# SPIN-dvEvo: Exploration of vast functional sequence space by directed virtual evolution from a local sequence cluster

Zhihang Chen[1,2,§], Jinle Tang[3,§], Tingkai Zhang[3,4,§], Xing Zhang[3], Qinghui Nie[2,3], Jian Zhan[3,5,6,*], and Yaoqi Zhou[1,3,*]

[1]Tsinghua University, Beijing 100084, China

[2]Shenzhen Medical Academy of Research and Translation (SMART), Shenzhen, 518107, China

[3]Institute of Systems and Physical Biology, Shenzhen Bay Laboratory, Shenzhen, 518107, China

[4]School of Medicine, Southern University of Science and Technology, Shenzhen, 518055, China

[5]Ribopeutic (Shenzhen) Co., Ltd., Futian, Shenzhen, 518000, China

[6]Ribopeutic Inc., Qiantang, Hangzhou, 310018, China

[§]Co-first authors. These authors are contributed to the manuscript equally.

*Corresponding authors: Yaoqi Zhou, +86-(755) 2684 6275, zhouyq@szbl.ac.cn; Jian Zhan, +86-(755) 2684 6275, zhanjian@szbl.ac.cn

## Abstract

Both natural and directed evolution are powerful in improving protein functions but they are slow in exploring the nearly endless sequence space. Here, we present SPIN-dvEvo that couples few-shot low-rank adaptation (LoRA) of an ESM-2 protein language model with a genetic algorithm to quickly evolve functional remote homologs from a local cluster of highly-homologous, binary-labeled sequences. We experimentally tested SPIN-dvEvo on an enzyme (the core deaminase component of adenine base editors, TadA) and an intrinsically disordered protein (antitoxin CcdA). In TadA, virtually evolved sequences with low sequence identity to the starting sequences achieved a 38% success rate (23/60) in the first round and a 51% success rate along with a one-order-of-magnitude improvement in enzymatic activity in the second round, for which SPIN-dvEvo was retrained on first-round labels. Virtual evolution of the disordered protein CcdA was also successful, albeit at low success rate of 2.6%. Thus, SPIN-dvEvo can simulate billions of years of evolution in just minutes, rapidly creating new functional clusters.

## Introduction

Directed evolution is a central strategy for engineering functional proteins, enabling stepwise improvement of enzymes and binders directly in the laboratory through iterative cycles of mutagenesis, selection, and amplification. It has produced catalysts with enhanced activity, altered specificity, and improved stability for applications ranging from therapeutics to industrial biocatalysis[1,2]. Despite these successes, both natural and laboratory evolution remain intrinsically

constrained by the locality of accessible mutational steps: most variants that can be efficiently sampled and screened in practice differ from their progenitors by only a small number of substitutions[3], and optimization therefore proceeds as a "local walk" on a rugged fitness landscape[3]. Because protein function is shaped by epistasis and higher-order constraints, such local searches can become trapped on suboptimal peaks, leaving distant but potentially functional better-fitness regions of sequence space systematically underexplored[3, 4]. While deep mutational scanning and deep sequencing expanded access to large-scale sequence–function measurements [5-7], they were also largely limited to local sequence space.

More recently, directed virtual evolution has emerged as a sequence-first alternative that learns surrogate fitness landscapes from sequence–function data and improves sequences by in silico search. Early work established probabilistic surrogate modeling, with Gaussian processes trained on measured activities to guide navigation of fitness landscapes.[8, 9] In parallel, structure-first pipelines such as AiCE integrate inverse-folding models with structural and evolutionary constraints to propose fold-compatible variants when reliable templates are available.[10] These approaches for directed virtual evolution can be grouped by how they obtain supervision and how far they can reliably search. One class relies on regression-style predictors trained on quantitative measurements (activity, fitness, binding, or other continuous phenotypes). These methods include deep supervised models or protein language models (PLM) with evolutionary context and active-learning pipelines for multi-round optimization such as ECNet[11] and EVOLVEpro[12], active-learning evolution of artificial metalloenzymes by Vornholt and colleagues[13], iterative deep-learning–guided directed evolution described by Li and colleagues[14], and Active Learning-assisted Directed Evolution (ALDE) by Yang et al.[15] A second class focuses on improving label efficiency, exemplified by Low-N protein engineering, which trains data-efficient predictors from small labeled sets and then screens large virtual libraries. [16] A third class formalizes "design–test–learn" iteration and can be coupled to automated platforms or biofoundries for higher-throughput cycles, exemplified by the STAR web server[17] and biofoundry-integrated PLM workflows for automated protein evolution.[18] Despite substantial gains in these systems, the learned surrogates continue to be most reliable near the starting scaffold. That is, the search remains implemented as local exploration (for example, small-step mutation moves, Bayesian optimization, or iterative screening) rather than sampling remote, low-identity functional sequences directly.

Another limitation of the above methods is their reliance on strong, quantitative supervision. Many successful workflows train regression-style surrogate landscapes on continuous activity measurements, kinetic readouts, or well-calibrated phenotypes, often requiring hundreds to thousands of assayed variants per round to achieve predictive accuracy that is sufficient to guide search.[12-14, 16, 19] When the available signal is weaker—binary functional labels or enrichment counts from pooled selections—performance can degrade because the training objective becomes less informative per variant and experimental noise from low counts and sampling variance becomes a dominant factor that must be modeled carefully.[6, 7, 20] In practice, this data requirement can be the primary bottleneck early in a campaign, when only a small number of positives exist and quantitative characterization is costly or not yet available.

Here we introduce SPIN-dvEvo (Sequence Prediction with Integrated Neural networks – directed

81   virtual Evolution), a directed virtual evolution framework that unifies discovery and multi-round
82   refinement in a single, sequence-first, low-data workflow. In contrast to conventional generative
83   models that necessitate vast functional datasets, SPIN-dvEvo maintains high precision in label-
84   scarce regimes by harnessing the evolutionary priors embedded within ESM-2[18] and fine-tuning its
85   trajectory through lightweight LoRA.[21] The resulting model was employed as a scorer for directed
86   virtual evolution powered by a genetic algorithm. SPIN-dvEvo was applied to virtually evolve TadA
87   adenosine deaminase activity and the intrinsically disordered antitoxin CcdA. In both systems, only
88   sparse binary functional labels were used to train a LoRA-based activity scorer on a frozen ESM-2
89   encoder, and candidate sequences were then generated by iterative mutation–crossover search
90   guided by the fixed scorer, with additional feedback round performed for TadA by updating the
91   LoRA head with newly obtained experimental labels. Experimental results confirmed the ability of
92   SPIN-dvEvo to quickly evolve from a local cluster of a few highly homologous sequences to remote
93   functional clusters.
94

## Results

### Directed virtual evolution by SPIN-dvEvo

97   The SPIN-dvEvo framework consists of two tightly coupled components: (i) a LoRA-adapted ESM-
98   2 scoring model that trained task-specific functional scores using a small set of positive and negative
99   sequences, and (ii) a genetic algorithm (GA) that directs sequence sampling and virtual evolution
100   under this learned landscape. Each run starts from an initial seed pool generated by applying 20%
101   random substitutions to the starting sequences. Here, SPIN-dvEvo was fine-tuned using only
102   qualitatively labeled sequences with binary activity labels (1 for active, 0 for inactive) **(Fig. 1)**.
103

### Directed virtual evolution from the neighborhood of an enzyme: TadA

105   To evaluate the enzyme-evolution capability of SPIN-dvEvo, we selected the tRNA-specific
106   adenosine deaminase TadA as our model enzyme. TadA, originally evolved to target tRNA, has been
107   engineered into adenine base editors that catalyze A•T→G•C conversions in DNA[22]. This system
108   utilizes an R67 DHFR-based codon reversion reporter to rapidly detect the intracellular DNA-
109   editing activity of evolved TadA variants, as in prior studies[23, 24]. In this codon reversion assay, an
110   active variant reverts a premature TAG stop codon to TGG in the reporter, enabling growth under
111   trimethoprim (TMP) selection **(Fig. 2A)**. We quantified intracellular DNA-editing activity as the
112   mutation frequency $f = N_1 / N_0$—the number of TMP-resistant revertants ($N_1$) divided by the total
113   number of viable cells plated without TMP ($N_0$)—and converted it into $\mu_{\text{s.p.b.}}$ (per base per
114   generation; See **Methods**).
115

116   We compiled a compact set of 10 TadA sequences spanning the wild type from *E. coli* (UniProt ID
117   P68398) **(Supplementary Table S1)** and previously engineered active variants from *E. coli* with 6-
118   20 mutations, (>88.6% sequence identity) and labeled all these sequences as 1. That is, we started
119   with a sequence cluster of close functional neighbors. An equal number of 10 hypothetical inactive
120   sequences were obtained by performing random mutations at 20% of positions in these TadA
121   sequences (See **Methods**).
122

123  We then employed SPIN-dvEvo to produce (evolved virtually) 1,000 sequences by starting from the
124  inactive sequences pool (20% random mutations). We confirmed that such virtual evolution started
125  from a tightly clustered sequence region (in red) and quickly expands to other regions according to
126  the t-SNE projection of ESM-2 sequence embeddings (**Fig 2B**).

127

128  To examine whether the TadA function was preserved during the virtual evolution, we obtained
129  sequence logos from 1000 natural TadA homologs compiled as in Ref [25] with a median sequence
130  identity of 34.1% and compared them to sequence logos from 1000 evolved sequences (median
131  identity 55.8%) in **Fig 2C**. The sequence motifs found previously[26] in the TadA family such as HAE
132  and PCXXC zinc-dependent deaminase motifs and structural-core signatures EVP and TLE were
133  also conserved in the evolved sequences while allowing substantial variation elsewhere. Thus,
134  essential sequence information preserved in the natural sequences evolved over billions of years
135  was captured by SPIN-dvEvo in a short 12 minutes of computing time with AMD EPYC 9654/ RTX
136  4090 (24 GB) starting from a local sequence cluster around *E. coli* TadA.

