
BIOREASON: Incentivizing Multimodal Biological Reasoning within a DNA-LLM Model

Adibvafa Fallahpour*^{1,2,3,5} **Andrew Magnuson***^{1,2}
adibvafa.fallahpour@mail.utoronto.ca andrew.magnuson@mail.utoronto.ca

Purav Gupta*^{1,2} **Shihao Ma**^{1,2,3}
purav.gupta@mail.utoronto.ca shihao.ma@mail.utoronto.ca

Jack Naimer^{1,2,3} **Arnav Shah**^{1,2,3}
jack.naimer@mail.utoronto.ca arnav.shah@mail.utoronto.ca

Haonan Duan^{1,2} **Omar Ibrahim**³
haonan.duan@mail.utoronto.ca omar.ibrahim2@uhn.ca

Hani Goodarzi^{†4,6} **Chris J. Maddison**^{†1,2,7}
hani.goodarzi@ucsf.edu cmaddis@cs.toronto.edu

Bo Wang^{†1,2,3}
bowang@vectorinstitute.ai

¹University of Toronto ²Vector Institute ³University Health Network (UHN)
⁴Arc Institute ⁵Cohere ⁶University of California, San Francisco ⁷Google DeepMind

Abstract

Unlocking deep, interpretable biological reasoning from complex genomic data is a major AI challenge hindering scientific discovery. Current DNA foundation models, despite strong sequence representation, struggle with multi-step reasoning and lack inherent transparent, biologically intuitive explanations. We introduce **BIOREASON**, a pioneering architecture that, for the first time, deeply integrates a DNA foundation model with a large language model (LLM). This novel connection enables the LLM to directly process and reason with genomic information as a fundamental input, fostering a new form of multimodal biological understanding. **BIOREASON**'s sophisticated multi-step reasoning is developed through supervised fine-tuning and targeted reinforcement learning, guiding the system to generate logical, biologically coherent deductions. Across biological reasoning benchmarks, **BIOREASON** significantly improves performance, raising accuracy on KEGG-based disease pathway prediction from 86% to 98% and delivering an average 15% gain over strong single-modality baselines in variant effect prediction tasks. **BIOREASON** reasons over unseen biological entities and articulates decision-making through interpretable, step-by-step biological traces, offering a transformative approach for AI in biology that enables deeper mechanistic insights and accelerates testable hypothesis generation from genomic data. Data, code, and checkpoints are publicly available at <https://github.com/bowang-lab/BioReason>.

*Equal contribution. †Equal advising.

1 Introduction

Biological data, spanning genomics, transcriptomics, biomedical literature, and more, is expanding at an unprecedented rate, creating immense opportunities for scientific discovery. This data explosion has catalyzed the development of foundation models (FMs), deep networks trained on vast datasets that enable a wide array of downstream tasks. In genomics, DNA foundation models [6, 10, 31, 45, 13] have demonstrated remarkable capabilities by learning dense sequence representations that drive splice site identification, variant effect prediction, and regulatory element characterization.

Despite these advances, a critical challenge with foundation models still persists: effectively translating these learned representations into mechanistic insights and falsifiable hypotheses. Current DNA foundation models, while powerful in their representational capacity, typically function as "black boxes" that lack the inherent ability to generate transparent, biologically intuitive explanations [4, 26]. These limitations are prominent in complex biological problems requiring mechanistic understanding, such as gene pathway analysis, phenotype prediction, and disease mechanism elucidation [9].

Large language models (LLMs) [32, 2, 12, 34] have rapidly advanced in reasoning capabilities, problem-solving, and knowledge depth. Through sophisticated training methods including reinforcement learning and supervised fine-tuning, these models demonstrate increasingly sophisticated multi-step reasoning across domains from mathematical problem-solving to logical deduction [14, 28, 30, 18]. However, LLMs alone lack the specialized architecture to effectively process raw genomic sequences and often fail to capture nuanced biological patterns in genetic data.

