

---

# In silico evaluation of pre-training strategies based on synthetic data for functional DNA generation

---

Anonymous Authors<sup>1</sup>

## Abstract

Despite their broad therapeutic and biotechnological potential, deoxyribozymes remain challenging to design rationally due to few consensus sequences, context-dependent activity, target inaccessibility, and a lack of large curated datasets for machine learning. This work proposes a robust multi-stage computational and experimental pipeline for DNAzyme discovery to address these limitations. The pipeline introduces domain-grounded parametric constraints for synthetic data generation, model training, and candidate evaluation. Pre-training strategies are based on two synthetic data approaches: sequences with minimal free energy boundaries matching stable structures, and sequences matching property distributions that produce structures similar to real DNAzymes. A hierarchical screening approach using machine learning models and structural priors then selects candidates for laboratory validation under conditions that include plausible cofactors crucial for catalytic activity. Overall, this pipeline provides a resource-efficient strategy for exploring novel DNAzyme catalytic cores and enabling predictable DNAzyme activity.

## 1. INTRODUCTION

Beyond the canonical role of genetic data storage, DNA molecules possess the capacity to function as catalysts, which was discovered by Breaker and Joyes who synthesized the first deoxyribozyme (DNAzyme) in 1994<sup>1</sup>. These catalytic DNA molecules are typically single-stranded sequences obtained through combinatorial *in vitro* selection techniques<sup>2,3</sup>. Similar to their RNA counterparts (ribozymes), both DNA and RNA can exhibit the ability to bind ligands with high affinity and specificity, and catalyze

a broad spectrum of chemical reactions<sup>4</sup>. Initially conceived as potential gene-silencing agents due to their capacity for the sequence-specific cleavage of target mRNA, subsequent research has led to the isolation of numerous DNAzymes that facilitate a diverse array of chemical transformations. This functional diversity has significantly expanded their scope of application, firmly establishing them as genuine biocatalysts alongside ribozymes and proteinaceous enzymes<sup>5,6</sup>.

The DNAzymes 8-17 and 10-23 are among the most well-characterized catalytic DNA molecules, renowned for their ability to mediate the sequence-specific cleavage of RNA<sup>7</sup>. The therapeutic applications of DNAzymes extend to a broad range of pathologies, including viral infections, cancer, and cardiovascular diseases<sup>8-11</sup>. Furthermore, DNAzymes represent a promising strategy for anti-inflammatory therapy. A promising avenue is anti-inflammatory therapy, exemplified in allergic asthma by the DNAzyme SB010, which targets GATA3 mRNA and has been shown to attenuate asthmatic responses in patients<sup>12, 13</sup>. DNAzymes also serve as valuable research tools for probing gene function, such as using a Twist-targeting DNAzyme to elucidate its role in apoptotic pathways<sup>14-16</sup>. Despite this promise, a central design challenge persists. While modern methods like *in vitro* selection and machine learning have revolutionized the optimization of binding arms for high-affinity, specific target recognition, the catalytic core remains a bottleneck. Designers are largely constrained to a limited set of known consensus sequences (e.g., from 8-17 or 10-23), whose efficacy is highly context-dependent and unpredictable. Optimizing this core is a laborious, empirical process that offers no guarantee of success and requires extensive resource-intensive screening. This highlights a critical need for novel, rational approaches to design *de novo* catalytic cores with high and predictable activity.

Target site inaccessibility is a common limitation for many nucleic acid-based tools, including DNAzymes. Strategies to overcome this face significant constraints. While reaction conditions can be modified, this often lacks therapeutic relevance. Extending binding arms can impede product release, hindering rapid diagnostics. Chemically modified nucleotides enhance hybridization and stability<sup>17-19</sup> but re-

<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

quire laborious optimization, costly synthesis, and are difficult to generalize for long, structured RNA (lsRNA) targets. Computational prediction of RNA accessibility loses accuracy with longer transcripts<sup>20</sup>, and experimental mapping is often too resource-intensive for rapid development<sup>21–23</sup>. Using antisense oligonucleotides (ASOs) to displace obstructive structures near the cleavage site is promising but not a universal solution.

While machine learning (ML) holds transformative potential for molecular design, its application to DNazymes remains a largely untapped frontier. Despite the introduction of foundational resources like the comprehensive DNAmoreDB<sup>24</sup>, the field still suffers from a critical scarcity of large, well-curated, and functionally characterized datasets that are essential for training robust, predictive ML models. This data paucity is the primary bottleneck hindering the development of sophisticated computational generative tools capable of moving beyond exploratory analysis to true de novo design and reliable activity prediction.

The existing landscape of tools is notably sparse. A prominent exception is SequenceCraft<sup>25</sup>, employs ML for the exploratory analysis of RNA-cleaving deoxyribozymes. However, while tools like SequenceCraft represent a crucial first step, the field still lacks dedicated, high-throughput platforms for the systematic generation and in silico validation of novel DNzyme candidates. The challenge extends beyond mere sequence generation; it requires predicting complex structure-activity relationships, catalytic efficiency, and specificity within a given biological context—a task for which current datasets and algorithms are insufficient.

