cumulative signal retention coefficient $\bar{\alpha}_t$,
425 defined over normalized time $u \in [0,1]$ as:

$$426 \quad \bar{\alpha}_u = \frac{f(u)}{f(0)}, \quad f(u) = \cos^2 \left(\frac{u}{s} + 1 \cdot \frac{\pi}{2} \right)$$

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

431

432 **Model Training**

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

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

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

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

446 Training was performed for 100 epochs using the AdamW optimizer³⁵ with a learning
447 rate of 1×10^{-3} and automatic mixed precision. Each epoch comprised 10,000 unique
448 queries sampled from the whole training set, totaling $\sim 10^6$ optimization steps. Training
449 required approximately 25 h on a single NVIDIA A100 GPU and was implemented in
450 PyTorch 2.2³⁶.

451

452 **Single-sequence conditional sampling via probability absorption**

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

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

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

461 The sampled sequence is projected back into latent space as a centered one-hot
462 representation, scaled by a conditioning strength parameter γ :

$$463 \quad \hat{x}_0 = \gamma \cdot (\text{one_hot}(S^0) - 0.5)$$

464 The latent state is then updated via the Gaussian reverse transition, blending the absorbed
465 identity signal with stochastic noise. The parameter γ controls the trade-off between
466 evolutionary adherence and exploratory breadth: higher values promote conservative,
467 high-confidence homologs, whereas lower values allow broader exploration of remote
468 sequence space.

469

470 **Benchmark test sets**

471 **RecentPDB-monomer.** To evaluate structure prediction performance, an independent
472 test set was curated from PDB entries released between January 1 and July 1, 2024,
473 following AlphaFold3 benchmarking principles¹⁴. Only protein-only monomers solved by
474 X-ray crystallography at ≤ 2.0 Å resolution were retained. Chains were filtered to lengths
475 of 30–500 residues and subjected to 40% sequence-identity filtering both within the set
476 and against the training/validation data. The final set comprised 159 non-redundant
477 monomers.

478 **CASP15.** To assess generalization beyond RecentPDB-monomer, we evaluated structure
479 prediction performance on the CASP15 monomer benchmark, an independent community
480 test set of experimentally determined protein structures. The 71 monomeric protein

481 targets in this set were used as an external benchmark to evaluate the robustness and
482 generalizability of EvoGUD-assisted structure prediction.

483 **RecentPDB-multimer.** Protein–nucleic-acid complexes were derived from the
484 AlphaFold3 benchmark dataset¹⁴. After excluding entries lacking protein chains or
485 containing non-canonical residues, the final set comprised 165 complexes. vMSAs were
486 generated only for protein components, while nucleic-acid sequences were held fixed.

487

488 **Virtual MSA settings for EvoGUD + AF3-SS**

489 A grid search over conditioning strength γ and vMSA depth on the RecentPDB-monomer
490 benchmark identified a broad but non-monotonic regime in which EvoGUD-derived
491 vMSAs substantially improved structure prediction accuracy relative to AF3-SS
492 (Supplementary Fig. S5). Within this regime, appropriate combinations of γ
493 and vMSA depth produced marked gains in both predicted confidence (pTM¹⁴) and
494 structural accuracy (TM-score), indicating that EvoGUD-generated sequence sets can
495 recover much of the evolutionary information normally supplied by natural homologs.
496 However, the benefit depended on the combination of these two parameters: shallow
497 vMSAs were often insufficient, whereas overly deep or overly strongly conditioned
498 vMSAs could reduce performance on some targets, indicating that the optimal virtual
499 alignment regime is target dependent.

500 We therefore implemented a genetic algorithm to adaptively select a target-specific
501 subset of EvoGUD-generated sequences for AlphaFold3 inference. Candidate sequences
502 generated at $\gamma = 1, 2,$ and 4 (1,024 sequences per setting) were pooled, and the algorithm
503 iteratively searched for improved vMSA subsets using AlphaFold3 confidence as the

504 optimization signal (see Supplementary Methods for details). This GA-based procedure
505 was used as the default EvoGUD+AF3-SS setting for monomer prediction and was
506 applied to the RecentPDB-monomer and CASP15 benchmarks.

507 For applications requiring lower computational cost, we also implemented a fast
508 inference mode based on a small fixed ensemble of 10 EvoGUD parameter settings ($\gamma = 1$
509 with $\#vMSA \in \{2,4,8,16\}$ and $\gamma = 2$ with $\#vMSA \in \{2,4,8,16,32,64\}$) without iterative
510 AlphaFold3 optimization. This lightweight mode recovered most of the performance gain
511 of EvoGUD+AF3-SS at substantially lower cost (Supplementary Fig. S12), and was
512 therefore used for protein–nucleic-acid complex prediction.