This disconnect between powerful sequence representations of DNA foundation models and sophisticated reasoning capabilities of LLMs creates a significant barrier to developing AI systems that provide deep mechanistic insights comparable to biology domain experts. To bridge this gap, we present BIOREASON: a novel architecture that fundamentally integrates a DNA foundation model with an LLM, enabling a new paradigm of multimodal biological understanding and reasoning.

BIOREASON is distinguished by its ability to create a unique flow of information between genomic and natural language. This architecture enables the system to process raw DNA sequences while leveraging the reasoning capabilities of modern LLMs to generate biologically coherent explanations and predictions. Through a training methodology combining supervised fine-tuning and reinforcement learning, BIOREASON develops the capacity for sophisticated multi-step reasoning over genomic data; a capability that neither DNA foundation models nor LLMs can achieve independently.

Contributions. Our key contributions include:

- **Novel multimodal architecture.** The first successful integration of a large DNA foundation model with an LLM, establishing a new methodology for AI-driven biological studies.
- **Advanced reasoning.** A systematic training approach combining supervised fine-tuning and reinforcement learning that incentivizes multi-step biological reasoning.
- **New biological reasoning benchmarks.** Development and curation of novel benchmarks for evaluating biological reasoning capabilities, including an annotated reasoning dataset for gene pathway and disease prediction dataset from KEGG [20].
- **Empirical performance improvements.** Demonstration that BIOREASON outperforms both DNA foundation models and LLMs with average performance gains of 15%+ over baseline.
- **Interpretable reasoning traces.** A mechanism for generating step-by-step biological reasoning traces that provide interpretable predictions, enhancing scientific insight and hypothesis generation

2 Background & Related Work

2.1 DNA Foundation Models

Recent years have witnessed the emergence of DNA foundation models [6, 10, 43, 31] that have significantly accelerated discovery throughout the biological sciences. These models extract meaningful representations directly from nucleotide sequences by pre-training on vast genomic datasets.

Moreover, comprehensive benchmarking studies [15] have demonstrated their proficiency across various genomics tasks in both zero-shot and fine-tuned settings.

Evo2 [6], in particular, represents a significant advancement as one of the largest genomic foundation models to date, enabling extremely long-range context windows and predictions. Its ability to generate complete bacterial and yeast genomes underscores the potential of these models to capture complex genomic patterns [22]. However, a critical limitation persists: these foundation models operate as "black boxes," lacking the interpretability necessary to explain how they derive conclusions from their embeddings. This opacity hampers the advancement of biological knowledge by obscuring the mechanistic insights that could otherwise be derived from model predictions.

2.2 Large Language Models for Biological Reasoning

LLMs have demonstrated remarkable capabilities in understanding and generating human-like text, with substantial success in interpreting and reasoning over complex biomedical data. Recent reviews [43] highlight their success across diverse domains, from clinical applications involving patient notes to biological research contexts. The development of specialized models pre-trained on biomedical literature [25], has further enhanced their domain-specific performance.

Genomics-focused LLMs such as GeneGPT [19], agentic models such as TxGemma [40], and scientific reasoning models such as rbiol [18] represent initial attempts to integrate language models with genomic databases. ChatNT [35] takes a step further by integrating DNA foundation model representations with language models in a multimodal framework. However, no previous work has trained such systems for complex biological reasoning. We present the first multimodal framework designed and trained to perform biological reasoning by combining textual knowledge from LLMs with nucleotide-level representations from DNA foundation models.

2.3 Genomics Benchmarks

DNA foundation models are typically evaluated on established benchmarks encompassing diverse prediction tasks, including regulatory element identification, variant effect prediction, transcription factor binding site prediction, and splice site classification. Comprehensive benchmarking frameworks like BEND [29] provide standardized evaluation protocols that enable meaningful comparisons between models across these supervised tasks.