137

138  As protein structures play an essential role in enzymatic functions, we predicted structures of these
139  evolved sequences and compared them to the structure of wild-type TadA (PDB ID:2B3J). We
140  employed PLM-based OmegaFold[27] to make predictions because it does not require homologous
141  sequences for input, and therefore permits fast, large-scale calculations for all 1000 evolved
142  sequences. We obtained the distribution of structural accuracy (measured by TM-score[28], 1 for
143  perfect match and 0 for no match) for predicted structures of those evolved SPIN-dvEvo sequences
144  and compared it to two baseline models PLM-based sequence generators Pinal[29, 30] and structure-
145  based protein-design method ProteinMPNN[31]. ProteinMPNN employed a native structure template;
146  Pinal was prompted with a natural-language TadA functional description (adenosine
147  deaminase/base-editor context; EC 3.5.4.33) together with the wild-type TadA sequence **(Methods).**
148  The results show that most evolved sequences given by SPIN-dvEvo adopted near-native structures
149  (TMscore ~0.8, 89.6% sequences with TM-score>0.5), and was only slightly worse than the
150  structure-based method ProteinMPNN (TMscore ~0.95) (**Fig. 2D**). The baseline sequence-based
151  method Pinal shows a bimodal TM-score distribution, with one major peak at low TM-scores (~0.2–
152  0.3, 53.9% sequences with TM-score<0.5) and another in the near-native range (~0.8–0.9),
153  indicating a mixture of largely off-fold sequences and a smaller subset that retains the TadA fold.
154  An example of a predicted structure for a SPIN-dvEvo sequence is compared to the native structure
155  in **Supplementary Fig. 1**, highlighting near-perfect match, particularly in the regions interacting
156  with a DNA substrate and near catalytic core.

157

158  We further selected 60 evolved sequences to validate their enzymatic functions experimentally with
159  the R67 DHFR–based codon reversion assay (**Fig. 2A**). These 60 sequences were selected from the
160  above 1000 evolved sequences according to the high structure-confidence scores (normalized
161  pLDDT> 0.9 given by AlphaFold 3[32] with a single natural MSA to save computing time) and low
162  sequence identity ($\leqslant$0.5) to the wild type (as shown in **Fig. 2E**). Functional validation identified 23
163  active variants out of 60 tested (38.3% success rate). Activities spanned more than three orders of
164  magnitude, with several variants matching or exceeding the reference activity of *E. coli* TadA **(Fig.**
165  **2F, Supplementary Fig. 2 A, Table S3, Table S4)**. More importantly, these individually validated
166  functional sequences span 39–79% amino-acid identity to the *E. coli* TadA wild type, confirming

167   the ability of SPIN-dvEvo to find functional solutions by going significantly beyond the immediate
168   neighborhood of the starting sequences within the identity neighborhood of ≥88% *E. coli* TadA (**Fig.**
169   **2E**).

170

171   Given 60 newly experimentally tested sequences, we re-trained the LoRA model with the enlarged
172   binary-labelled dataset and performed sequence evolutions again by GA. The newly 1000 evolved
173   sequences (Round II) are now forming new sequence clusters (Fig. 2B). The TMscore distribution
174   of predicted structures for the second-round sequences improves over that of the first round. All
175   predicted structures (100%) are now with TMscore > 0.78 and the highest peak located at TMscore
176   of 0.88, compared to 0.80 in the first round (**Fig. 2D**). We tested 60 new variants chosen according
177   to high AlphaFold 3's pLDDT and low sequence similarity. In this second round, 31 of 60 new
178   variants were active. The higher success rate in Round II than in Round I (51% versus 38.3%)
179   indicates that incorporating new experimental labels with definitive inactive sequences improved
180   the classifier-guided evolution **(Supplementary Fig. 2B, Table S3, Table S5)**. Moreover, the
181   measured activity for the functional sequences in the second round shifted upward relative to the
182   first-round actives by one order of magnitude (**Fig. 2F**). These validated evolved sequences in
183   Round II are more divergent from wild type (29–54% identity, compared to 39–79% in the first
184   round; **Fig. 2E**), confirming the formation of new functional clusters with improved activity (**Fig.**
185   **2B**). This is remarkable considering the fact that only binary labels were employed to train SPIN-
186   dvEvo.

187

188   A few selected variants are illustrated along with positive and negative controls by plating on TMP-
189   selective medium (dvTadA-55 and dvTadA-56 from round 1; dvTadA-2-02 from round 2). These
190   evolved sequences produced markedly more TMP-resistant colonies than the negative control of
191   expressing only an Xten linker-T7RNAP cassette in place of TadA and thus lacking deaminase
192   activity and were comparable to the positive control (*E. coli* TadA) (**Fig. 2G**), consistent with robust
193   in vivo editing activity.

194

## Directed virtual evolution of intrinsically-disordered binder: anti-toxin CcdA

196

197   To test whether SPIN-dvEvo can generalize beyond enzymes with well-defined structures to
198   intrinsically disordered binding proteins, we applied it to the CcdA–CcdB toxin–antitoxin system.
199   In *E. coli*, the antitoxin CcdA is a 72-residue protein. Here we only engineered its C-terminal
200   segment (CcdA[36–72], 36 residues), which mediates binding to CcdB and thereby blocks CcdB
201   binding to GyrA to neutralize toxicity[33]. This 36-residue C-terminal domain is intrinsically
202   unstructured prior to binding to CcdB[33, 34]. We started from the canonical *E. coli* CcdA (P62552),
203   retrieved CcdA family homologs from closely related *Enterobacterales/ Gammaproteobacteri*a,
204   removed incomplete or atypical entries as well as those sequences at 100% sequence identity cutoff.
205   This yielded 22 close homologs (**Supplementary Table S2**) at 55.2–97.2% sequence identity. A
206   LoRA head on a frozen ESM-2 encoder was fine-tuned on this curated set and then coupled to the
207   GA to generate candidate binders, without introducing any CcdB sequence or structural information
208   during training or sampling. We chose this CcdA-CcdB system because bacterial growth is
209   directly correlated to the ability of the CcdA evolved by SPIN-dvEvo to bind and neutralize CcdB,
210   enabling straightforward functional selection (**Fig. 3A**).

211

212 As in the TadA case, we evolved 1000 CcdA variants by SPIN-dvEvo. As shown in **Fig 3B**, these
213 sequences moved far away from the original sequence cluster and formed multiple clusters
214 according to the t-SNE projections of the base ESM-2 embeddings. When we generated the
215 sequence-logo from SPIN-dvEvo sequences (with a median sequence identity of 50.2%), it has
216 similar sequence motifs as those from 100 natural homologs collected by querying the canonical
217 'Antitoxin CcdA' and filtering to a non-redundant set with a median sequence identity of 38.7%
218 from UniProtKB, suggesting that key binding determinants preserved such as W44, E54,[35] G63,
219 S64, F65, D71 and W72[36, 37] (**Fig. 3C,** blue box) in natural CcdA homologs were captured during
220 virtual evolution by SPIN-dvEvo, despite that it was started from a highly local seed set.

221

222 To test those sequences experimentally, we synthesized a library of 3,041 evolved CcdA variants
223 and evaluated them using a pooled bacterial growth selection, because the ability for the bacterium
224 to grow is correlated to the ability of the evolved CcdA to neutralize CcdB by binding **(Fig. 3A)**.
225 That is, the fitness of activity of CcdA variants can be measured by counting the number of a specific
226 variant pre- and post-selections from high-throughput sequencing[38] (**Fig. 3A**). We estimated
227 enrichment and uncertainty with the DiMSum pipeline [39,40] with Poisson–Delta variance modeling
228 and overdispersion correction. Among 3,041 synthesized CcdA variants, only 2,363 variants were
229 found with >30 reads and a minimum frequency of $10^{-6}$ in both the pre-selection and post-selection
230 libraries from high-throughput-sequencing data. Further application of an FDR-controlled filter
231 relative to internal stop-codon negative controls of $q\_value < 10^{-3}$ yielded 155 statistically
232 significant functional variants (a 6.6% hit rate). We further employed an effect-size threshold to
233 define more robust positives as those variants with $\log_2(\text{fitness}) > 3.0$, resulting in 62 active CcdA
234 variants (a 2.6% hit rate, **Fig. 3D**). These variants contain 26 with $\log_2(\text{fitness}) > 5$ and some
235 comparable to the fitness of *E. coli* CcdA ($\log_2(\text{fitness}) = 8.5$).

236

237 To validate the above high-throughput result, we selected four variants around the stringent
238 threshold of 3.0 with $\log_2(\text{fitness}) = 3.3, 3.3, 3.2,$ and $3.0,$ respectively, along with two positive
239 controls *E. coli* CcdA and an evolved variant with $\log_2(\text{fitness}) = 5.3$ for in vivo functional testing
240 (**Supplementary Table S7**). As shown in **Fig. 3E** by serial 10-fold dilution spot assays, we
241 confirmed that all variants with $\log_2(\text{fitness}) \geq 3.0$ are functional and the variant 878 with a larger
242 fitness value has stronger growth. In particular, the variant 1654 with $\log_2(\text{fitness}) = 3.0$ showed
243 weak growth only at the dilution factor of $10^2$. It is noted that sequences with $\log_2(\text{fitness}) \geq 3.0$
244 retained only ~60–70% sequence identity to the *E. coli* CcdA (**Supplementary Fig. 3**), indicating
245 substantial novelty among functional hits, given that only 36 residues were targeted for virtual
246 evolution.