Generative design represents a paradigm shift in DNzyme discovery, potentially enabling the exploration of novel catalytic centers and providing a powerful alternative to traditional SELEX procedures. However, the development of such models for functional nucleic acids remains a significant challenge. While pioneering frameworks for biomolecular design exist such as RFDiffusion<sup>26</sup> for protein structures and EvoDiff for generative sequence modeling, specialized<sup>27</sup> models for functional nucleic acids, and DNzymes in particular, are absent. The path to developing such models is fraught with obstacles. It requires the curation of large, high-quality datasets of DNzyme sequences with validated catalytic activity and structural information. Furthermore, model architecture must be capable of learning the complex relationship between primary sequence, higher-order structure, and catalytic function, which is poorly understood. Finally, the computational cost of training such a model and the subsequent experimental validation of its generated candidates present a substantial resource barrier that has yet to be overcome.

Therefore, there is a pressing need for two parallel advancements: first, the creation of expanded, high-quality gen-

erative datasets that capture a wider diversity of validated sequences; and second, the development of next-generation ML pipelines specifically engineered to leverage this data for the generative design of potential DNzyme candidates. Our work proposes a robust methodology for DNzyme discovery that comprises a multi-stage computational and experimental pipeline. The process initiates with the generation of a large-scale, general-purpose datasets for primary pre-training, followed by the selection of optimal model architectures for this task. Subsequent DNzyme-specific fine-tuning is augmented by discriminative model filtering to select high-probability candidates. Promising sequences are further analyzed through secondary and tertiary structure prediction to identify functionally viable motifs. To infer catalytic potential, we employ clustering with known DNzymes to hypothesize cofactor requirements, which are then evaluated via molecular docking and binding affinity simulations. The final stage involves experimental validation of the top-ranked candidates.

## 2. RESULTS AND DISCUSSION

### 2.1. Domain-Grounded Priors for DNzyme Candidate Evaluation

Given the scarcity of experimentally validated DNzyme sequences, we established qualitative and quantitative criteria to guide synthetic data generation, model training, and candidate filtering. These domain-grounded rules—analogueous to Lipinski’s rules in drug discovery—enable systematic evaluation of generative model outputs and provide a structured protocol for selecting high-fidelity DNzyme candidates.

#### 2.1.1. SEQUENCE CRAFT DATASET

The Sequence Craft platform enables predictive modeling of DNzyme activity, shifting from stochastic discovery to bioinformatic design. Its comprehensive dataset contains 349 unique catalytic cores across 549 systems, annotated with 71 reaction condition parameters. A predictive model trained on this data achieved  $R^2 > 0.93$  for estimating catalytic rate constants ( $k_{obs}$ ) prior to synthesis.

#### 2.1.2. STRUCTURAL ANALYSIS

DNzyme therapeutic function depends on three-dimensional folding that defines the catalytic core. Structural analysis identifies active sites, substrate recognition mechanisms, cofactor coordination, and motifs that enhance in vivo stability. By transitioning from primary sequences to structural models, we derived benchmark criteria for minimal free energy (MFE), sequence length, GC content, and secondary structure elements (hairpins, loops).

**MFE calculations.** Catalytic efficiency is governed by thermodynamic stability. Low MFE ensures a resilient catalytic core and stable metal-ion binding at 37°C, while high MFE indicates structural instability. In the Sequence Craft dataset, mean MFE is  $-7.5 \pm 4.47$  kcal/mol, with  $Mn^{2+}$ -associated sequences showing the lowest mean ( $-8.7$  kcal/mol) due to  $Mn^{2+}$ 's superior coordination versatility compared to  $Mg^{2+}$ . To ensure robust catalytic conformations in our generative pipeline, we applied a conservative MFE threshold of  $-10.0$  kcal/mol for pre-training data generation.

**Secondary structure analysis.** Using dot-bracket notation (parentheses for paired nucleotides, dots for unpaired regions), we characterised structural features across the dataset. Figure 1 presents distributions of sequence lengths, GC composition, hairpin center positions, and hairpin lengths.

The dataset displays a highly optimised structural profile. Sequence lengths (Fig. 1a) cluster near 40 nucleotides—sufficient for a complex catalytic core yet efficient for synthesis. GC content (Fig. 1b) centres at 45–50

From these observations we derived heuristic filtering rules. Although fully unpaired sequences and stems longer than 9 nucleotides appear as rare exceptions in the reference set, their functional reliability is uncertain. We therefore excluded candidates with: (i) no paired nucleotides (unable to form a stable catalytic core), and (ii) stems  $> 9$ nt (excessive rigidity may impair conformational flexibility during catalysis). Additionally, the normalised hairpin start and end positions exhibited a conserved distribution; we adopted the 5th and 95th percentiles (0.232 and 0.831) as selection thresholds, retaining only candidates whose hairpins fall within this window.

## 2.2. Synthetic Data Generation and Curation

The scarcity of experimentally validated DNAszymes limits the application of deep generative models to catalytic DNA design. Although DNA foundation models perform well in genomic annotation tasks, they are typically trained on natural genomic sequences and therefore fail to capture the biochemical and thermodynamic constraints underlying artificial catalysis. To address this limitation, we generated large-scale synthetic datasets guided by established heuristic rules and structural priors.

### 2.2.1. SAMPLE PREPARATION FOR MODEL PRE-TRAINING

Two independent synthetic corpora were generated to evaluate alternative strategies for modeling functional DNAszyme space. The first approach prioritized thermodynamic stability by selecting sequences with low minimum free energy

(MFE), while the second reproduced the statistical characteristics of experimentally derived datasets.