While these benchmarks effectively measure performance on specific downstream applications, they inadequately evaluate a model’s capacity for higher-order reasoning or hypothesis generation, capabilities essential for advancing scientific understanding. This represents a critical conceptual gap between current evaluation metrics and the sophisticated reasoning abilities desired from next-generation foundation models. The field requires benchmarks that challenge models to perform multi-step logical reasoning and predict potential biological mechanisms. This need motivated our curation of the KEGG pathway database [20] to create a multi-step reasoning, variant effect prediction dataset that specifically evaluates a model’s capacity for mechanistic biological reasoning.

3 BioReason Model

We introduce BIOREASON, a multimodal framework designed to unlock deep, interpretable biological reasoning by synergistically integrating genomic and language data. BIOREASON operates on two primary input streams: (i) one or more genomic sequences, denoted S_{DNA} ; and (ii) textual queries, Q_{TEXT} . These queries are processed by an LLM-specific tokenizer, $T_{\text{LLM}}(\cdot)$, into a sequence of M tokens (w_1, \dots, w_M) from the LLM’s vocabulary \mathcal{V}_{LLM} . Current methods often fall short in this domain: LLMs treat S_{DNA} as simple strings, thereby missing rich genomic features, while DNA Foundation Models (f_{DNA}) capture these features but primarily yield task-specific discriminative outputs (e.g., classification or regression scores) rather than interpretable natural language. BIOREASON bridges this gap by deriving contextualized DNA embeddings from the S_{DNA} input(s) and integrating them with the tokenized Q_{TEXT} to form a unified multimodal input sequence, X_{LLM} , for its core LLM. This direct integration enables the generation of explanatory text, $Y_{\text{OUT}} = (y_1, \dots, y_K)$, grounded in genomic nuances. The output Y_{OUT} presents biological reasoning and the final response. Figure 1 depicts the overall architecture.

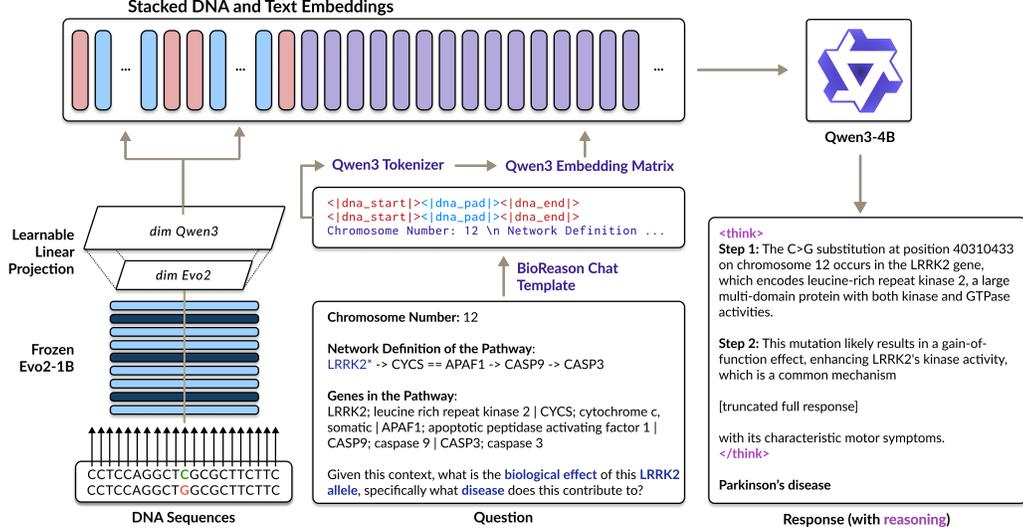


Figure 1: **BIOREASON Architecture**. Schematic representation of our novel multimodal framework that integrates a DNA foundation model with a Large Language Model.