247

## Discussion

249

250 SPIN-dvEvo directly addresses a practical gap in current directed virtual evolution: most existing
251 methods either require substantial labelled datasets to optimize a single scaffold locally, or function
252 as one-shot generators whose sequences are not coupled to an explicit score-and-search loop that
253 can be iterated with newly acquired labels. In contrast, SPIN-dvEvo mimics natural evolution by

254 employing a LoRA adaptor on the top of a frozen ESM-2 encoder to learn functional restraints. We
255 showed that the functional restraints learned from a few dozen positive, binary-labeled samples of
256 a highly homologous sequence cluster are sufficient to drive virtual evolution from dysfunctional
257 sequences to functionally active proteins that are substantially away from original positive
258 sequences by using a genetic algorithm. Some of these sequences, despite low sequence identity,
259 are experimentally validated for their functions on two illustrative cases: enzymatic activity (TadA
260 adenosine deaminase) and toxin-binding intrinsically disordered protein CcdA.

262 For virtual evolution of TadA enzyme, no structural information of was used to train SPIN-dvEvo
263 and to drive evolution. Yet most evolved TadA variants have TadA structural folds (**Fig 2D,**
264 **Supplementary Fig. 1**) in the first round (89.7% of sequences with predicted structural
265 accuracy >0.5 in TMscore). A minor peak with TMscore<0.5 in the first round was eliminated after
266 including experimental results from 60 variants (still in binary coding). The improved structural
267 similarity to the wild type highlights the importance of a larger and cleaner dataset because in the
268 first round, negatives represented by 20% random mutations may not be negatives. Interestingly, the
269 second-round success rate increased from 38% to 51% along with a one-order-of-magnitude
270 improvement in enzymatic activity, indicating that adding new experimental labels can improve
271 classifier-guided search even for enzymatic activity, despite lacking quantitative labels.

273 We have selected sequences with high confidence in predicted structures for experimental
274 validations. The high (38% in Round I) but not yet >90% success rate for TadA's virtual evolution
275 illustrates that the structural fold alone is not sufficient as an indicator of enzymatic activity. This is
276 because enzyme function not only requires highly precise active-site geometry and transition-state
277 stabilization, but also depends on compatible conformational dynamics and kinetics that enable
278 efficient substrate binding and product release on a productive timescale.[41-43] More studies are
279 needed to search for a better activity indicator as well as improving scoring for virtual evolution of
280 enzymes.

282 SPIN-dvEvo evolved functional TadA starting from a 20% randomly mutated (inactive) seed. We
283 kept starting sequences close to the TadA family where the LoRA scorer remains informative. We
284 also tried to start from fully random sequences and found that evolution from these sequences is not
285 productive according to analysis of their predicted structures. This indicates that the sequence space
286 is too large to be located by starting from purely random sequences within practical GA generations.
287 Nevertheless, it can start from one neighborhood of an active sequence to locate other
288 neighborhoods far away from the original sequence cluster as shown in **Fig. 2B** and **Fig 3B**.

290 However, the success rate of SPIN-dvEvo for a disordered protein CcdA is only 2.6%. This is much
291 lower than virtual evolution of TadA enzyme. Designing an intrinsically disordered protein is a
292 challenging task because activity is typically encoded in an ensemble of rapidly interconverting
293 conformations and mediated by weak, context-dependent interactions, so improvements in fold
294 stability or a single "best" structure provide little guidance. Recent progress has come from
295 explicitly optimizing ensemble-level objectives, for example by using sequence-to-ensemble
296 predictors for IDRs and by combining generative models with biophysical/simulation-based
297 forward models to design sequences that realize targeted disordered-state properties, as well as from

298   diffusion-based binder design strategies that focus the objective on functional binding constraints
299   rather than enforcing an ordered fold.[44] Here, we achieved a success (albeit low success rate) without
300   relying on any information from binding partner CcdB or predicted complex structures.

301

302   It is of interest to know how new functional clusters would have been evolved naturally if they were
303   mixed with natural homologs when building phylogenetic trees(see SI). As shown in **Fig. 4A** and
304   **Fig. 4B**, both virtually evolved TadA and CcdA are forming several phylogenetically distinct
305   clusters but do share common ancestors with naturally occurring sisters at different time points. For
306   TadA, this split corresponds to an evolutionary timescale on the order of ~0.2–1.2 Ga, based on
307   TimeTree-derived lineage-age estimates for these taxa[45-47]. Similarly, the estimates for the virtually
308   evolved CcdA clade dating to approximately 2.508 Ga as diverging from a Gammaproteobacteria-
309   associated branc. By comparison, these virtual evolutions took only 713 seconds for TadA and 761
310   seconds for CcdA by SPIN-dvEvo on a workstation equipped with an AMD EPYC 9654 (96-core,
311   2.4 GHz) CPU and an NVIDIA RTX 4090 GPU (24 GB).

312

313   SPIN-dvEvo was purposefully trained on binary-labeled sequences (1 for functional and 0 for
314   nonfunctional). This is because most proteins with known functions do not have a quantitative
315   functional label. One immediate improvement for SPIN-dvEvo is to employ a regression head,
316   rather than a classification head, when quantitative functional data such as a fitness score, binding
317   affinity, or enzymatic activity is available for a small dataset. A regression head would contain a
318   more accurate evolution direction than a classification head. This is a subject of an ongoing study.

319

320   One limitation of SPIN-dvEvo is its reliance on the ESM-2 650M. While ESM-2 is one of the best
321   protein language models available, we did not have the resource to test other language models or
322   utilization of multiple language models that could be potentially more beneficial than ESM-2 in
323   directed virtual evolution. Moreover ESM-2 may be inherently biased toward some protein
324   sequences with large family of homologous sequences as it was indiscriminately trained on all
325   protein sequences.[48, 49] Further studies in this area are needed.

326

327   Moreover, current implementation of SPIN-dvEvo is optimized for a single functional objective. A
328   multi-objective model, where functional objectives are optimized alongside other property
329   objectives such as stability, pH tolerance, and thermostability, can be easily implemented. This
330   research is also currently ongoing.

331

## Methods

### Data Collection and Curation

334   For TadA, we compiled 10 functional sequences from previously engineered DNA-editing TadA
335   variants[22] (listed in **Supplementary Table S1**). For CcdA, we constructed the 22-sequence set by
336   sequence-identity clustering of UniProtKB CcdA homologs. Starting from the canonical *E. coli*
337   CcdA (P62552; 36 aa) as the query, we retrieved annotated CcdA family homologs from closely
338   related *Enterobacterales/Gammaproteobacteria*. We then removed incomplete/aberrant entries (e.g.,
339   truncated sequences or atypical lengths) and identical sequences (100% sequence identity). This

340   yielded a deduplicated set by keeping only unique amino-acid sequences, yielding 22 non-redundant

341   homologs (accessions in **Supplementary Table S2**). To balance classes during few-shot training,

342   we generated synthetic decoys by randomly mutating 20% of residues in each positive sequence.

343   All positive sequences were labeled as 1 (functional), and all negative sequences—whether

344   randomly generated or literature-confirmed—were labeled as 0 (non-functional).

345

## LoRA-Based Model Adaptation

347   We adapted ESM-2 (650M parameters) to each task using low-rank adapters (LoRA) while keeping

348   all base model weights frozen. This model size offered a practical trade-off between representation

349   quality and computational cost, allowing training on a single 24–40 GB GPU.

350

351   LoRA modules were inserted into the self-attention Q/K/V projection layers of every transformer

352   block. For each pretrained projection $W \in \mathbb{R}^{d \times d}$, LoRA adds a trainable low-rank update $\Delta W = $

353   $s\,AB$ with rank $r$ and scaling $s = \alpha/r$:

$$\widetilde{W} = W + s\,AB, A \in \mathbb{R}^{d \times r}, B \in \mathbb{R}^{r \times d}, s = \alpha/r.$$

354

356   We employed $(r, \alpha) = (16,16)$. This setting adds 4,055,040 **LoRA** trainable parameters

357   (excluding the final linear head), corresponding to ~0.62% of the ~650M-parameter ESM-2 base

358   model, and was used throughout this work.

359

### Classification head (binary activity)

361   For binary activity prediction $y_i \in \{0,1\}$, the frozen ESM-2 produces a sequence representation

362   $h \in \mathbb{R}^d$ (pooled from token embeddings), which is mapped to a scalar logit

$$z = u \top h + b, score = f(x) = \sigma(z) \in [0,1]$$

363

365   The classifier was trained with binary cross-entropy:

$$\mathcal{L}_{\text{BCE}} = -\frac{1}{N} \sum_{i=1}^{N} [y_i \log p_i + (1 - y_i) \log (1 - p_i)].$$

366

368   Only the LoRA parameters $(A, B)$ and the classification head parameters $(u, b)$ were updated

369   during training; all ESM-2 weights remained frozen,

370

371   Sequences were truncated to 1,000 amino acids and fine-tuned for 5 epochs using AdamW (learning

372   rate $5 \times 10^{-4}$, weight decay $10^{-3}$) with a cosine schedule and gradient clipping ($\| \nabla \|_{\text{max}} = 0.5$).

373   LoRA adapters targeted the attention Q/K/V projections (rank $r = 16$, $\alpha = 16$, dropout 0.2; base

374   model frozen) with batch size 4.

375

## Genetic Algorithm Sampling

377   We performed an iterative mutation–crossover search guided by a fixed LoRA activity scorer.

378   Diversity arose implicitly from uniform parent sampling and stochastic point mutations, and exact

379   duplicate children were removed during population construction. In each generation, parent

380   sequences were sampled uniformly from the current mating pool and recombined to produce a child.