### MFE based approach

The first strategy focused on generating structurally stable sequences with MFE values below  $-10.0$  kcal/mol. The optimization process, summarized in Algorithm 1, uses a genetic algorithm in which sequence fitness is determined by predicted secondary-structure stability. Starting from a random population, sequences with lower MFE values are preferentially retained, while poorly performing individuals are iteratively replaced through crossover and mutation operations. Selection pressure is gradually increased across generations, concentrating the search on energetically favorable regions of sequence space. The procedure continues until 250,000 sequences are generated.

---

#### Algorithm 1 MFE based approach

---

**Initialize:** Population  $P$  of random sequences,  $k \leftarrow 10$ ,  $G \leftarrow 100$ , Target  $N \leftarrow 250,000$   
**Compute Fitness:**  $MFE(s)$  for all  $s \in P$   
**while** Total sequences  $< N$  **and**  $gen < G$  **do**  
  Sort  $P$  by  $MFE$  (ascending: more stable first)  
  Identify  $P_{bottom}$  as the worst 90% of  $P$   
   $k \leftarrow k \times 1.1$   
   $P_{parents} \leftarrow$  Select  $k$  least-fit individuals  
  Apply crossover ( $p = 0.75$ ) and mutation ( $p = 0.25$ )  
  Recompute  $MFE$  for offspring  
  Replace population and reduce its size  
   $gen \leftarrow gen + 1$   
**end while** Stable sequences with low MFE

---

**Distribution based approach** The second strategy aimed to reproduce the statistical profile of the Sequence Craft dataset. As outlined in Algorithm 2, reference sequences were cleaned from duplicates and outliers, after which key descriptors were extracted, including GC content, Shannon entropy, and mean Levenshtein distance. Synthetic sequences were evolved to minimize the Wasserstein distance between generated and reference feature distributions using tournament selection, crossover, and mutation, while a hall-of-fame mechanism preserved the best candidates. This approach produced 250,000 sequences closely matching the compositional properties of the reference corpus.

The two strategies differ in computational complexity. The MFE-based method requires RNA/DNA folding with complexity  $\mathcal{O}(N \cdot L^3)$ , where  $N$  is the number of sequences and  $L$  is sequence length. The distribution-based method additionally computes pairwise Levenshtein distances against a reference set of size  $M$ , leading to  $\mathcal{O}(N \cdot M \cdot L^2)$  complexity. For typical DNAszyme settings ( $L \approx 40$ ,  $M \approx 1000$ ), this makes it roughly 25× more expensive per sequence. Thus, the MFE-based strategy is more scalable, whereas the

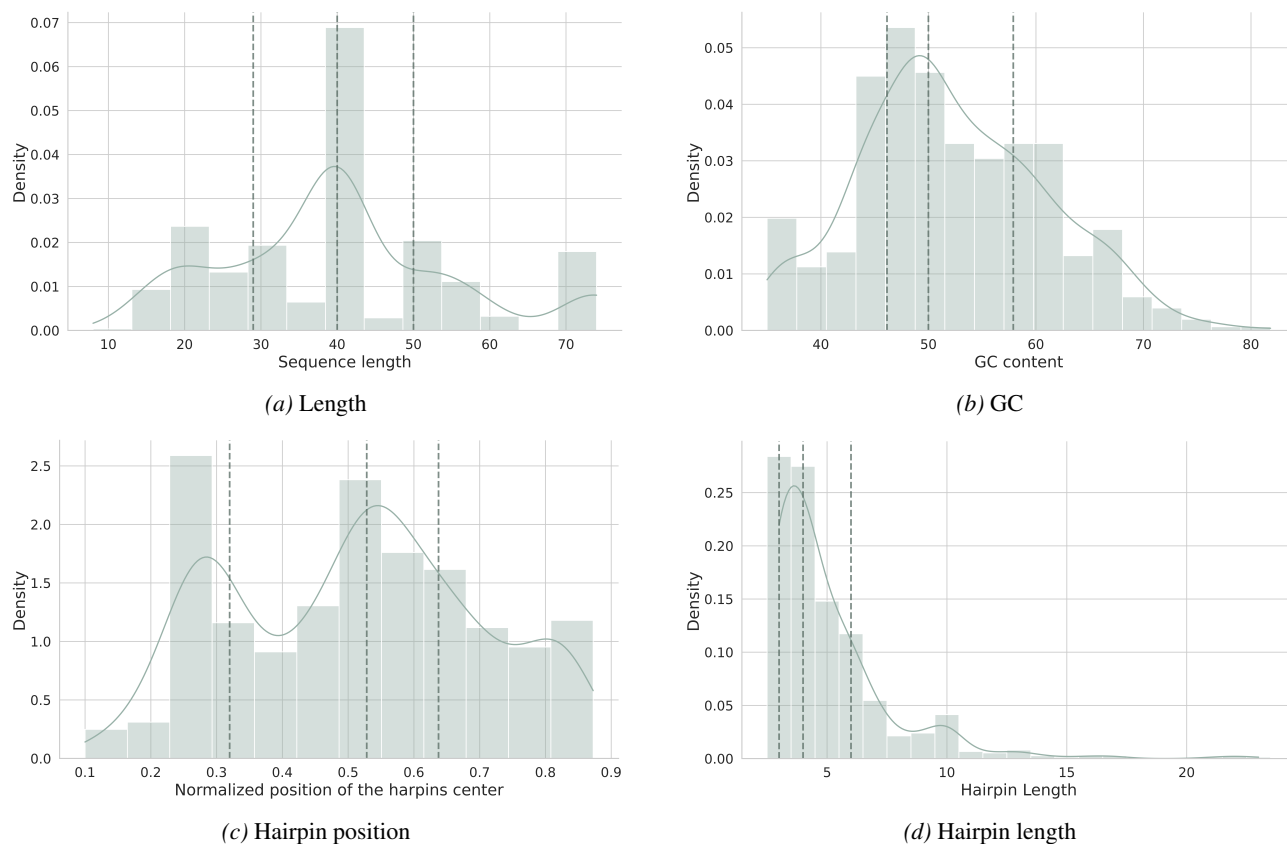


Figure 1. Structural parameters of Sequence Craft dataset.

distribution-based approach requires heuristics or hardware acceleration for practical use.