513

514 **Generation and validation of TadA variants**

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

521 TadA activity was quantified using a trimethoprim resistance reversion assay in *E. coli*
522 based on a premature stop-codon reporter. Editing activity was defined as per-base
523 mutation rates estimated from the observed frequency of TMP-resistant colonies, using a
524 Luria–Delbrück approximation³⁷. Full experimental protocols are provided in the
525 Supplementary Information.

526

527 **Generation and validation of CcdA variants**

528 CcdA variants were generated using a conditional single-chain strategy in which the
529 antitoxin sequence was embedded within a symmetrized CcdB–G₅₀–CcdA–G₅₀–CcdB
530 fusion, with wild type CcdA and CcdB sequences from *E. coli* (PDB: 3HPW)¹⁶. Here, G₅₀
531 denotes a 50-amino-acid poly-glycine linker that spatially separates the CcdA and CcdB
532 domains while preserving sequence-level context³⁸. Residues outside the CcdA region
533 were fixed during sampling. Conditioning strength $\gamma = 2$ was selected according to
534 predicted foldability, novelty, and diversity (Supplementary Fig. S4).
535 A pooled library of 10,000 CcdA variants was generated under a co-generation identity
536 constraint of $\leq 75\%$ relative to wild-type CcdA and subsequently clustered at 90%
537 sequence identity, yielding 5,623 unique sequences. Functional selection was performed
538 using a toxin-rescue assay in *E. coli*, followed by deep sequencing.
539 Variant fitness was estimated using a Poisson-based \log_2 enrichment model and false-
540 discovery-rate (FDR) correction²⁹. Variants passing $\text{FDR} \leq 0.01$ in two independent
541 experiments were considered functional. Full experimental protocols are provided in the
542 Supplementary Information.
543