381  Each sequence was scored by the LoRA-adapted ESM-2 classifier, with the positive-class
382  probability computed from the logits as

383
$$p_{\text{act}}(x) = \frac{\exp(z_1)}{\exp(z_0) + \exp(z_1)}.$$

384  **Initialization.**
385  The initial population consisted of $N$ sequences (equal to the size of the seed pool), generated by
386  applying 20% random substitutions to a set of positive sequences (natural homologs or previously
387  engineered variants).
388

389  **Embedding & activity model.**
390  Each sequence was scored by a LoRA-tuned binary activity classifier on a frozen ESM-2 (650M),
391  returning $p_{\text{act}}(x)$ . (Sequence embeddings $\phi(x)$ were computed when needed for
392  visualization/analysis, by mean-pooling the last hidden state over non-special tokens followed by
393  L2 normalization.)
394

395  **Variation & constraints.**
396  Children were generated using a one-point crossover plus point-mutation operator
397  (mutate_crossover). One parent was first chosen as the base; a crossover point $c \in [1, \min(|p_1|, |$
398  $p_2|) - 1]$ was sampled, and the suffix was swapped with the other parent, yielding a recombinant
399  whose length follows the suffix donor. After crossover, each position was independently mutated
400  with probability 0.02 by substituting a uniformly sampled amino acid from the 20 standard residues.
401  Candidate sequences were filtered with NCBI segmasker[50] to reject sequences containing low-
402  complexity segments longer than 5 residues.
403

404  **Selection & replacement.**
405  For each parent sequence $x$ with score $p_{\text{act}}(x)$, a child $x'$ was proposed and evaluated to obtain
406  $p_{\text{act}}(x')$. The acceptance ratio was computed as

408
$$r = \frac{p_{\text{act}}(x')}{p_{\text{act}}(x)}.$$

407

409  The child was accepted if $r \geq 1$; otherwise, it was accepted with probability $0.125 \times r$. After
410  iterating this accept/reject update across the population, sequences were ranked by score (by $p_{\text{act}}$ in
411  probability-only mode) and the top 25% sequences (ranked by score) were retained as the mating
412  pool for the next generation. Unless stated otherwise, virtual evolutions were conducted for a pre-
413  specified number of generations (default is 100) and the per-generation mean score was logged.
414

415  **Parallel runs.** Each run outputs $N$ evolved sequences (set by the seed pool size). Larger libraries
416  were obtained by launching multiple independent runs in parallel with different random seeds and
417  by aggregating the resulting sequences.
418

419  **Sequence sampling of baseline models: ProteinMPNN and Pinal**

420  **ProteinMPNN**
421  A structure-templated baseline library was generated using ProteinMPNN in fixed-backbone design
422  mode with the experimental TadA reference structure as the input template (PDB: **2B3J**, Chain **A**).

423    The structure file was preprocessed to retain only the designed protein chain (non-protein atoms

424    were removed) and was provided to ProteinMPNN to compute per-position amino-acid distributions

425    conditioned on the backbone coordinates. 1,000 sequences were then sampled stochastically from

426    the model using temperature-controlled decoding (temperature = 0.1) with otherwise default

427    ProteinMPNN settings. Sampled sequences were post-processed to remove exact duplicates and

428    were written to FASTA for downstream structure prediction and evaluation.

429

430    **Prompt for Pinal Sequence Generation**

431    TadA Prompt：TadA (tRNA adenosine deaminase) is an enzyme that catalyzes the deamination of

432    adenosine to inosine at the wobble position (A34) of tRNA molecules, thereby expanding codon

433    recognition during translation, $adenosine_{34}$ in tRNA + H2O + $H^+$ = $inosine_{34}$ in tRNA + $NH4^+$.

434    EC:3.5.4.33. Through the introduction of two key mutations, A106V and D108N, the substrate

435    specificity of *E.coli* TadA has been reprogrammed, enabling the enzyme to catalyze adenosine and

436    cytosine deamination directly on DNA substrates. These engineered TadA variants are incorporated

437    into adenine base editors (ABEs), facilitating the precise conversion of A•T base pairs to G•C in

438    DNA without introducing double-strand breaks. This strategy offers an efficient and high-fidelity

439    tool for genome editing, particularly for the correction of disease-associated point mutations.

440

441    CcdA Prompt：CcdA is a bacterial antitoxin protein that functions as part of the CcdA–CcdB type

442    II toxin-antitoxin system encoded by the F plasmid in *Escherichia coli*. The CcdA protein

443    comprises 72 amino acids and adopts a two-domain structure: an N-terminal dimerization and DNA-

444    binding domain, followed by a C-terminal domain that binds to the CcdB toxin. In the absence of

445    CcdB, the C-terminal domain of CcdA is intrinsically disordered. Upon binding to CcdB, CcdA

446    undergoes a conformational change, forming a stable CcdA–CcdB complex that neutralizes the

447    toxicity of CcdB. This complex also acts as a transcriptional repressor of the ccd operon by binding

448    to the operator region. The CcdA–CcdB interaction is dynamic, with varying stoichiometries leading

449    to different complex formations, including (CcdA)2–(CcdB)2 and (CcdA)2–(CcdB)4 complexes.

450    The balance between CcdA and CcdB concentrations regulates the stability of the complex and the

451    repression of the operon. CcdA is subject to degradation by the Lon protease, which modulates the

452    levels of the antitoxin and, consequently, the activity of the toxin.",

453

454    # Structure prediction for SPIN-dvEvo sequences

455    SPIN-dvEvo sequences were evaluated by two complementary structure-prediction pipelines with

456    distinct roles. For high-throughput, distribution-level benchmarking across large libraries, we used

457    the MSA-free, PLM-based OmegaFold (v2.3.2)[27] to predict structures for all sequences, and

458    quantified global fold similarity to experimental references using TM-align (TM-score). For TadA,

459    PDB 2B3J (tRNA adenosine deaminase from *Staphylococcus aureus* in complex with RNA) was

460    used as the reference structure, because it provides a substrate-bound, catalytically relevant

461    conformation for a consistent TM-score fold-similarity benchmark; in contrast, the available *E. coli*

462    TadA structure PDB 1Z3A is apo and does not capture the RNA-engaged state[51]. TM-scores reported

463    in the main text refer to alignments between the native structure (PDB 2B3J) and OmegaFold-

464    predicted structures for SPIN-dvEvo-evolved variants.

465

466 Separately, we used AlphaFold3 (AF3) to obtain model confidence estimates for experimental
467 prioritization. To reduce the computational time for MSA retrieval, sequences were clustered at 80%
468 pairwise identity; a representative sequence per cluster was used to query the AF3 MSA database,
469 and the resulting MSAs were reused for all members of that cluster during the batch inference. For
470 TadA, per-chain pLDDT was used as the confidence metric.
471

## TadA experimental methods

### Reagents and Strains

474 All PCR reactions for cloning restriction sites and generating recombineering targeting cassettes
475 were performed using 2 × Phanta UniFi Master Mix DNA Polymerase (Vazyme, Nanjing, China,
476 P516-02). Colony PCR reactions for subsequent sequencing were conducted using Premix Taq™
477 DNA Polymerase (Takara, Dalian, China, R901A). Homologous recombination was performed
478 using the CloneExpress II One Step Cloning Kit (Vazyme, C112-02). All primers were synthesized
479 by GENEWIZ (Suzhou, China). Gene sequences for R67, which confers resistance to trimethoprim
480 (TMP), and engineered TadA variants were synthesized by GENERAL BIOL (Anhui, China).
481 Antibiotics, including ampicillin sodium (Sangon Biotech, Shanghai, China, A100339-0025) and
482 chloramphenicol, along with L-arabinose, were obtained from commercial sources. Chemically
483 competent *E. coli* DH5α cells were purchased from AlpalifeBio (Beijing, China), and chemically
484 competent *E. coli* DH10B cells were obtained from Biomed (Beijing, China).
485

### Plasmid construction

487 Engineered TadA variants used in this study are detailed in **Tables S3.** Expression plasmids for these
488 variants and T7 RNA polymerase (T7RNAP) were constructed using the pMuta088 vector backbone.
489 This backbone, derived from pDae079, carries the tandem PmCDA1-T7 RNA polymerase and uracil
490 glycosylase inhibitor (UGI).
491 For this study, expression plasmids for the engineered TadA variants were constructed by replacing
492 the PmCDA1 gene in the pMuta088 scaffold with the specific TadA sequences via homologous
493 recombination. A negative control plasmid (pT7RNAP-ΔTadA), expressing only an Xten-linker–
494 T7RNAP cassette, was constructed using the same strategy.[23]
495

496 TadA editing activity was quantified by measuring the frequency of trimethoprim-resistant
497 revertants following the general MutaT7/eMutaT7 workflow with minor modifications as detailed
498 below.[52] To characterize the A•T-to-G•C editing activity of TadA variants via antibiotic resistance
499 reversion, a reporter plasmid was developed. The *R67* gene, encoding dihydrofolate reductase
500 (DHFR) which confers resistance to trimethoprim (TMP), was cloned into a low-copy-number
501 plasmid (T7 promoter + terminators reporter plasmid). This was achieved by replacing the existing
502 *neoR/ kanR* gene (from Tn5) in a precursor plasmid via homologous recombination. In the final
503 reporter construct (pReporter-R67), expression of the *R67* gene is driven by a T7 promoter and
504 transcription is terminated by a tandem array of ten T7 terminators. Subsequently, site-directed
505 mutagenesis was employed to convert the tryptophan codon (TGG) at position 23 into a premature
506 stop codon (TAG), resulting in the final reporter construct pReporter-R67[W23*]. In this system, TadA-
507 mediated adenine deamination reverts the stop codon to wild-type, thereby restoring functional R67
508 expression and conferring TMP resistance.
509

**Evaluation of TadA Variant Activity in *E. coli***

To quantitatively characterize intracellular DNA-editing activity, the mutation (editing) frequency was defined as the ratio of the total TMP-resistant revertants to the total viable cell population.