---

**Algorithm 2** Distribution based approach
 

---

**Load Dataset:**  $D_{ref} \leftarrow$  Sequence Craft  
**Extract Target Distributions:**  $T \leftarrow \{GC, Entropy, Levenshtein\}$   
**Initialize:** Population  $P$  with sampled sequence lengths  
**while** Total sequences  $< N$  **and**  $gen < 30$  **do**  
  **for** each sequence  $s \in P$  **do**  
    Compute  $GC(s)$ , entropy, and mean Levenshtein distance  
    Evaluate Wasserstein distance to  $T$   
  **end for**  
  Tournament selection  
  Two-point crossover and shuffle-index mutation  
  Preserve best-performing individuals  
   $gen \leftarrow gen + 1$   
**end while** Synthetic sequences matching  $D_{ref}$

---

**Classifier**

The training set consisted of 161 experimentally validated DNAzymes with RNA cleavage rates above  $10^{-7} \text{ min}^{-1}$  and 1,000 randomly generated inactive sequences. The class

imbalance was intentionally introduced to minimize false positives and improve selection reliability.

Among all tested configurations, LightGBM with DNA-BERT embeddings achieved the best performance. Generalization was further assessed on an independent subset of the Sequence Craft dataset excluded from training. On this benchmark, the model achieved a ROC-AUC of 0.935, MCC of 0.788, Precision of 0.919, Recall of 0.845, and an F1-score of 0.880, with only 11 false-positive predictions. These results confirm the robustness and discriminative capability of the selected classifier.

**Regressor**

As an additional validation stage, generated sequences passing the classifier filter were evaluated using the regression model introduced in the Sequence Craft study. The model incorporated sequence information together with environmental parameters, including temperature, cofactor, and buffer composition. To ensure consistency, experimental conditions were fixed to the most common settings observed in the original dataset.

**Clustering**

To predict metal ion specificity, a  $k$ -nearest neighbors

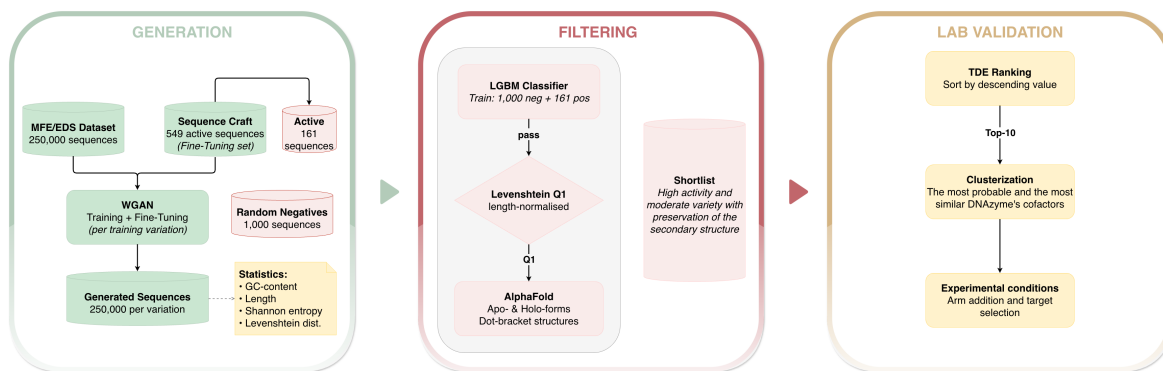


Figure 2. Hierarchical screening pipeline.

(KNN) classifier was trained on HyenaDNA embeddings of Sequence Craft sequences. The dataset was divided into 449 training and 100 test samples.

The model achieved a micro-averaged F1-score, Precision, and Recall of 0.740, indicating strong performance on dominant classes. However, macro-averaged metrics were substantially lower (F1 = 0.505, Precision = 0.523, Recall = 0.504), demonstrating reduced performance on minority classes. To compensate for this imbalance, predicted cofactors were additionally compared with those of the nearest known sequences.

Predicted metal ion labels and corresponding confidence scores were appended to each generated sequence and exported for downstream analysis.

### 2.2.2. SCREENING PIPELINE

A hierarchical screening pipeline was developed to reduce and prioritize generated DNAzyme candidates (Figure 2). The workflow sequentially applied: (i) classification filtering, (ii) Levenshtein distance stratification, (iii) secondary-structure filtering, and (iv) Tree Distance Editing (TDE) ranking.

The classifier served as the primary filtering stage and removed the majority of low-probability candidates. The remaining sequences were stratified by normalized Levenshtein distance to experimentally validated DNAzymes. This step balanced novelty and functional plausibility by excluding both trivial near-duplicates and highly divergent sequences. Candidates from the first quartile were retained as the closest representatives of the known functional space.