544 **Reference**

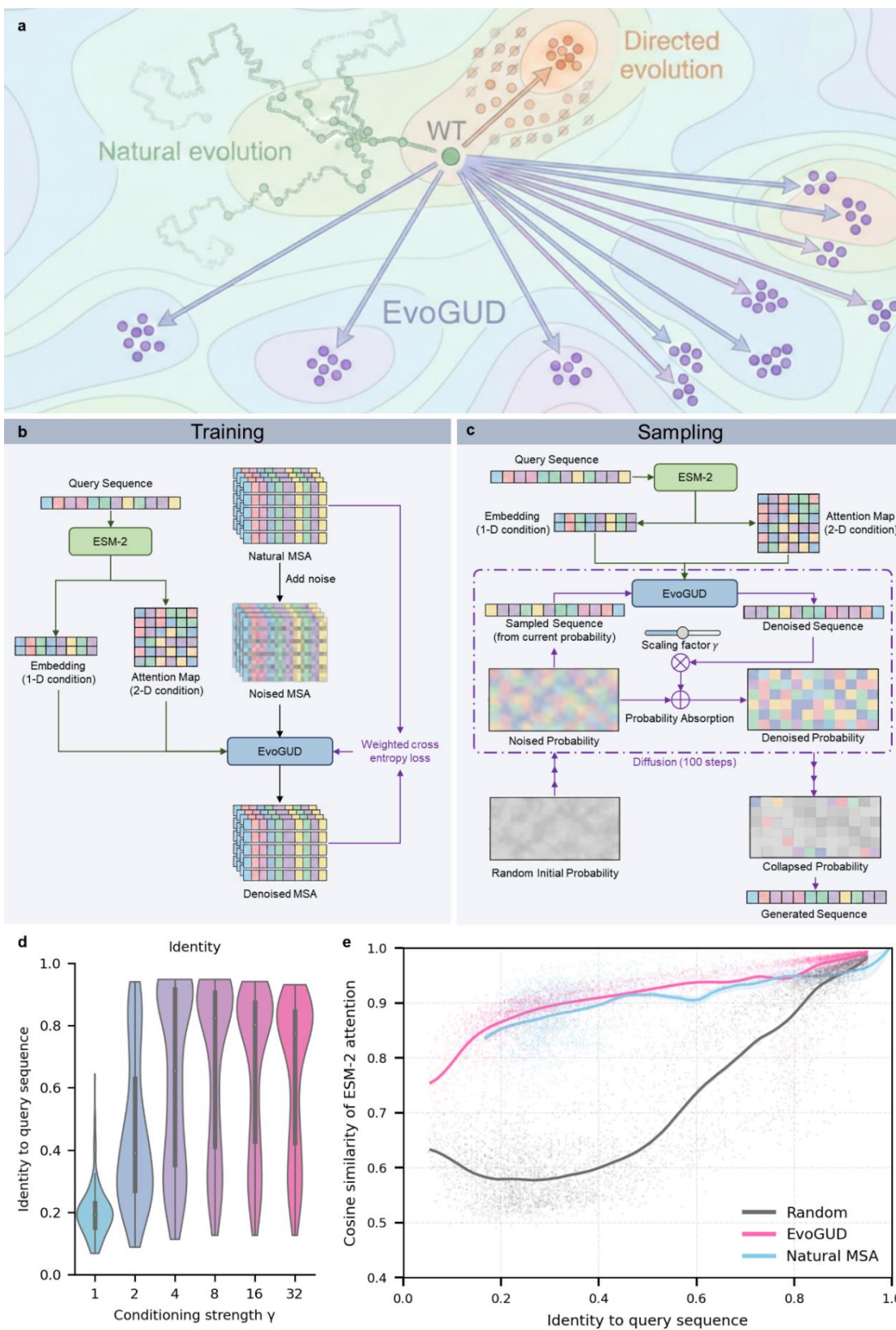
- 545 1. Arnold, F. H. Design by directed evolution. *Accounts of chemical research* **31**, 125–131
546 (1998).
- 547 2. Romero, P. A. & Arnold, F. H. Exploring protein fitness landscapes by directed
548 evolution. *Nature reviews Molecular cell biology* **10**, 866–876 (2009).
- 549 3. Poelwijk, F. J., Tănase-Nicola, S., Kiviet, D. J. & Tans, S. J. Reciprocal sign epistasis is a
550 necessary condition for multi-peaked fitness landscapes. *Journal of theoretical biology*
551 **272**, 141–144 (2011).
- 552 4. Huang, P.-S., Boyken, S. E. & Baker, D. The coming of age of de novo protein design.
553 *Nature* **537**, 320–327 (2016).
- 554 5. Dauparas, J. *et al.* Robust deep learning–based protein sequence design using
555 ProteinMPNN. *Science* **378**, 49–56 (2022).
- 556 6. Rives, A. *et al.* Biological structure and function emerge from scaling unsupervised
557 learning to 250 million protein sequences. *Proceedings of the National Academy of*
558 *Sciences* **118**, e2016239118 (2021).
- 559 7. Madani, A. *et al.* Large language models generate functional protein sequences across
560 diverse families. *Nature biotechnology* **41**, 1099–1106 (2023).
- 561 8. Nijkamp, E., Ruffolo, J. A., Weinstein, E. N., Naik, N. & Madani, A. Progen2: exploring
562 the boundaries of protein language models. *Cell systems* **14**, 968-978. e3 (2023).
- 563 9. Zhang, J. *et al.* Unsupervisedly prompting AlphaFold2 for accurate few-shot protein
564 structure prediction. *Journal of Chemical Theory and Computation* **19**, 8460–8471
565 (2023).

- 566 10. Ho, J., Jain, A. & Abbeel, P. Denoising diffusion probabilistic models. *Advances in*
567 *neural information processing systems* **33**, 6840–6851 (2020).
- 568 11. Trippe, B. L. *et al.* Diffusion probabilistic modeling of protein backbones in 3d for the
569 motif-scaffolding problem. *arXiv preprint arXiv:2206.04119* (2022).
- 570 12. Alamdari, S. *et al.* Protein generation with evolutionary diffusion: sequence is all you
571 need. *BioRxiv* 2023.09. 11.556673 (2023).
- 572 13. Sgarbossa, D., Lupo, U. & Bitbol, A.-F. Generative power of a protein language model
573 trained on multiple sequence alignments. *Elife* **12**, e79854 (2023).
- 574 14. Abramson, J. *et al.* Accurate structure prediction of biomolecular interactions with
575 AlphaFold 3. *Nature* **630**, 493–500 (2024).
- 576 15. Gaudelli, N. M. *et al.* Programmable base editing of A•T to G•C in genomic DNA
577 without DNA cleavage. *Nature* **551**, 464–471 (2017).
- 578 16. De Jonge, N. *et al.* Rejuvenation of CcdB-poisoned gyrase by an intrinsically
579 disordered protein domain. *Molecular cell* **35**, 154–163 (2009).
- 580 17. Lin, Z. *et al.* Evolutionary-scale prediction of atomic-level protein structure with a
581 language model. *Science* **379**, 1123–1130 (2023).
- 582 18. Austin, J., Johnson, D. D., Ho, J., Tarlow, D. & Van Den Berg, R. Structured
583 denoising diffusion models in discrete state-spaces. *Advances in neural information*
584 *processing systems* **34**, 17981–17993 (2021).
- 585 19. Berman, H. M. *et al.* The protein data bank. *Nucleic acids research* **28**, 235–242
586 (2000).