3.1 DNA Foundation Model (f_{DNA}) Encoder

f_{DNA} transforms each input S_{DNA} sequence into contextualized embeddings. We utilize established DNA foundation models such as StripedHyena2 (e.g., Evo2) [6, 33], or the Nucleotide Transformer (NT), [11], as the f_{DNA} . Each S_{DNA} is first processed by its respective DNA-specific tokenizer, $T_{\text{DNA}}(\cdot)$, which segments it into a sequence of L' DNA tokens, $D = (d_1, \dots, d_{L'})$; each token d_j can represent one or more nucleotides. If an input S_{DNA} sequence, after tokenization by T_{DNA} , exceeds a defined context length (e.g., 2048 DNA tokens), it is truncated. The chosen f_{DNA} architecture then maps each token sequence D to a sequence of high-dimensional per-token embeddings $E_{\text{DNA}} = (e_1, \dots, e_{L'}) \in \mathbb{R}^{L' \times d_{\text{dna}}}$. These d_{dna} -dimensional embeddings capture context-dependent genomic features. The weights of the f_{DNA} are kept frozen during BIOREASON’s training and inference.

3.2 Large Language Model (f_{LLM}) Backbone

The f_{LLM} is the primary reasoning engine and text generator. We employ Qwen3 [41, 42], an autoregressive Transformer-based LLM, initialized with its original pre-trained weights. This model receives the multimodal input sequence X_{LLM} and is trained to predict the next token y_i in the sequence Y_{OUT} , conditioned on the preceding tokens $y_{<i}$ and X_{LLM} . Mathematically, we optimize the parameters θ_{LLM} of the f_{LLM} by maximizing the log-likelihood of the observed sequences:

$$\mathcal{L} = \sum_i \log P(y_i | y_{<i}, X_{\text{LLM}}; \theta_{\text{LLM}}) \quad (1)$$

The f_{LLM} utilizes special tokens to structure conversational interactions and reasoning within its textual output Y_{OUT} . These include tokens defining user and assistant roles (e.g., `<|im_start|>user/assistant ...<|im_end|>`), structuring reasoning steps (e.g., `<think> ...</think>`), alongside standard padding tokens (`<|endoftext|>`).

3.3 Multimodal Genomic Integration

Genomic information, as DNA embeddings from f_{DNA} , is integrated into the f_{LLM} ’s input by stacking these with embeddings of the user’s query Q_{TEXT} and special tokens such as `<dna_start>` and `<dna_end>`. Key to this integration is the preparation of the DNA embedding block, E'_{DNA} , formed from one or more input DNA sequences. For each sequence $S_{\text{DNA},k}$, its f_{DNA} -generated embedding sequence $E_{\text{DNA},k} \in \mathbb{R}^{L'_k \times d_{\text{dna}}}$ (where L'_k is its tokenized length) is projected by a learnable linear layer, $\text{Proj} : \mathbb{R}^{d_{\text{dna}}} \rightarrow \mathbb{R}^{d_{\text{lm}}}$, to yield $E'_{\text{DNA},k}$ of dimension d_{lm} . The resulting E'_{DNA} block is obtained by stacking all $E'_{\text{DNA},k}$ sequences along the sequence dimension.

Concurrently, the user’s tokenized query $Q_{\text{TEXT}} = (w_1, \dots, w_M)$ is mapped to its embedding sequence $\mathbf{E}_{Q_{\text{text}}} = (E(w_1), \dots, E(w_M))$ by the f_{LLM} ’s input embedding layer, $E(\cdot)$. Similarly, the special tokens `<dna_start>` and `<dna_end>` are embedded via $E(\cdot)$ to produce $e_{\text{<dna_start>}}$ and $e_{\text{<dna_end>}}$. These components are then stacked along the sequence dimension to form the multimodal input X_{LLM} for f_{LLM} .

$$X_{\text{LLM}} = (e_{\text{<dna_start>}}, \mathbf{E}'_{\text{DNA}}, e_{\text{<dna_end>}}, \mathbf{E}_{Q_{\text{text}}}) \quad (2)$$

All constituent embedding vectors within X_{LLM} receive positional information via Rotary Position Embedding (RoPE) [39], applied according to their final sequence positions. This strategy enables f_{LLM} fine-grained attention over both genomic and textual components within a unified modality.