Secondary-structure filtering was then applied using rules derived from the Sequence Craft dataset. Sequences composed entirely of unpaired nucleotides were excluded due to the absence of a stable catalytic core. Candidates containing stems longer than 9 nucleotides were also removed, since excessive stabilization may reduce catalytic flexibility and increase misfolding risk. In addition, normalized hairpin

coordinates were constrained to the 5<sup>th</sup>–95<sup>th</sup> percentile interval (0.232–0.831), preserving the structural organization characteristic of functional DNAzymes.

Finally, candidates were ranked using TDE. Lower TDE scores indicate secondary structures more similar to experimentally validated DNAzymes and therefore higher structural plausibility.

The multi-stage design reduces false positives associated with any single metric and ensures that selected candidates satisfy probabilistic, sequence-level, and structural constraints simultaneously.

### 2.2.3. THE COMPARATIVE ANALYSIS OF APPROACHES TO TRAINING

Five WGAN training configurations were compared to evaluate the effect of synthetic pre-training on DNAzyme generation. All results of the comparison are captured in Table 3. The main objective was to determine whether pre-training on synthetic datasets generated by either MFE optimization or evolutionary distribution synthesis (EDS) improves performance relative to training only on experimentally validated DNAzymes. Generated sequences from each configuration were processed through the screening pipeline, and the resulting statistics are summarized in Tables 1 and 2.

Among all configurations, EDS + FT showed the closest agreement with the validated DNAzyme distribution. It achieved the lowest Jensen–Shannon divergence (JSD) values for GC-content (0.006) and entropy (0.005), demonstrating that statistically matched synthetic pre-training effectively captures the compositional properties of functional DNAzymes.

Fine-tuning substantially improved both synthetic pre-training strategies. For MFE pre-training, fine-tuning reduced GC-content divergence from 0.100 to 0.042. In the EDS pipeline, fine-tuning reduced 3-mer JSD from 0.028 to 0.004, indicating improved local motif fidelity. Structural diversity remained stable across all models, with APJD and

Table 1. Effect of screening stages on candidate counts.

Method	EDS FT	MFE FT	NS	EDS	MFE	Sequence Craft
Classifier	1462	947	520	367	343	2056
Levenshtein	422	254	132	94	99	516
Dot-bracket	332	220	108	79	76	412

Table 2. Metrics after screening pipeline.

Setup	APJD	MPD	JSD (3-mer)	JSD (GC)	JSD (Entropy)
MFE + FT	0.702	0.682	0.014	0.044	0.035
EDS + FT	0.701	0.670	0.017	0.016	0.059
Negative	0.718	0.717	0.007	0.034	0.015

MPD values consistently ranging between 0.70 and 0.73.

Negative samples (RS) were used as a control for screening performance. Surprisingly, models trained only on synthetic datasets without fine-tuning did not outperform random sequences in the number of candidates passing the screening stages. In contrast, both fine-tuned configurations produced substantially more high-quality candidates, confirming the importance of adaptation to experimentally validated DNAzymes.

The screening pipeline reduced differences between models by selecting structurally and compositionally conservative candidates. After filtering, diversity metrics decreased for all setups, although MFE + FT retained higher structural diversity (MPD 0.682 vs 0.670 for EDS + FT).

EDS + FT maintained the best GC-content agreement (JSD 0.016), but no longer dominated all metrics after screening. MFE + FT showed lower divergence for 3-mer distributions and entropy, indicating a more balanced trade-off between compositional realism and diversity.

Overall, synthetic pre-training consistently improved generation quality compared to training only on scarce experimental data. EDS + FT achieved the highest statistical fidelity to validated DNAzymes, while MFE + FT preserved greater structural diversity. Together, these results identify EDS + FT and MFE + FT as the most effective and complementary configurations for DNAzyme generation.

Visualizing the sequence space with Principal Component Analysis (PCA) on HyenaDNA embeddings (Figure 3) provides qualitative validation of the statistical trends described above and enables comparison between generated and natural DNAzyme distributions.

All models produced broad continuous clusters without collapse into isolated regions, indicating the absence of overfitting and successful learning of general DNAzyme patterns rather than memorization of training sequences.

The EDS + FT (Fig. 3a) and MFE + FT (Fig. 3b) configurations showed the strongest overlap with the natural DNAzyme distribution. Sequences passing the screening pipeline were concentrated within dense regions of the natural embedding space, confirming the effectiveness of the filtering strategy.

The two leading models demonstrated complementary be-

havior. EDS + FT showed the closest agreement with the empirical distribution, consistent with its low JSD values. In contrast, MFE + FT displayed broader dispersion in PCA space, reflecting higher structural diversity and wider exploration of thermodynamic configurations. Negative Samples formed sparse and disorganized clusters, highlighting the inefficiency of random sequence generation.

### 3. Conclusion

The systematic comparison of WGAN training strategies demonstrates that pre-training on large-scale synthetic datasets decisively outperforms training solely on the scarce collection of experimentally validated DNAzymes. Pre-trained models achieved Jensen–Shannon divergences up to an order of magnitude lower than the DNAzyme-only baseline for both global sequence properties and local motif distributions, confirming that synthetic data provide a robust statistical prior that effectively compensates for the limited availability of real catalytic sequences.