- 587 20. Zhang, Y. & Skolnick, J. Scoring function for automated assessment of protein
588 structure template quality. *Proteins: Structure, Function, and Bioinformatics* **57**, 702–710
589 (2004).
- 590 21. Kryshchak, A., Schwede, T., Topf, M., Fidelis, K. & Moult, J. Critical assessment
591 of methods of protein structure prediction (CASP)—Round XV. *Proteins: Structure,*
592 *Function, and Bioinformatics* **91**, 1539–1549 (2023).
- 593 22. Williams, C. J. *et al.* MolProbity: more and better reference data for improved all-atom
594 structure validation. *Protein Science* **27**, 293–315 (2018).
- 595 23. Mariani, V., Biasini, M., Barbato, A. & Schwede, T. IDDT: a local superposition-free
596 score for comparing protein structures and models using distance difference tests.
597 *Bioinformatics* **29**, 2722–2728 (2013).
- 598 24. Chang, C., Lee Luo, C. & Gao, Y. In crystallo observation of three metal ion promoted
599 DNA polymerase misincorporation. *Nature Communications* **13**, 2346 (2022).
- 600 25. Richter, M. F. *et al.* Phage-assisted evolution of an adenine base editor with improved
601 Cas domain compatibility and activity. *Nature biotechnology* **38**, 883–891 (2020).
- 602 26. Lapinaite, A. *et al.* DNA capture by a CRISPR-Cas9–guided adenine base editor.
603 *Science* **369**, 566–571 (2020).
- 604 27. Losey, H. C., Ruthenburg, A. J. & Verdine, G. L. Crystal structure of *Staphylococcus*
605 *aureus* tRNA adenosine deaminase TadA in complex with RNA. *Nature structural &*
606 *molecular biology* **13**, 153–159 (2006).
- 607 28. Loris, R. *et al.* Crystal structure of CcdB, a topoisomerase poison from *E. coli*. *Journal*
608 *of molecular biology* **285**, 1667–1677 (1999).

- 609 29. Benjamini, Y. & Hochberg, Y. Controlling the false discovery rate: a practical and
610 powerful approach to multiple testing. *Journal of the Royal statistical society: series B*
611 (*Methodological*) **57**, 289–300 (1995).
- 612 30. Kumar, S., Stecher, G., Suleski, M. & Hedges, S. B. TimeTree: a resource for
613 timelines, timetrees, and divergence times. *Molecular biology and evolution* **34**, 1812–
614 1819 (2017).
- 615 31. Ahdritz, G. *et al.* OpenProteinSet: Training data for structural biology at scale.
616 *Advances in Neural Information Processing Systems* **36**, 4597–4609 (2023).
- 617 32. Ahdritz, G. *et al.* OpenFold: Retraining AlphaFold2 yields new insights into its
618 learning mechanisms and capacity for generalization. *Nature methods* **21**, 1514–1524
619 (2024).
- 620 33. Peebles, W. & Xie, S. Scalable diffusion models with transformers. in *Proceedings of*
621 *the IEEE/CVF international conference on computer vision* 4195–4205 (2023).
- 622 34. Nichol, A. Q. & Dhariwal, P. Improved denoising diffusion probabilistic models. in
623 *International conference on machine learning* 8162–8171 (PMLR, 2021).
- 624 35. Loshchilov, I. & Hutter, F. Decoupled weight decay regularization. *arXiv preprint*
625 *arXiv:1711.05101* (2017).
- 626 36. Paszke, A. *et al.* Pytorch: An imperative style, high-performance deep learning library.
627 *Advances in neural information processing systems* **32**, (2019).
- 628 37. Luria, S. E. & Delbrück, M. Mutations of bacteria from virus sensitivity to virus
629 resistance. *Genetics* **28**, 491 (1943).
- 630 38. Chen, X., Zaro, J. L. & Shen, W.-C. Fusion protein linkers: property, design and
631 functionality. *Advanced drug delivery reviews* **65**, 1357–1369 (2013).

632 **Figures**



633

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

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

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

641 c, Sampling procedure: starting from a random probability distribution, EvoGUD
642 iteratively denoises sequence probabilities under query conditioning, with a scaling factor
643 γ controlling the strength of probability absorption before collapsing to a discrete
644 sequence.