To perform this assay, chemically competent *E. coli* DH10B cells were co-transformed with two plasmids: (1) The reporter plasmid (AmpR) pReporter-R67$^{W23*}$; (2) a chloramphenicol-resistant (CmR) expression plasmid (pDae079 derivative) encoding either pT7RNAP-ΔTadA (negative control), wild-type TadA (positive control), or an engineered TadA variant.

Transformants were selected on LB agar plates containing 100 μg/mL ampicillin and 25 μg/mL chloramphenicol, followed by incubation at 37℃ for 12–16 hours. Individual colonies were then inoculated directly into 10 mL of LB broth supplemented with 100 μg/mL ampicillin, 25 μg/mL chloramphenicol, and 0.2% (w/v) L-arabinose, followed by overnight incubation (16 h) at 37℃ with shaking at 220 rpm to initiate TadA expression and mutation accumulation.

On the following day, the overnight cultures were diluted 1:100 into fresh LB medium containing the same concentrations of ampicillin, chloramphenicol, and L-arabinose. To promote the fixation of mutations during active growth, these cultures were incubated for 4 hours at 37℃ with shaking at 220 rpm.

**Editing activity Assay**

At the endpoint, cultures were serially diluted (10-fold). To determine the total viable cell population ($N_0$), 10 μL aliquots of each serial dilution were spotted onto a single non-selective LB agar plate (containing 100 μg/mL ampicillin and 25 μg/mL chloramphenicol). To enumerate the TMP-resistant population ($N_1$), 300 μL aliquots of undiluted culture were spread onto three selective LB agar plates containing 20 μg/mL TMP (supplemented with the same antibiotics). Plating for $N_1$ was performed in triplicate. Colony counts were extrapolated to the full 10 mL culture volume to derive the total viable cells ($N_0$, scaled from the 10 μL spot and dilution factors) and total TMP-resistant revertants ($N_1$, scaled from the 300 μL spread). The frequency $f$ was calculated as the ratio $N_1 / N_0$.

**Mutation-rate calculation.**

For cross-study comparison to prior eMutaT7 reports, endpoint TMP-reversion frequencies were converted to per-base, per-generation mutation rates using the Luria–Delbrück rare-mutation approximation, where the expected mutant frequency satisfies $E[f] \approx \mu \ln(R_{\text{eff}})$. Although induction was maintained for 16 h, the calculation was normalized to the effective population expansion of the final outgrowth step, as mutation fixation is replication-dependent. This single 4 h propagation propagation round used a 1:100 reinoculation followed by regrowth to saturation, corresponding to $\sim 6.6$ generations ($G$). Assuming binary fission ($R_{\text{eff}} = 2^G$), $\ln(R_{\text{eff}}) = G\ln 2 \approx 4.57$. Because TMP-resistance restoration of the R67 reporter requires a single-base reversion, the effective target size was set to $S = 1$ and rates were reported as site-specific values (not normalized by the 192-bp reporter length):

$$\mu_{s.p.b.} = \frac{f}{G \ln 2} \approx \frac{f}{4.57} \ (per\ base\ per\ generation)$$

**Verification of R67 Gene Reversion**

To confirm that TMP resistance resulted from the targeted A•T-to-G•C edit in the *R67* gene, colony PCR was performed. For a representative subset of TadA variants tested, five independent TMP-resistant colonies were randomly picked from the selective agar plates for each selected variant. The *R67* gene locus was PCR-amplified from these colonies. The resulting amplicons were purified and subjected to Sanger sequencing (GENEWIZ, Suzhou, China). The obtained sequences were aligned with the reference $R67^{W23*}$ sequence and the wild-type *R67* gene sequence to identify the specific A-to-G reversion at codon 23 and any other potential off-target mutations within the amplified region.

## CcdA library generation, selection, and validation

**The Plasmids Construction**

The pUC57-Kan-ccdA/B plasmid was constructed to co-express the $CcdA^{36-72}$ domain and ccdB in *E. coli*. In this generation, the forward strand carries the J23119 promoter–driven $CcdA^{36-72}$ cassette, and the reverse strand carries the AmpR promoter–driven ccdB gene. A 21-bp spacer was inserted between the two stop codons to facilitate PCR amplification. Both $ccdA^{36-72}$ and ccdB were codon-optimized for *E. coli*, synthesized by General Biosystems, and subcloned into pUC57-Kan using PciI and NdeI restriction sites. For construction of the ccdA mutant library, we generated pUC57-Kan-2BspQI-ccdB by inserting two BspQI sites using primers BspQI-FP and BspQI-RP (**Supplementary Table S4**); this cloning step was performed in DB3.1 competent cells, which are resistant to ccdB toxicity. All plasmids were verified by Sanger sequencing, and complete vector and primer sequences are provided in **Supplementary Table S4**.

**Library Construction, Selection and High-Throughput Sequencing**

The SPIN-dvEvo-evolved $ccdA^{36-72}$ variants, codon-optimized for *E. coli*, were synthesized as an oligo pool containing the BspQI site by GenScript (China). The oligo pool was first amplified using PrimerSTAR HS DNA polymerase (Takara) and subsequently digested with BspQI. The digested fragments were then ligated into the BspQI-linearized pUC57-Kan-2BspQI-ccdB vector using T4 DNA ligase (Takara). Finally, the ligation products were purified and eluted in nuclease-free water, ready for electroporation.

The ligation products were electroporated into electrocompetent DB3.1 cells using a Bio-Rad Micropulser according to the manufacturer's protocol. Transformants were recovered in 10 mL of LB medium at 37°C for 1 hour. To estimate the library size, a portion of the culture was serially diluted, plated on LB agar containing kanamycin, and incubated for colony counting. Meanwhile, kanamycin was added to the main culture to a final concentration of 50 µg/mL, followed by incubation at 37°C for 10 hours. Subsequently, 100 µL of this culture was inoculated into 10 mL of fresh LB medium for amplification and subsequent plasmid extraction. The remainder of the overnight culture was harvested, resuspended in LB medium with 15% glycerol, and stored at -80°C. The initial, unselected ccdA library consisted of plasmids extracted from the CcdB-resistant DB3.1 strain. To perform functional selection, this library was electroporated into the CcdB-sensitive DH5α strain. Plasmids successfully recovered from DH5α transformants then represented the selected ccdA library. The $CcdA^{36-72}$ gene was PCR-amplified from both libraries using INDEX-containing

596 primers. The amplicons were gel-purified and sequenced by Salus Pro platform (ShenZhen Salus
597 Biomed Ltd)..

598

599 **In vivo functional analysis of the SPIN-dvEvo-evolved CcdA variants**
600 Selected CcdA variants (see **Supplementary Table S7**), encompassing a range of fitness scores,
601 were cloned into a pUC57-Kan-ccdA/B expression vector. All gene sequences were synthesized and
602 subsequently confirmed by DNA sequencing (General Biol). To evaluate *in vivo* function, 80 ng of
603 each plasmid construct was transformed into the ccdB-sensitive *Escherichia coli* strain DH5α.
604 Transformants were selected on LB agar plates supplemented with kanamycin. A ten-fold serial
605 dilution series of each transformation was plated to enable quantitative assessment. After incubation
606 (37 °C, 20 h), colony-forming units (CFUs) were counted at matched dilution factors and reported
607 as relative survival/growth under co-expression of ccdB, where functional CcdA variants rescue
608 colony formation (**Supplementary Fig. 6**).

609

610 **Sequencing data processing**
611 Raw reads were demultiplexed, adapter-trimmed, and quality-filtered. Reads were assigned to
612 SPIN-dvEvo-evolved variants by matching the variable region to the SPIN-dvEvo-evolved
613 dictionary (allowing ≤1 mismatch to tolerate sequencing error; ambiguous matches were discarded).
614 For each variant $i$ counts $c_i^{\text{pre}}$ and $c_i^{\text{post}}$ were tabulated. Samples with $<10^6$ total mapped reads
615 were excluded. Unless noted, a small pseudocount (α=0.5) was used only for descriptive
616 normalization of very low counts; final fitness estimates and uncertainty were obtained from
617 DiMSum.[39]

618

619 **Fitness estimation and statistical analysis**
620 After read mapping and quality filtering, 2,363 SPIN-dvEvo-evolved variants were retained for
621 downstream analysis. For each variant $s$, we denote the pre-selection and post-selection read counts
622 as $c_{\text{pre}}(s)$ and $c_{\text{post}}(s)$, with total library depths

624
$$N_{\text{pre}} = \sum_s c_{\text{pre}}(s), N_{\text{post}} = \sum_s c_{\text{post}}(s).$$

623

625 Counts were library-size normalized, and per-variant enrichment was defined as

627
$$ES(s) = \frac{c_{\text{post}}(s)/N_{\text{post}}}{c_{\text{pre}}(s)/N_{\text{pre}}}.$$

626

628 Variant fitness was then defined as the $\log_2$ enrichment without any wild-type normalization:

630
$$F(s) = \log_2 ES(s) = \log_2\left(\frac{c_{\text{post}}(s)}{c_{\text{pre}}(s)}\right) - \log_2\left(\frac{N_{\text{post}}}{N_{\text{pre}}}\right).$$

629

631 Fitness (log2 enrichment) and associated uncertainty were estimated with DiMSum (Poisson–Delta
632 model with overdispersion correction), consistent with the definition above.