From an algorithmic perspective, the two pre-training strategies differ meaningfully in their data-generation cost. The MFE-based dataset is constructed using thermodynamic folding calculations with a time complexity of  $\mathcal{O}(N \cdot L^3)$ , which for short DNAzyme sequences ( $L \sim 40$ ) is an extremely rapid operation. In contrast, the EDS pipeline requires computing Levenshtein distances to all reference sequences, scaling as  $\mathcal{O}(N \cdot M \cdot L^2)$  and incurring a substantially larger computational burden as the reference set  $M$  grows. Thus, the MFE-based approach offers a compelling practical advantage: it produces a large, structurally grounded training corpus notably faster, making it the preferred option when computational efficiency and thermodynamic stability are the primary concerns.

Despite this speed benefit, the EDS-based strategy, which explicitly matches the property distributions of the Sequence-Craft dataset and incorporates subsequent fine-tuning, ultimately yielded candidates that are statistically indistinguishable from natural DNAzymes across all examined metrics, establishing EDS + FT as the method of choice for maximum biological realism. The hierarchical screening pipeline developed in this work preserved the relative quality ranking of the models: the superior fidelity of pre-trained variants persisted after filtering, with EDS + FT continuing to deliver the most realistic candidates. This confirms that the pipeline

330 acts as a robust amplifier of inherent candidate quality rather  
 331 than distorting model comparisons, while simultaneously  
 332 ensuring that only the most promising sequences advance  
 333 to resource-intensive experimental validation.

334 In summary, the synergy between computationally ef-  
 335 ficient pre-training (MFE), high-fidelity distribution-  
 336 matching (EDS), targeted fine-tuning, and a principled  
 337 multi-metric screening workflow constitutes a versatile,  
 338 resource-efficient, and generalizable framework for the gen-  
 339 erative design of functional nucleic acids.  
 340

## 341 Accessibility

342 All code exists in the repository:  
 343 <https://anonymous.4open.science/r/DnaZymeXHSE-55B5/>  
 344  
 345  
 346

## 347 Software and Data

348 All data exists in the repository:  
 349 <https://anonymous.4open.science/r/DnaZymeXHSE-55B5/>  
 350  
 351  
 352

## 353 References

- 354 1. Breaker, R. R. & Joyce, G. F. A DNA enzyme that cleaves  
 355 RNA. *Chem. Biol.* **1**, 223–229 (1994).  
 356
- 357 2. Hollenstein, M. DNA Catalysis: The Chemical Reper-  
 358 toire of DNAzymes. *Molecules* **20**, 20777–20804 (2015).  
 359
- 360 3. Forty Years of In Vitro Evolution - Joyce - 2007 - Ange-  
 361 wandte Chemie International Edition - Wiley Online Library.  
 362 <https://onlinelibrary.wiley.com/doi/full/10.1002/anie.200701369>  
 363
- 364 4. Silverman, S. K. Catalytic DNA: Scope, Applications,  
 365 and Biochemistry of Deoxyribozymes. *Trends Biochem. Sci.*  
 366 **41**, 595–609 (2016).  
 367
- 368 5. DNAzyme technology and cancer therapy:  
 369 cleave and let die | Molecular Cancer Therapeu-  
 370 tics | American Association for Cancer Research.  
 371 [https://aacrjournals.org/mct/article/7/2/243/93030/DNAzyme-  
 372 technology-and-cancer-therapy-cleave-and](https://aacrjournals.org/mct/article/7/2/243/93030/DNAzyme-technology-and-cancer-therapy-cleave-and).
- 373 6. Peracchi, A. DNA catalysis: potential, limitations, open  
 374 questions. <https://doi.org/10.1002/cbic.200500098> (2005)  
 375 doi:10.1002/cbic.200500098.  
 376
- 377 7. Santoro, S. W. & Joyce, G. F. A general purpose RNA-  
 378 cleaving DNA enzyme. *Proc. Natl. Acad. Sci.* **94**, 4262–  
 379 4266 (1997).
- 380 8. DNAzymes and cardiovascular disease - Benson - 2008 -  
 381 British Journal of Pharmacology - Wiley Online Library.  
 382 <https://bpspubs.onlinelibrary.wiley.com/doi/full/10.1038/bjp.2008.145>  
 383  
 384

9. Inhibition of bcr-abl Oncogene Expression by Novel  
 Deoxyribozymes (DNAzymes) | Human Gene Therapy.  
<https://www.liebertpub.com/doi/abs/10.1089/10430349950016573>.

10. Goila, R. & Banerjea, A. C. Sequence specific cleav-  
 age of the HIV-1 coreceptor CCR5 gene by a hammer-head  
 ribozyme and a DNA-enzyme: inhibition of the corecep-  
 tor function by DNA-enzyme. *FEBS Lett.* **436**, 233–238  
 (1998).

11. Yan, J., Ran, M., Shen, X. & Zhang, H. Ther-  
 apeutic DNAzymes: From Structure Design to Clinical  
 Applications. <https://doi.org/10.1002/adma.202300374>  
 doi:10.1002/adma.202300374.

12. Gata-3-specific Dnazyme As An Approach For  
 Asthma-therapy - Journal of Allergy and Clinical  
 Immunology. [https://www.jacionline.org/article/S0091-  
 6749\(06\)02366-9/fulltext](https://www.jacionline.org/article/S0091-6749(06)02366-9/fulltext).

13. Krug, N. *et al.* Allergen-Induced Asthmatic Responses  
 Modified by a GATA3-Specific DNAzyme. *N. Engl. J. Med.*  
**372**, 1987–1995 (2015).