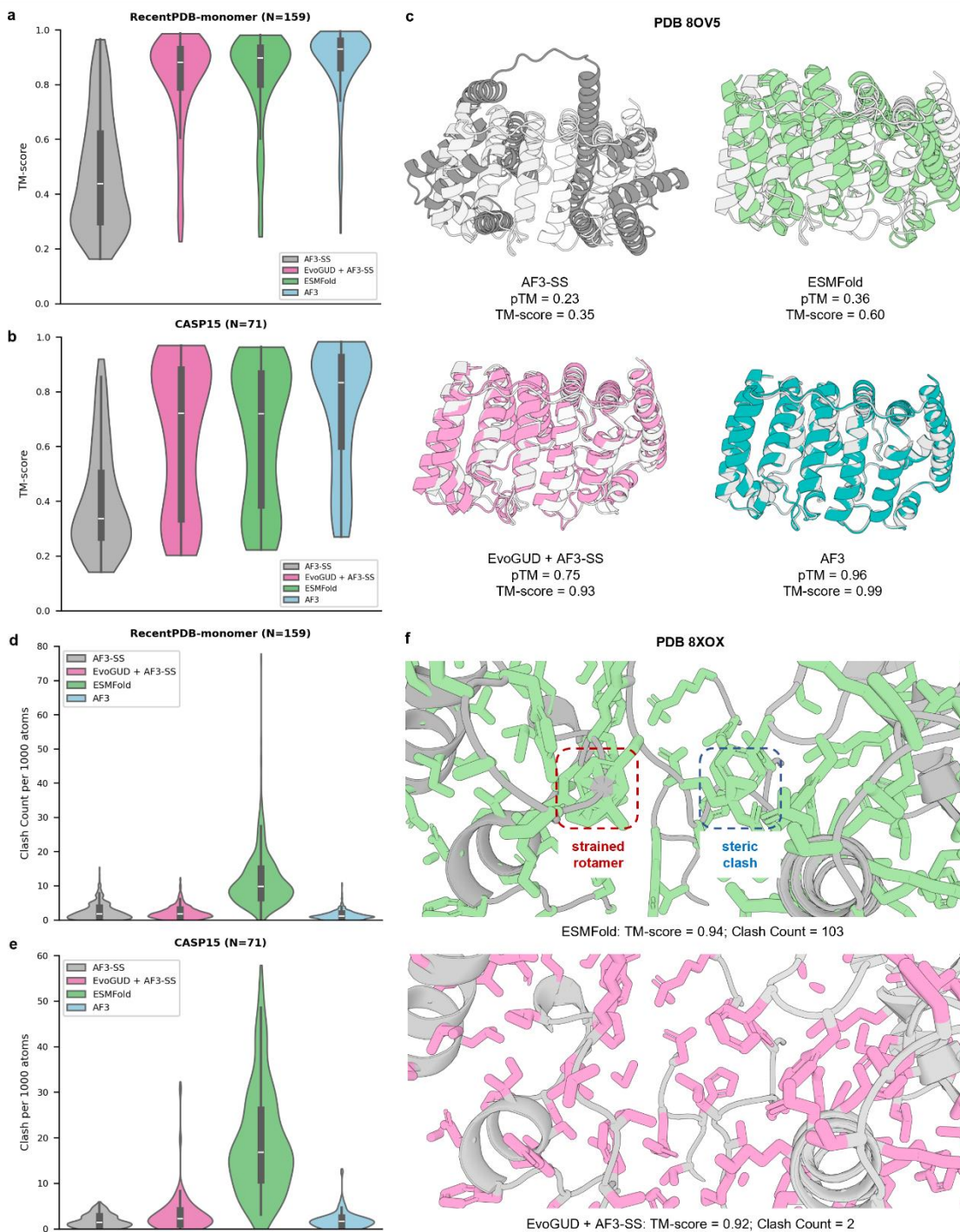
645 d, Distribution of sequence identity to the conditioning query as a function of γ ,
646 demonstrating tunable control over sampling distance in sequence space. Additional
647 validation metrics are shown in Supplementary Figures S1.

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

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

657 confidence intervals.

658



659

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

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

667 b, TM-score distributions on the CASP15 monomer benchmark (N = 71), evaluated using
668 the same EvoGUD ensemble settings selected on RecentPDB-monomer. EvoGUD +
669 AF3-SS again substantially improves backbone accuracy relative to AF3-SS,
670 demonstrating robust performance on an independent community benchmark.