633

634 To identify significantly enriched variants, we applied an FDR-controlled significance filter based
635 on DiMSum-reported $q$-values:
636
$$q\_\text{value} < 10^{-3},$$

637

638     For effect-size stratification, we labeled variants with log2 enrichment $F(s) > 3.0$ as functional and

639     those with $F(s) > 5.0$ as wild-type-like

## Code availability

The SPIN-dvEvo source code and the LoRA model weights for TadA and CcdA will be soon publicly available

## Data availability

All data generated or analyzed in this study are included in the main text and Supplementary Information. Input and output sequence files (including training seeds, natural homolog sets, and evolved sequence libraries), as well as analysis-ready intermediate results, are publicly available at https://zhouyq-lab.szbl.ac.cn/download/. Additional materials are available from the corresponding authors upon reasonable request.

## Author Contributions

ZC collected data, built the models, and performed sequence-based computational evolution. JT, TZ, and QN designed experiments and performed experimental validations. XZ helped with computational design. JZ and YZ initiated and supervised the project. YZ provided the funding support. YZ and ZC drafted the initial manuscript. All authors contributed to subsequent manuscript revision and approved the final version.

## Acknowledgements

## Conflict of Interest

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

## Reference

1. Arnold, F.H. Directed Evolution: Bringing New Chemistry to Life. *Angew Chem Int Ed Engl.* **57**, 4143-4148 (2018).
2. Bloom, J.D. & Arnold, F.H. In the light of directed evolution: Pathways of adaptive protein evolution. *Proc. Natl. Acad. Sci. U. S. A.* **106**, 9995-10000 (2009).
3. Tracewell, C.A. & Arnold, F.H. Directed enzyme evolution: climbing fitness peaks one amino acid at a time. *Curr. Opin. Chem. Biol.* **13**, 3-9 (2009).
4. Romero, P.A. & Arnold, F.H. Exploring protein fitness landscapes by directed evolution. *Nat. Rev. Mol. Cell Biol.* **10**, 866-876 (2009).

674   5.   Fowler, D.M. & Fields, S. Deep mutational scanning: a new style of protein science. *Nature*
675        *Methods* **11**, 801-807 (2014).
676   6.   Wrenbeck, E.E., Faber, M.S. & Whitehead, T.A. Deep sequencing methods for protein engineering
677        and design. *Curr. Opin. Struct. Biol.* **45**, 36-44 (2017).
678   7.   Wong, T.S., Roccatano, D., Zacharias, M. & Schwaneberg, U. A Statistical Analysis of Random
679        Mutagenesis Methods Used for Directed Protein Evolution. *J. Mol. Biol.* **355**, 858-871 (2006).
680   8.   Romero, P.A., Krause, A. & Arnold, F.H. Navigating the protein fitness landscape with Gaussian
681        processes. *Proc Natl Acad Sci U S A* **110**, E193-201 (2013).
682   9.   Bedbrook, C.N. et al. Machine learning-guided channelrhodopsin engineering enables minimally
683        invasive optogenetics. *Nat. Methods* **16**, 1176-1184 (2019).
684   10.  Fei, H. et al. Advancing protein evolution with inverse folding models integrating structural and
685        evolutionary constraints. *Cell* **188**, 4674-4692.e4619 (2025).
686   11.  Luo, Y. et al. ECNet is an evolutionary context-integrated deep learning framework for protein
687        engineering. *Nat Commun* **12**, 5743 (2021).
688   12.  Jiang, K. et al. Rapid in silico directed evolution by a protein language model with EVOLVEpro.
689        *Science* **387**, eadr6006 (2025).
690   13.  Vornholt, T. et al. Enhanced Sequence-Activity Mapping and Evolution of Artificial
691        Metalloenzymes by Active Learning. *ACS Central Science* **10**, 1357-1370 (2024).
692   14.  Li, X. et al. An iterative deep learning-guided algorithm for directed protein evolution. *iScience*
693        **28**, 113324 (2025).
694   15.  Yang, J. et al. Active learning-assisted directed evolution. *Nat Commun* **16**, 714 (2025).
695   16.  Biswas, S., Khimulya, G., Alley, E.C., Esvelt, K.M. & Church, G.M. Low-N protein engineering
696        with data-efficient deep learning. *Nature Methods* **18**, 389-396 (2021).
697   17.  Yang, L., Liang, X., Zhang, N. & Lu, L. STAR: A Web Server for Assisting Directed Protein
698        Evolution with Machine Learning. *ACS Omega* **8**, 44751-44756 (2023).
699   18.  Zhang, Q. et al. Integrating protein language models and automatic biofoundry for enhanced
700        protein evolution. *Nat Commun* **16**, 1553 (2025).
701   19.  Yang, J. et al. Active learning-assisted directed evolution. *Nat Commun* **16**, 714 (2025).
702   20.  Rubin, A.F. et al. A statistical framework for analyzing deep mutational scanning data. *Genome*
703        *Biol.* **18**, 150 (2017).
704   21.  Hu, J.E. et al. LoRA: Low-Rank Adaptation of Large Language Models. *arxiv* **abs/2106.09685**
705        (2021).
706   22.  Gaudelli, N.M. et al. Programmable base editing of A•T to G•C in genomic DNA without DNA
707        cleavage. *Nature* **551**, 464-471 (2017).
708   23.  Seo, D., Koh, B., Eom, G.-e., Kim, H.W. & Kim, S. A dual gene-specific mutator system installs
709        all transition mutations at similar frequencies in vivo. *Nucleic Acids Res.* **51**, e59-e59 (2023).
710   24.  Moore, C.L., Papa, L.J., III & Shoulders, M.D. A Processive Protein Chimera Introduces
711        Mutations across Defined DNA Regions In Vivo. *J. Am. Chem. Soc.* **140**, 11560-11564 (2018).
712   25.  Zhang, S. et al. TadA orthologs enable both cytosine and adenine editing of base editors. *Nat*
713        *Commun* **14**, 414 (2023).
714   26.  Yokobori, S.-i., Kitamura, A., Grosjean, H. & Bessho, Y. Life without tRNAArg–adenosine
715        deaminase TadA: evolutionary consequences of decoding the four CGN codons as arginine in
716        Mycoplasmas and other Mollicutes. *Nucleic Acids Res.* **41**, 6531-6543 (2013).
717   27.  Wu, R. et al. High-resolution de novo structure prediction from primary sequence. *bioRxiv*,

718      **2022.07.21.500999** (2022).

719    28.   Zhang, Y. & Skolnick, J. TM-align: a protein structure alignment algorithm based on the TM-
720      score. *Nucleic Acids Res.* **33**, 2302-2309 (2005).

721    29.   Dai, F. et al. Pinal: Toward *De Novo* Protein Design from Natural Language. *bioRixv*,
722      2024.2008.2001.606258 (2025).

723    30.   Yuan, J.S.a.C.H.a.Y.Z.a.J.S.a.X.Z.a.F. SaProt: Protein Language Modeling with Structure-aware
724      Vocabulary. *ICLR 2024* (2024).

725    31.   Dauparas, J. et al. Robust deep learning–based protein sequence design using ProteinMPNN.
726      *Science* **378**, 49-56 (2022).

727    32.   Abramson, J. et al. Accurate structure prediction of biomolecular interactions with AlphaFold 3.
728      *Nature* **630**, 493-500 (2024).

729    33.   Aghera, N.K. et al. Mechanism of CcdA-Mediated Rejuvenation of DNA Gyrase. *Structure* **28**,
730      562-572.e564 (2020).

731    34.   De Jonge, N. et al. Rejuvenation of CcdB-Poisoned Gyrase by an Intrinsically Disordered Protein
732      Domain. *Molecular Cell* **35**, 154-163 (2009).

733    35.   Bajaj, P., Manjunath, K. & Varadarajan, R. Structural and functional determinants inferred from
734      deep mutational scans. *Protein science : a publication of the Protein Society* **31**, e4357 (2022).

735    36.   Chandra, S., Manjunath, K., Asok, A. & Varadarajan, R. Mutational scan inferred binding
736      energetics and structure in intrinsically disordered protein CcdA. *Protein science : a publication of*
737      *the Protein Society* **32**, e4580 (2023).

738    37.   De Jonge, N. et al. Rejuvenation of CcdB-poisoned gyrase by an intrinsically disordered protein
739      domain. *Mol. Cell* **35**, 154-163 (2009).

740    38.   Chandra, S. et al. The High Mutational Sensitivity of ccdA Antitoxin Is Linked to Codon
741      Optimality. *Mol. Biol. Evol.* **39** (2022).

742    39.   Faure, A.J., Schmiedel, J.M., Baeza-Centurion, P. & Lehner, B. DiMSum: an error model and
743      pipeline for analyzing deep mutational scanning data and diagnosing common experimental
744      pathologies. *Genome Biol.* **21**, 207 (2020).

745    40.   Faure, A.J. et al. Mapping the energetic and allosteric landscapes of protein binding domains.
746      *Nature* **604**, 175-183 (2022).

747    41.   Lienhard, G.E. Enzymatic catalysis and transition-state theory. *Science*    **180**, 149-154 (1973).

748    42.   Hanson, J.A. et al. Illuminating the mechanistic roles of enzyme conformational dynamics. *Proc*
749      *Natl Acad Sci U S A* **104**, 18055-18060 (2007).

750    43.   Acevedo, O. & Jorgensen, W.L. Advances in quantum and molecular mechanical (QM/MM)
751      simulations for organic and enzymatic reactions. *Acc. Chem. Res.* **43**, 142-151 (2010).