14. Caramori, G., Chung, K. F. & Barnes, P. J. Allergen  
 Responses Modified by a GATA3 DNAzyme. *N. Engl. J.*  
*Med.* **373**, 1176–1177 (2015).

15. Grimpe, B. *et al.* The critical role of basement  
 membrane-independent laminin gamma 1 chain during axon  
 regeneration in the CNS. *J. Neurosci. Off. J. Soc. Neurosci.*  
**22**, 3144–3160 (2002).

16. Hjiantonou, E., Iseki, S., Uney, J. B. & Phylactou, L.  
 A. DNAzyme-mediated cleavage of Twist transcripts and  
 increase in cellular apoptosis. *Biochem. Biophys. Res.*  
*Commun.* **300**, 178–181 (2003).

17. RNA cleaving ‘10-23’ DNAzymes with enhanced stabil-  
 ity and activity | Nucleic Acids Research | Oxford Academic.  
<https://academic.oup.com/nar/article/31/20/5982/1039490>.

18. LNAzymes: Incorporation of LNA-Type  
 Monomers into DNAzymes Markedly Increases RNA  
 Cleavage | Journal of the American Chemical Society.  
<https://pubs.acs.org/doi/full/10.1021/ja0276220>.

19. Locked nucleoside analogues expand  
 the potential of DNAzymes to cleave struc-  
 tured RNA targets | BMC Molecular Biology.  
[https://link.springer.com/article/10.1186/1471-2199-  
 7-19](https://link.springer.com/article/10.1186/1471-2199-7-19).

20. Bacterial Regulatory RNA: Methods and Protocols |  
 SpringerLink. [https://link.springer.com/book/10.1007/978-  
 1-61779-949-5](https://link.springer.com/book/10.1007/978-1-61779-949-5).

21. Rapid in vitro Method for Obtaining RNA Acces-  
 sibility Patterns for Complementary DNA Probes: Cor-  
 relation with an Intracellular Pattern and Known RNA

- 385 Structures | Nucleic Acids Research | Oxford Academic.  
386 <https://academic.oup.com/nar/article/25/24/5010/1748765>.
- 387 22. SHAPE-Seq 2.0: systematic optimization and  
388 extension of high-throughput chemical probing of  
389 RNA secondary structure with next generation sequenc-  
390 ing | Nucleic Acids Research | Oxford Academic.  
391 <https://academic.oup.com/nar/article/42/21/e165/2903199>.
- 392 23. Selective 2'-hydroxyl acylation analyzed by primer  
393 extension (SHAPE): quantitative RNA structure analy-  
394 sis at single nucleotide resolution | Nature Protocols.  
395 <https://www.nature.com/articles/nprot.2006.249>.
- 396 24. Ponce-Salvatierra, A., Boccaletto, P. & Bujnicki, J. M.  
397 DNAmoreDB, a database of DNAzymes. *Nucleic Acids Res.*  
398 **49**, D76–D81 (2021).
- 399 25. Eremeyeva, M., Din, Y., Shirokii, N. & Serov, N.  
400 SequenceCraft: machine learning-based resource for ex-  
401 ploratory analysis of RNA-cleaving deoxyribozymes. *BMC*  
402 *Bioinformatics* **26**, 2 (2025).
- 403 26. Watson, J. L. *et al.* Broadly applicable and accurate pro-  
404 tein design by integrating structure prediction networks and  
405 diffusion generative models. 2022.12.09.519842 Preprint at  
406 <https://doi.org/10.1101/2022.12.09.519842> (2022).
- 407 27. Protein generation with evolutionary dif-  
408 fusion: sequence is all you need | bioRxiv.  
409 <https://www.biorxiv.org/content/10.1101/2023.09.11.556673v1>.
- 410 28. Breaker, R. R. In vitro selection of catalytic polynu-  
411 cleotides. *Chem. Rev.* **97**, 371–390 (1997).
- 412 29. Silverman, S. K. Deoxyribozymes: selection design  
413 and serendipity in the development of DNA catalysts. *Acc.*  
414 *Chem. Res.* **42**, 1521–1531 (2009).
- 415 30. Zhao, Z. *et al.* Structural basis for catalysis by a  
416 DNAzyme. *Nat. Chem.* **8**, 684–690 (2016).
- 417 31. Chen, J. *et al.* Single-molecule dynamics of a DNAzyme  
418 reveals hidden states. *Nat. Commun.* **8**, 1–11 (2017).
- 419 32. Herschlag, D. RNA chaperones and the RNA folding  
420 problem. *J. Biol. Chem.* **270**, 20871–20874 (1995).
- 421 33. Jencks, W. P. *Catalysis in Chemistry and Enzymology*.  
422 (Dover Publications, 1987).
- 423 34. Fedor, M. J. & Williamson, J. R. The catalytic diversity  
424 of RNAs. *Nat. Rev. Mol. Cell Biol.* **6**, 399–412 (2005).
- 425 35. Liu, J. & Lu, Y. Rational design of functional DNA-  
426 based metal sensors. *Angew. Chem. Int. Ed.* **45**, 90–94  
427 (2006).
- 428 36. Li, J. & Breaker, R. R. Kinetics of RNA degradation by  
429 specific DNA enzymes. *J. Am. Chem. Soc.* **121**, 5364–5372  
430 (1999).
- 431 37. Peracchi, A. Origins of the temperature dependence of  
432 ribozyme catalysis. *Biochemistry* **44**, 10542–10552 (2005).
- 433 38. Takahashi, S., Hamada, M., Tateishi-Karimata, H. &  
434 Sugimoto, N. Fitness landscapes and thermodynamic ap-  
435 proaches to development of nucleic acids enzymes: from  
436 classical methods to AI integration. *RSC Chem. Biol.* **6**,  
437 1667–1685 (2025).
- 438 39. Zhang, F., Shi, W., Guo, L., Liu, S. & He, J. The Pro-  
439 grammable Catalytic Core of 8-17 DNAzymes. *Molecules*  
440 **29**, 2420 (2024).
- 441 40. Hollenstein, M., Hipolito, C., Lam, C., Dietrich, D. &  
442 Perrin, D. M. A self-cleaving DNA enzyme modified with  
443 amines, guanidines and imidazoles operates independently  
444 of divalent metal cations (M<sup>2+</sup>). *Nucleic Acids Res.* **37**,  
445 1638–1649 (2009).