3.4 Group Relative Policy Optimization (GRPO)

To further enhance BIOREASON ’s reasoning performance beyond supervised fine-tuning, we employ Group Relative Policy Optimization (GRPO) [36, 12], a reinforcement learning strategy tailored for refining reasoning generation in language models. GRPO leverages reward signals within groups of sampled outputs, eliminating the need for an explicit value estimator. We implement Dr. GRPO [27], an unbiased variant of GRPO that improves token efficiency while maintaining performance.

For the full formalism, including the composite reward design, advantage normalization, and the clipped surrogate objective with KL regularization, please refer to Appendix A.4.

4 Datasets

To develop a multimodal DNA-LLM model with reasoning capabilities, we curated three datasets: one novel dataset specifically designed to incentivize reasoning and two adapted from established benchmarks. The adapted datasets are derived from ClinVar [24] and OMIM [1], which are widely used for variant effect prediction tasks. Our novel dataset is based on KEGG Network Variants data [20] and enhanced with cross-linked metadata from several public variant repositories including ClinVar [24], OMIM [1], dbSNP [37], and COSMIC [38]. This novel dataset relies on the high-quality manual annotations and descriptions from the curators of KEGG, for gene pathway descriptions and downstream phenotypic effects, like disease.

4.1 KEGG-Derived Biological Reasoning Dataset

4.1.1 Dataset Integration and Statistics

We present a high-quality biological reasoning dataset derived from the Kyoto Encyclopedia of Genes and Genomes (KEGG) pathway database [20], consisting of 1,449 entries that elucidate the mechanistic connections between genetic variants and disease phenotypes. As seen in Figure 2, the dataset construction involved a rigorous multi-stage process that integrates structured pathway information with variant data to enable step-by-step reasoning across molecular networks.

For primary data integration, we extracted pathway network data from KEGG [20], focusing on disease-associated molecular interactions. Pathway data was augmented with variant information from clinical databases (ClinVar, dbSNP, OMIM, COSM) [24, 37, 1, 38] through a semi-automated mapping protocol [17, 21] that preserved relational integrity between genomic loci and functional elements within pathways. Each molecular network was represented using a standardized symbolic notation (e.g., "GENE1+GENE2 -> GENE3 -| GENE4") that encapsulates interaction types including activation, inhibition, complex formation, and transcriptional regulation.

To support variant interpretation, we included paired reference and variant sequences with precise alignment coordinates. These sequences have an average length of approximately 4,000 base pairs, with most variants differing by only 1–3 nucleotides from their reference sequences.

4.1.2 Reasoning Path Construction and Curation

A distinctive feature of this dataset is its inclusion of explicit causal reasoning paths connecting genetic variants to disease phenotypes via defined molecular mechanisms; these paths were constructed