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

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

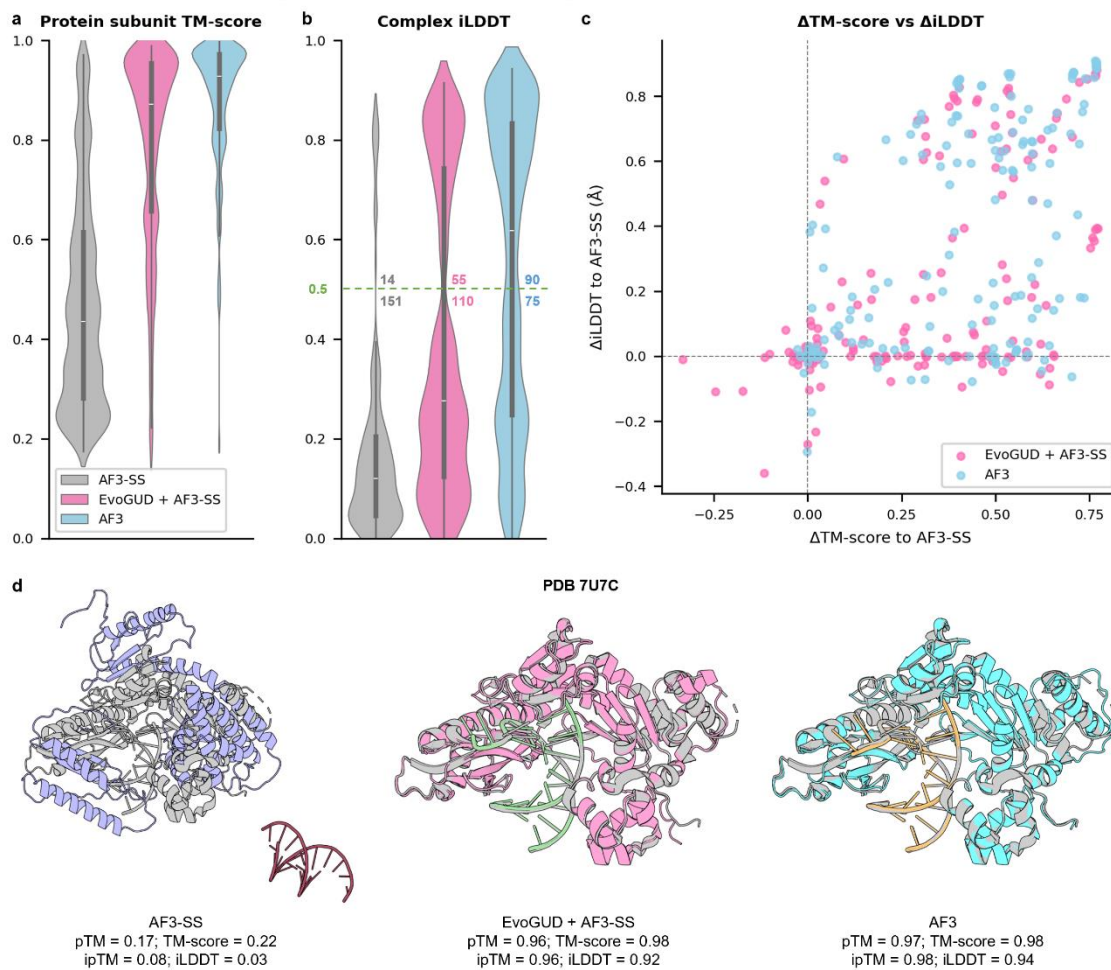
679 e, All-atom steric clash counts on the CASP15 monomer benchmark. EvoGUD + AF3-SS
680 maintains low clash rates comparable to AF3, indicating improved side-chain packing
681 without nMSAs search.

682 f, Local side-chain environment example (PDB: 8XOX). ESMFold (green sticks, top)
683 shows strained rotamers and steric clashes despite a high-quality backbone, whereas

684 EvoGUD + AF3-SS (pink sticks, bottom) produces well-packed, stereochemically

685 reasonable side chains.

686



687

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

690 a, TM-score distributions of protein subunits extracted from predicted complexes for a
 691 benchmark of 165 protein–DNA/RNA assemblies, comparing AF3-SS (single sequence),
 692 EvoGUD + AF3-SS (vMSAs generated from single sequences), and AF3 with nMSAs
 693 and PDB template search (“AF3”). TM-scores are computed on protein subunits only,
 694 quantifying the correctness of individual protein folds within the predicted complexes.