## A. Appendix

## B. Additional tables

Table 3. Performance metrics of WGAN architectures.

Training setup	APJD	MPD	JSD (3-mer)	JSD (GC)	JSD (Entropy)
RS	0.717	0.709	0.013	0.035	0.023
Sequence Craft	0.720	0.714	0.013	0.041	0.019
MFE	0.728	0.706	0.028	0.100	0.019
MFE + FT	0.729	0.729	0.015	0.042	0.016
EDS	0.728	0.706	0.028	0.099	0.019
EDS + FT	0.709	0.705	0.004	0.006	0.005

## C. Additional figures

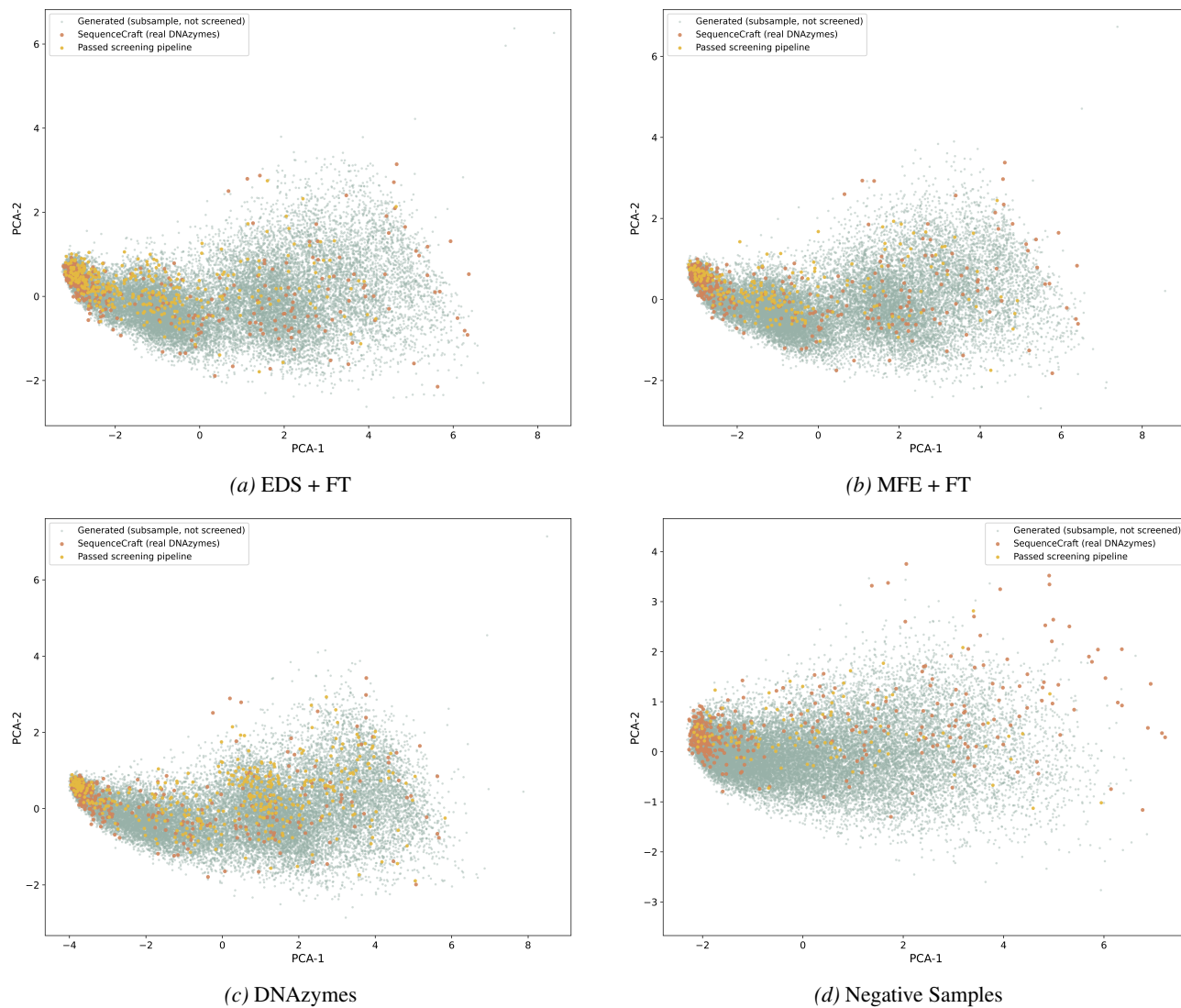


Figure 3. PCA projection of generated and reference DNAzyme sequences using HyenaDNA embeddings.