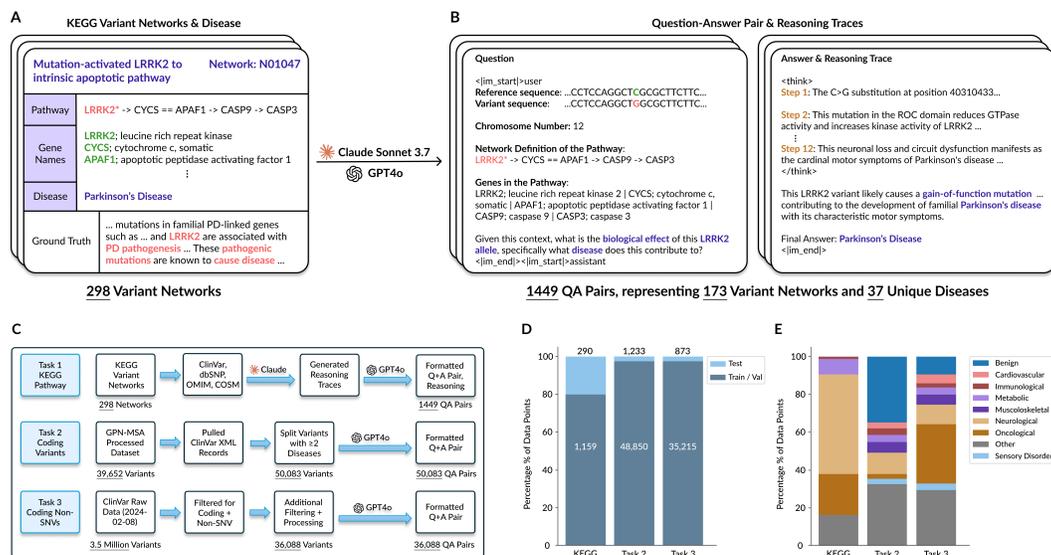


Figure 2: BIOREASON Dataset Curation and Composition. **A.** Representative example of a KEGG Variant Network element from the 298 networks utilized in our study, illustrating the relationship between genomic variants and their corresponding disease annotation that serves as ground truth for generating mechanistic reasoning traces. **B.** Exemplar of a structured question-answer pair with an accompanying multi-step reasoning trace demonstrating the expected logical progression from genomic variant to phenotypic outcome. **C.** Pipeline for data acquisition, integration, and curation across the three BIOREASON tasks. **D.** Distribution of train/test splits across the three curated datasets. 10% of train dataset was used for validation. **E.** Distribution of disease categories represented within the datasets, highlighting the diversity of variants and diseases represented in the datasets.

using the Claude 3.7 Sonnet model [2] and grounded with contextual disease information from the KEGG disease database [20]. For training and evaluation, the dataset is structured into standardized question-answer pairs: questions (illustrated in Figure 2B) incorporate variant details, network definitions, and gene descriptions, while answers provide concise mechanism-disease associations. The accompanying reasoning paths (mean length: 303.8 words) elaborate these mechanistic variant-to-phenotype links with precise molecular information.

4.2 Variant Effect Prediction of Coding Sequences

This dataset originated from the GPN-MSA [5] study. Affected gene names and disease phenotypes were extracted from ClinVar [24] XML records (via NCBI’s Entrez Direct tool [21]), while benign variants were sourced from gnomAD v3.1.2 [8] (requiring allele number $\geq 25,000$ and minor allele frequency (MAF) $> 5\%$). The data was split by chromosome (Chr 1–7, 9–22, X, Y for training; Chr 8 for testing). For training augmentation, GPT-4o [32] generated 50 semantically equivalent question variations per sample, prompting for pathogenic/benign classification and conditional disease phenotype prediction using chromosome and gene context; mutations linked to multiple diseases were treated as distinct samples for comprehensive phenotype coverage.

4.3 Variant Effect Prediction of Coding Non-SNVs

Coding non-SNVs were sourced from the ClinVar [24] database (2024-02-28 release). We filtered variants to retain only coding non-SNVs within the nuclear genome, affecting ≤ 64 base pairs, of certain significance, and with a review status of at least two stars, matched to GRCh38.p14 transcripts. After extracting affected gene names and disease phenotypes where available, a custom algorithm partitioned the dataset to ensure balanced disease representation in train/test splits. Finally, to augment training data, GPT-4o [32] generated 50 semantically equivalent question variations for each entry, using gene and chromosome number as context, prompting for pathogenic/benign classification and, if pathogenic, the associated disease phenotype.

5 Experiments

5.1 Datasets

BIOREASON’s performance is evaluated on three datasets (detailed in Section 4):

KEGG-Derived Biological Reasoning Dataset. This dataset (1,449 variants, 37 unique diseases) evaluates multi-step mechanistic reasoning. Input: paired reference and variant DNA sequences (S_{DNA}), and a textual query (Q_{TEXT}) with pathway/gene context. Task: predict the mutation’s effect and resulting disease by sequence-to-sequence generation of Y_{OUT} containing step-by-step reasoning between `<think>` (special tokens) and the final disease.