752    44.   Lotthammer, J.M., Ginell, G.M., Griffith, D., Emenecker, R.J. & Holehouse, A.S. Direct
753      prediction of intrinsically disordered protein conformational properties from sequence. *Nature*
754      *Methods* **21**, 465-476 (2024).

755    45.   Kumar, S. et al. TimeTree 5: An Expanded Resource for Species Divergence Times. *Mol. Biol.*
756      *Evol.* **39** (2022).

757    46.   Feng, D.-F., Cho, G. & Doolittle, R.F. Determining divergence times with a protein clock: Update
758      and reevaluation. *Proc Natl Acad Sci U S A* **94**, 13028-13033 (1997).

759    47.   Konaté, M.M. et al. Molecular function limits divergent protein evolution on planetary timescales.
760      *eLife* **8** (2019).

761    48.   Ding, F. & Steinhardt, J. Protein language models are biased by unequal sequence sampling across

762  the tree of life. *ICLR 2024 Workshop on Generative and Experimental Perspectives for*
763  *Biomolecular Design*, 2024.2003.2007.584001 (2024).

764 49. Notin, P. et al. in Proceedings of the 39th International Conference on Machine Learning,
765  *Proceedings of Machine Learning Research* **162**. 16990--17017 (2022).

766 50. Madden T, C.C. BLAST+ features. *National Center for Biotechnology Information (US)* (2008).

767 51. Rallapalli, K.L., Ranzau, B.L., Ganapathy, K.R., Paesani, F. & Komor, A.C. Combined
768  Theoretical, Bioinformatic, and Biochemical Analyses of RNA Editing by Adenine Base Editors.
769  *The CRISPR journal* **5**, 294-310 (2022).

770 52. Park, H. & Kim, S. Gene-specific mutagenesis enables rapid continuous evolution of enzymes in
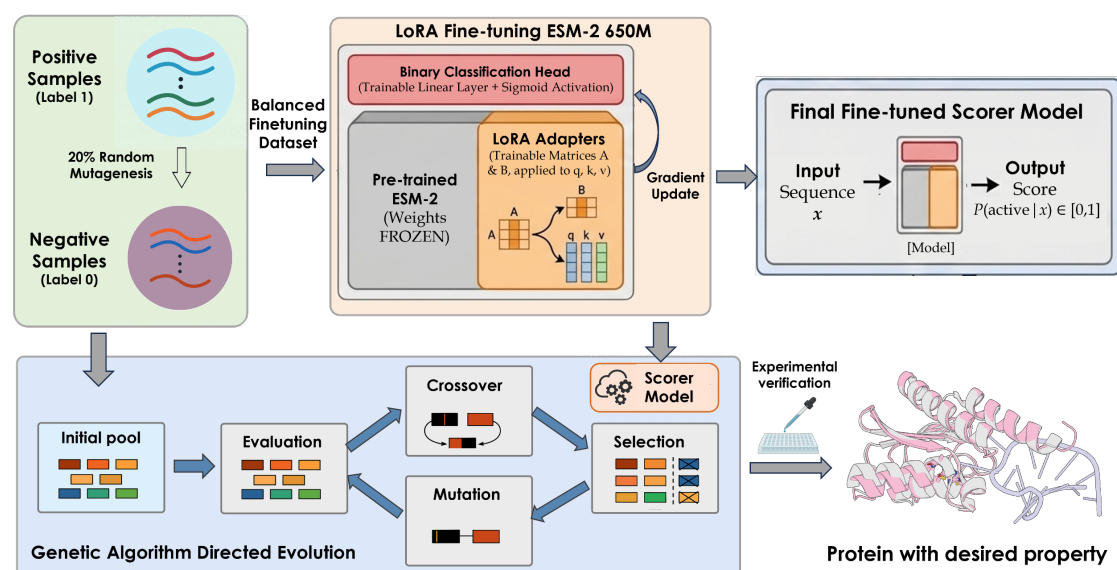771  vivo. *Nucleic acids research* **49**, e32 (2021).

772

773

774

775 **Figure 1**. **Schematic overview of the framework for directed virtual evolution: SPIN-dvEvo**. A
776 LoRA-adapted ESM-2 model is fine-tuned utilizing only a few curated positive and randomly
777 generated negative (binary) samples. The model is then integrated into a genetic algorithm as a
778 scorer to iteratively evolve sequences toward desired functionality but away from the original
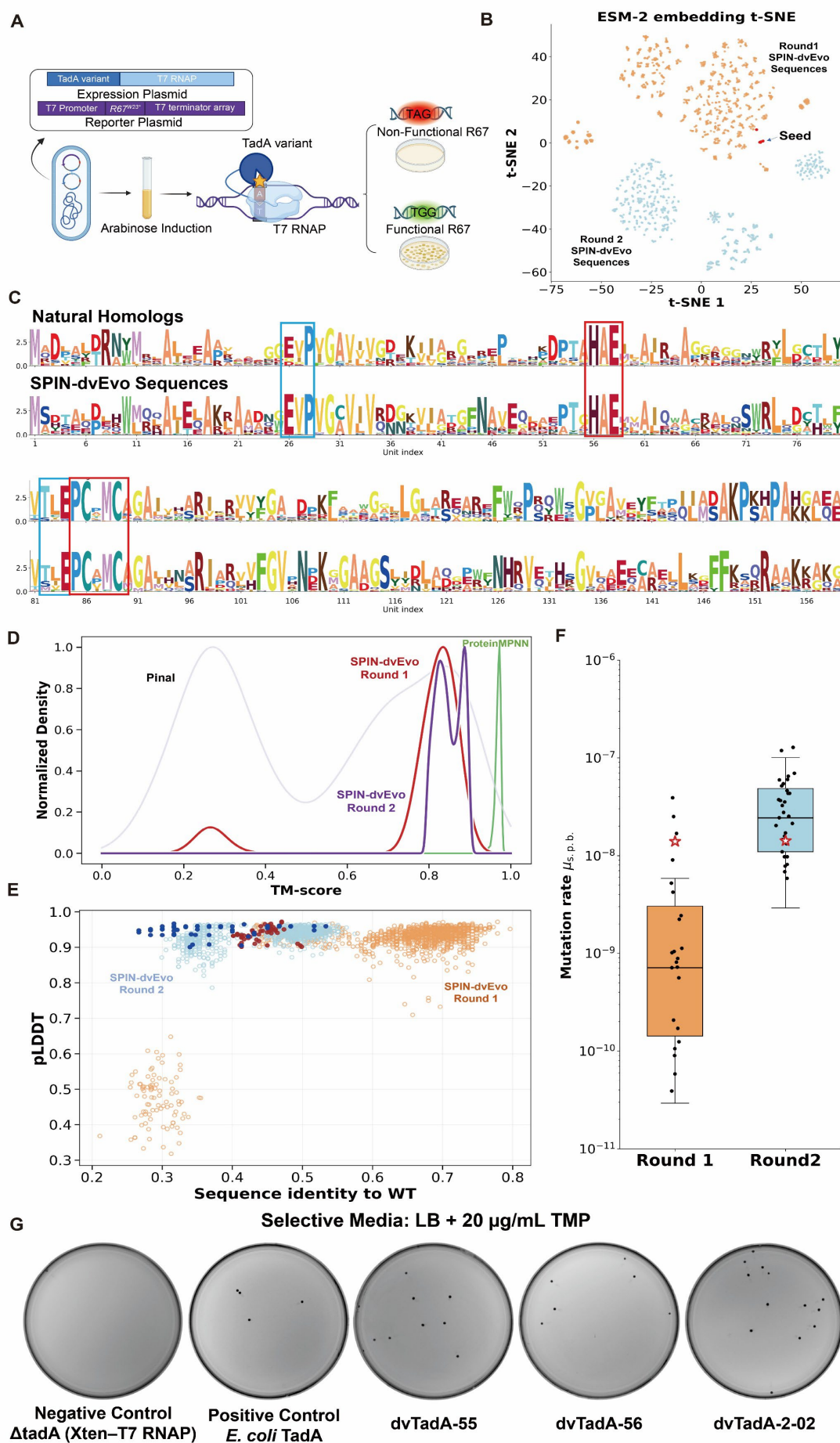779 sequence cluster.

780

781

782 **Figure 2.  Validation of virtually evolved enzyme TadA from sequence motifs, predicted structures**

783 **and experiments.** (**A**) Schematic of the experimental reporter system employed for quantifying A•T-

784 to-G•C editing activity. TadA-mediated reversion of a premature TAG stop codon to a TGG codon

785 in the *R67* gene confers resistance to trimethoprim (TMP), enabling selection of active variants. (**B**)

786 Newly emerged clusters from directed virtual evolution by SPIN-dvEvo according to the t-SNE

787 projections of the base ESM-2 embeddings of the 10 starting TadA sequences to 1000 evolved

788 sequences in Round 1 and Round 2. (**C**) Similar conserved functional and structural core motifs

789 between virtual evolved sequences and natural homologs (top). (**D**) The accuracy for the predicted

790 structures (according to TMscore) for 1000 TadA variants generated by four models (sequence generators

791 Pinal and structure-based designs ProteinMPNN) compared to those given by SPIN-dvEvo in two rounds.

792 (**E**) Scatter plot of the predicted confidence score pLDDT versus sequence identity to the wild type

793 (*E. coli* TadA) for 1000 evolved sequences by SPIN-dvEvo in Round 1 and Round 2. The 60

794 experimentally tested sequences selected from Round 1 and the 60 from Round 2 are highlighted as

795 filled points. (**F**) Boxplots comparing experimental activities of validated first- and second-round

796 evolved TadA sequences, showing an upward-shifted distribution after including the first-round