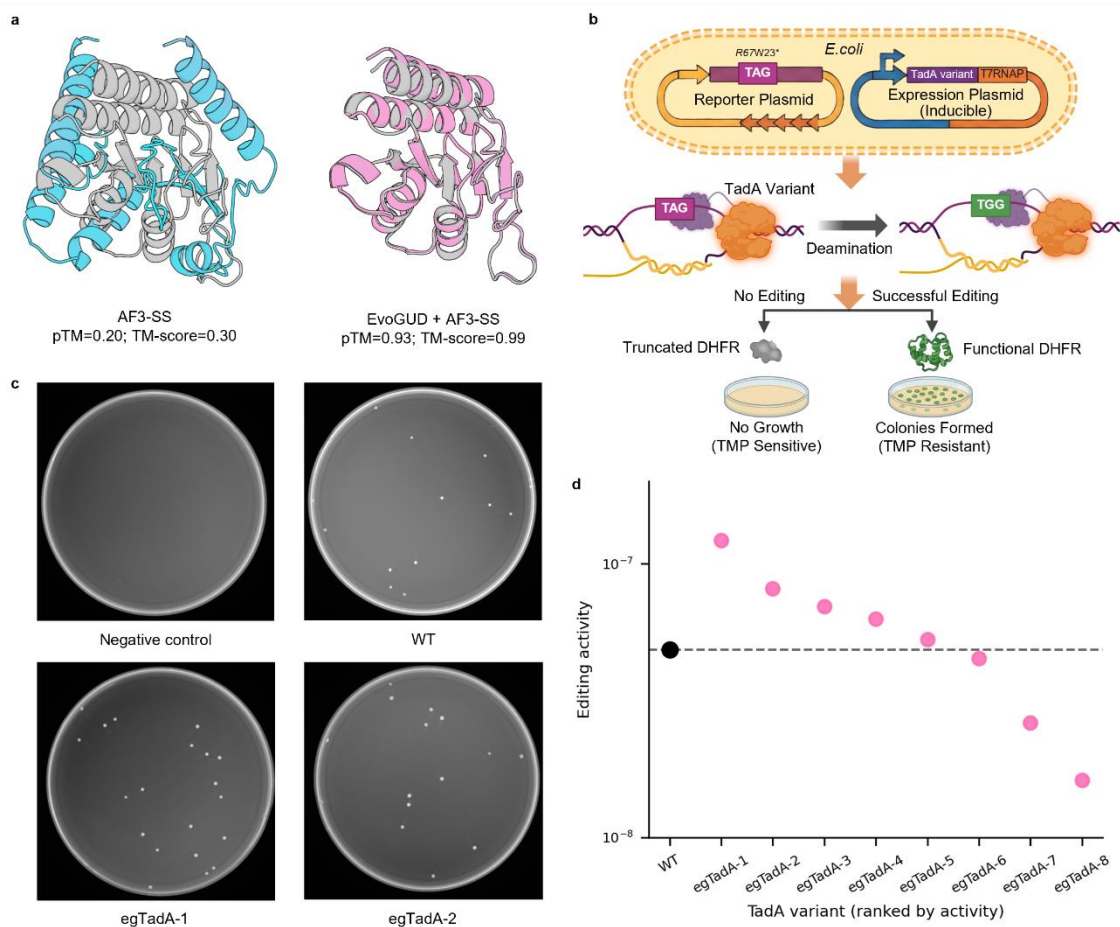
695 b, Interface LDDT (iLDDT) distributions for the same complexes, computed over
 696 protein–protein and protein–nucleic-acid interfaces following the AlphaFold3 evaluation

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

701 c, Per-complex changes in protein-subunit TM-score and interface $iLDDT$ relative to
702 AF3-SS. Each point corresponds to one complex, with ΔTM -score on the x-axis and
703 $\Delta iLDDT$ on the y-axis (positive values indicate improved interfaces).