Variante Effect Prediction of Coding Sequences (VEP-Coding). Comprising 50,083 core variant entries, this dataset tests classifying coding variants. Input: paired reference and variant DNA sequences (S_{DNA}), and a textual query (Q_{TEXT}) providing gene and chromosome context. Task: sequence generation to predict if a variant is benign, or pathogenic with its associated disease. Split: Chromosomes (Chr) 1–7, 9–22, X, Y for train/validation; Chr 8 for testing.

Variante Effect Prediction of Coding Non-SNVs (VEP-Non-SNV). Containing 36,088 core non-SNV entries, this dataset addresses non-SNV alterations (e.g., indels <64 bp). Input: paired reference and variant DNA sequences (S_{DNA}), and an augmented textual query (Q_{TEXT}) providing gene and chromosome context. Task: sequence generation to predict if a non-SNV is benign, or pathogenic with its associated disease(s). We used stratified train/test splits to ensure balanced disease representation.

5.2 Models and Baselines

To benchmark BIOREASON’s performance, we evaluated it against several baseline models, categorized as DNA foundation models (f_{DNA}) and Large Language Models (f_{LLM}).

For f_{DNA} baselines, we utilized pre-trained Evo2-1B [6] and Nucleotide Transformer (NT-500M) [11] models. For downstream predictions, f_{DNA} models were adapted with an attention head where a single learnable query vector attends to the sequence token embeddings to produce a final sequence representation. For (f_{LLM}) baselines, we fine-tuned pre-trained Qwen3 models of two sizes: Qwen3-1.7B and Qwen3-4B [41, 42, 34]. These models were trained to receive text queries and DNA sequences treated as plain text strings and generate text with reasoning steps and final predictions.

The proposed BIOREASON models, were evaluated in several f_{DNA} and f_{LLM} combinations. Specifically, we tested Evo2-1B and NT-500M as f_{DNA} encoders, each paired with Qwen3-1.7B and Qwen3-4B as f_{LLM} backbones. The primary training methodology for all BIOREASON configurations was Supervised Fine-Tuning (SFT). Reinforcement Learning (RL) fine-tuning using the GRPO algorithm was subsequently applied to select DNA-LLM models.

5.3 Experimental Setup

Our experimental setup varied by model architecture—BIOREASON, LLM-only, or f_{DNA} -only—and task. BIOREASON and LLM-only models underwent Supervised Fine-Tuning (SFT), with LLM parameters efficiently updated via Low-Rank Adaptation (LoRA) [16]. For f_{DNA} -only baselines, core DNA model weights were frozen; only a task-specific attention head and classifier were trained.

SFT objectives for these models differed: for the KEGG Dataset Task, models generated reasoning steps between `<think>` tokens and a final disease prediction. For VEP Datasets Tasks, they aimed for pathogenic/benign classification and conditional disease prediction for pathogenic variants. During SFT, a specialized attention mask restricted loss computation exclusively to the response between `<think>` tokens and final answer tokens, excluding those from the input query or DNA embeddings. Select BIOREASON models were further optimized with GRPO. Details for LoRA configurations, all SFT and RL hyperparameters, and GRPO reward functions are provided in Appendix A.1.

Performance evaluation metrics were task-specific. The KEGG Dataset Task utilized Accuracy, Macro F1-score, Macro Precision, and Macro Recall as a multi-class disease prediction assessment, considering potential class imbalances. For VEP Datasets Tasks, Accuracy and F1-score measured the binary pathogenic/benign classification. All LLM and DNA-LLM generations were deterministic with a decoding temperature of 0. We leveraged vLLM for fast inference. [23]

Table 1: Performance comparison of f_{DNA} -only, LLM-only, and DNA-LLM (BIOREASON) models on 290 test datapoints of the KEGG-derived biological reasoning task.