797 result in training. (**G**) Illustrative examples of the plates from the R67 DHFR–based *E. coli* reporter

798 assay on TMP-selective medium. Shown are the negative control (ΔTadA cells only expressing Xten

799 linker–T7RNAP), a positive-control TadA variant (*E. coli* TadA), and cells expressing SPIN-dvEvo-

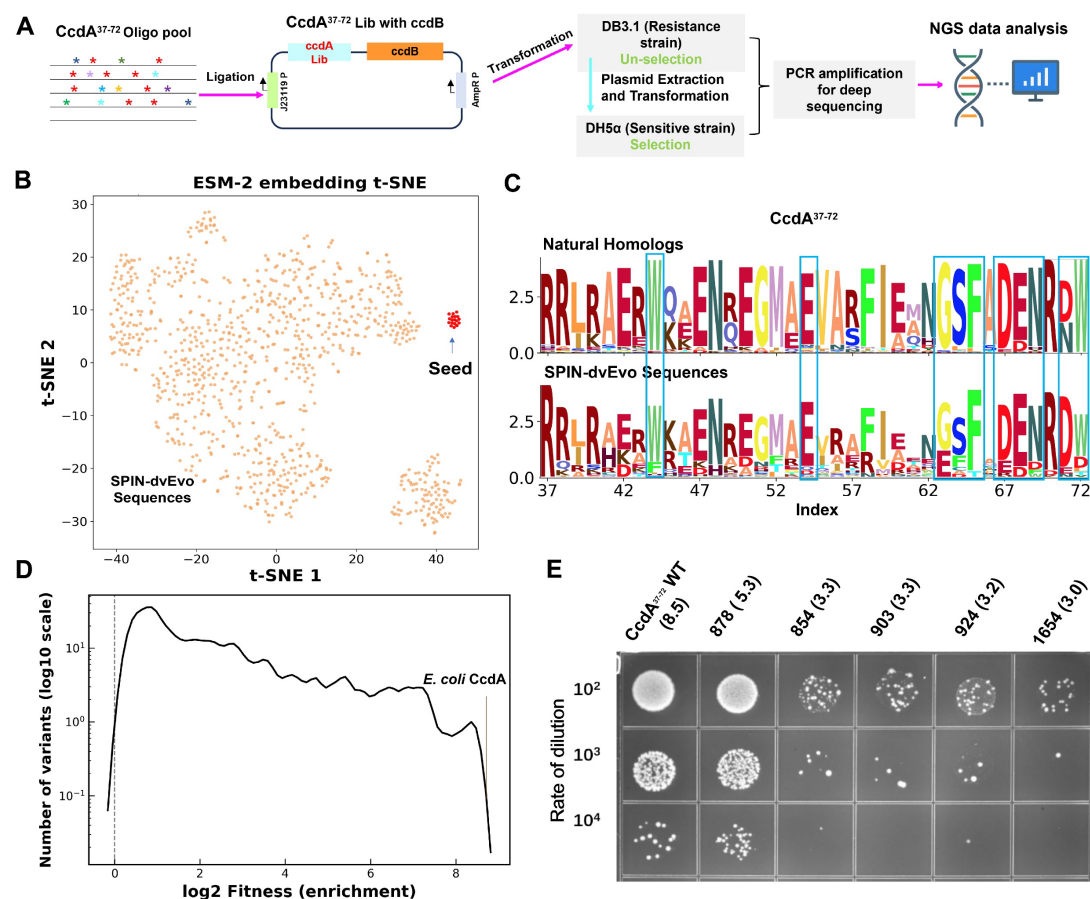800 evolved TadA variants dvTadA-55, dvTadA-56 and dvTadA-2-02.

801

**Figure 3**. **Experimental validation of evolved variant library of intrinsically disordered protein CcdA.**

(**A**) Schematic diagram for high-throughput validation of evolved CcdA according to the ability of a CcdA variant that can neutralize CcdB toxin in *E. coli* growth, measured by sequence counts pre and post selections. (**B**) Emergence of new clusters in SPIN-dvEvo sequences evolved from the starting 22 natural CcdA input sequences according to the t-SNE projections of the base ESM-2 embeddings. (**C**) Sequence motifs from SPIN-dvEvo sequences are highly similar to those obtained from natural homologs according to key conserved residues highlighted in blue boxes. (**D**) The distribution in number of variants as measured fitness scores (Log$_2$ fitness distributions normalized by the library size). (**E**) Activity confirmation of selected variants according to their fitness. Serial 10-fold dilution spot assay showing CcdA WT from *E. coli* and five CcdA variants (1654 (Log$_2$ fitness = 3.0), 924 (Log$_2$ fitness = 3.2), 903 (Log$_2$ fitness = 3.3), 854 (Log$_2$ fitness = 3.3), and 878 (Log$_2$ fitness = 5.3) along with the wild type (Log$_2$ fitness = 8.5)) for rescuing toxin CcdB at a dilution factor of $10^2$–$10^4$. Higher colony counts indicate stronger neutralization activity.
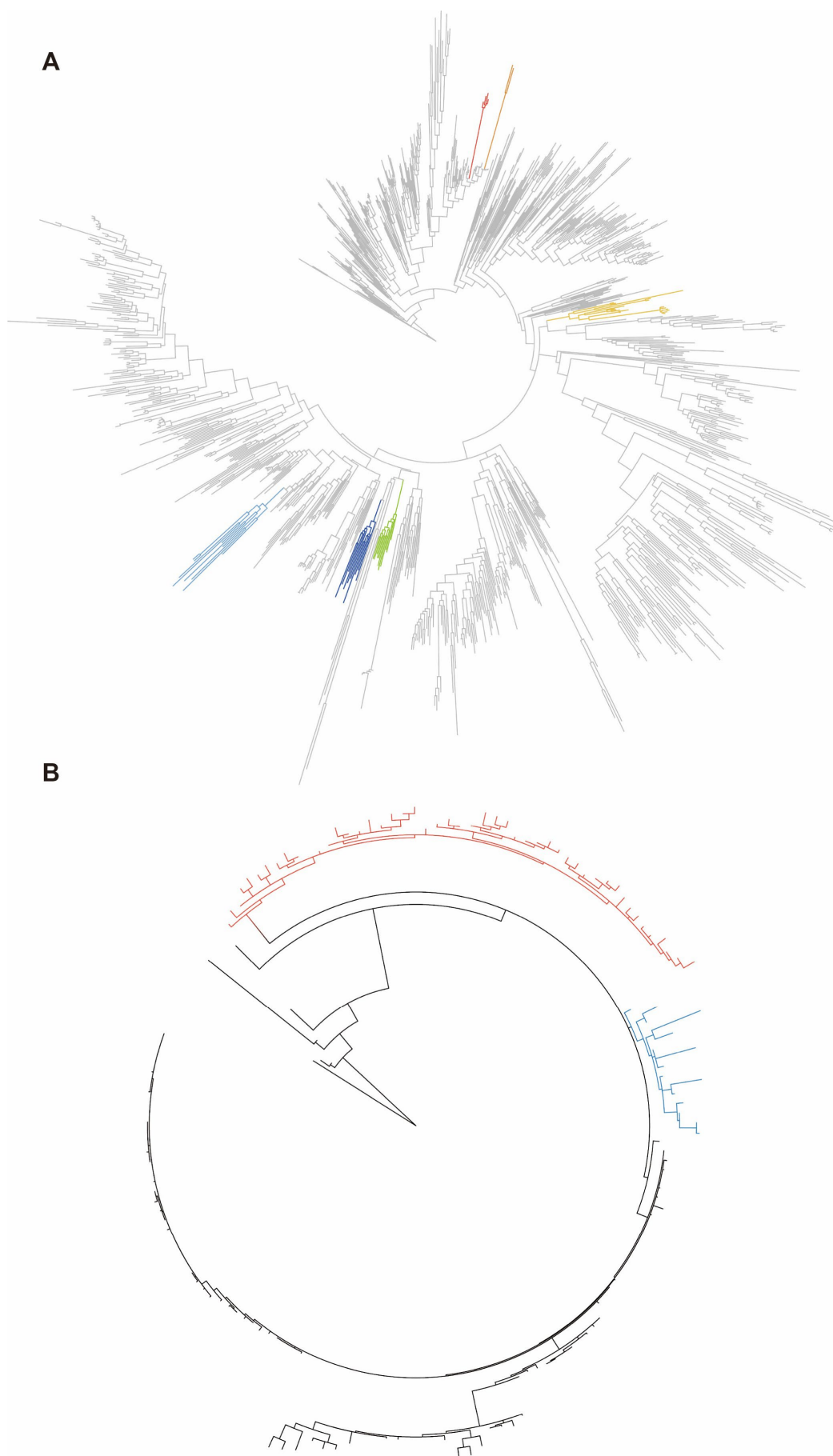
818

**Figure 4. Phylogenetic novelty of SPIN-dvEvo TadA and CcdA variants in joint natural–evolved trees.**

Maximum-likelihood phylogenies inferred from multiple sequence alignments containing natural homologs and experimentally validated SPIN-dvEvo evolved variants (sequences combined prior to alignment and tree building). Triangles denote nodes with bootstrap support in the 70–100 range. **(A)** TadA**:** alignment includes 1000 natural TadA homologs and 54 dvTadA variants. Highlighted sectors mark major, evolve-enriched dvTadA branches separated from dominant natural clades, supporting phylogenetically distinct lineages beyond the initial natural neighborhood. **(B)** CcdA: alignment includes 100 natural CcdA homologs and 62 dvCcdA variants. Light-blue and red sectors highlight two major evolved dvCcdA branches, indicating phylogenetically distinct lineages relative to the bulk of natural homologs.