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

705



706

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

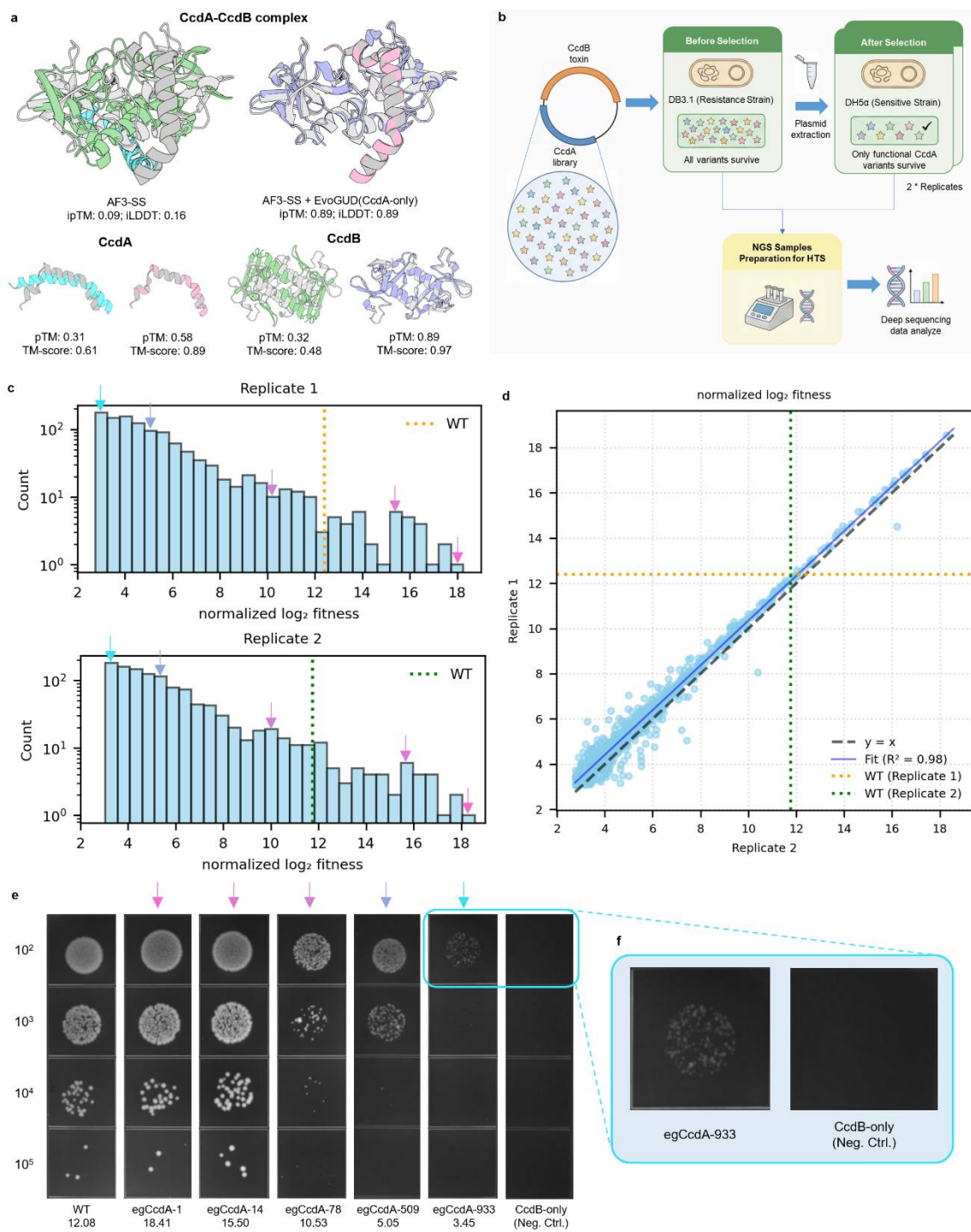
708 a, AlphaFold3 (AF3) structure predictions for wild type *Staphylococcus aureus* TadA
 709 with EvoGUD-generated variants as vMSA, compared with AF3 predictions obtained
 710 from single-sequence input without MSA (AF3-SS). Aligned on reference PDB (2B3J)
 711 structure (gray).

712 b, Schematic of the trimethoprim (TMP) resistance reversion assay used to evaluate
 713 TadA DNA-editing activity in *E. coli*. TadA-mediated A•T→G•C editing reverts a
 714 premature TAG stop codon in an R67 dihydrofolate reductase (DHFR) reporter gene,
 715 restoring the functional TGG codon and conferring TMP resistance.

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

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

723



724

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

726 CcdA and high-throughput functional screening.

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

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

740 c, Distributions of normalized \log_2 fitness for two biological replicates. Fitness is
741 computed from before/after sequencing using early-stop variants as negative controls to
742 estimate a baseline distribution (median and MAD), followed by BH-FDR filtering ($q \leq$
743 0.01) and normalization to obtain normalized \log_2 fitness. Dotted lines mark WT; arrows
744 indicate variants chosen for plate assays.

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

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

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

751

752 **Data availability**

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

755

756 **Author contributions**

757 XZ developed the computational methods, performed computational evaluation, and
758 managed the overall project. JT and TZ designed experiments and performed
759 experimental validations. ZC and ZZ performed NGS data processing together with XZ.
760 JZ and YZ initiated and supervised the project, and YZ provided funding support. YZ and
761 XZ drafted the initial manuscript. All authors contributed to manuscript revision and
762 approved the final version.

763

764 **Acknowledgements**

765 This work was supported by the National Natural Science Foundation of China (Grant
766 No. 92370202). We acknowledge the High-Performance Computing Cluster at Shenzhen
767 Bay Laboratory (SZBL) and the high-performance computing resources of the Shenzhen
768 Medical Academy of Research and Translation (SMART) for providing computational
769 support.

770

771 **Conflict of Interest**

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