Model	Accuracy	F1-Score	Precision	Recall
[DNA] NT - 500M	86.55	69.76	73.23	66.61
[DNA] Evo2 - 1B	88.28	72.43	75.23	69.83
[LLM] Qwen3 - 1B	85.17	65.71	71.39	64.19
[LLM] Qwen3 - 4B	90.00	79.66	88.24	75.08
[DNA-LLM] NT + Qwen3 - 1B	89.31	81.46	88.24	77.30
[DNA-LLM] NT + Qwen3 - 1B (+GRPO)	91.72	75.06	79.41	72.89
[DNA-LLM] NT + Qwen3 - 4B	95.86	86.25	88.24	84.95
[DNA-LLM] NT + Qwen3 - 4B (+GRPO)	98.28	90.15	91.18	89.62
[DNA-LLM] Evo2 + Qwen3 - 1B	90.42	75.62	77.42	73.91
[DNA-LLM] Evo2 + Qwen3 - 4B	95.17	86.14	91.18	83.33
[DNA-LLM] Evo2 + Qwen3 - 4B (+GRPO)	98.28	93.05	94.12	92.48

Table 2: Performance comparison of f_{DNA} -only, LLM-only, and DNA-LLM (BIOREASON) models on Variant Effect Prediction (VEP) benchmarks (VE-Coding with 1.23K and VE-Non-SNV with 873 test datapoints), evaluating pathogenic/benign classification.

Model	Variant Effect - Coding		Variant Effect - Non-SNV	
	Accuracy	F1-Score	Accuracy	F1-Score
[DNA] NT - 500M	60.91	45.20	67.93	65.97
[DNA] Evo2 - 1B	70.07	49.19	76.17	66.51
[LLM] Qwen3 - 1B	46.55	34.82	70.67	76.21
[LLM] Qwen3 - 4B	48.99	39.58	61.86	67.60
[DNA-LLM] NT + Qwen3 - 1B	55.58	54.50	72.82	76.93
[DNA-LLM] NT + Qwen3 - 4B	60.94	55.66	65.59	73.00
[DNA-LLM] Evo2 + Qwen3 - 1B	72.83	68.90	88.20	89.91
[DNA-LLM] Evo2 + Qwen3 - 4B	80.21	80.00	83.85	85.02

Beyond variant-effect tasks, we also experimented with BIOREASON using the supervised DNA foundation model Enformer [3] for chromatin accessibility prediction, see Appendix B.

5.4 Quantitative Results

BIOREASON’s DNA–LLM hybrids deliver consistent, substantial gains over single-modality baselines on the KEGG-derived reasoning benchmark (Table 1). The Evo2+Qwen3-4B model with GRPO reaches 98.28% accuracy and 93.05% F1, outperforming the standalone Qwen3-4B (90.00%/79.66%) and Evo2 DNA-only (88.28%/72.43%) models. Notably, scaling from Qwen3-1B to Qwen3-4B substantially amplifies both base performance and GRPO effectiveness: while the 1B backbone shows mixed results with GRPO (accuracy improves from 89.31% to 91.72%, but F1 actually declines from 81.46% to 75.06%), the 4B backbone demonstrates dramatic and consistent improvement with GRPO, jumping from 95.17% to 98.28% accuracy and from 86.14% to 93.05% F1. This pattern holds across DNA foundation models; NT+Qwen3-4B with GRPO also reaches 98.28% accuracy with 90.15% F1, confirming that larger LLM backbones provide a significantly more effective substrate for reinforcement learning refinement while simultaneously elevating overall hybrid performance.

On the VEP benchmarks (Table 2), the DNA-LLM hybrids maintain their advantage across variant effect prediction tasks. Evo2+Qwen3-4B achieves 80.21% accuracy and 80.00% F1 in coding variant classification, far surpassing DNA-only (70.07%/49.19%) and LLM-only (48.99%/39.58%) baselines. For non-SNV classification, Evo2+Qwen3-1B leads with 88.20% accuracy and 89.91% F1, surpassing DNA-only (76.17%/66.51%) and LLM-only (70.67%/76.21%). Class-wise breakdowns for the KEGG benchmark are reported in Appendix C.

	Question	Ground Truth (KEGG)	BioReason's Output
KEGG Data Point	<p>Question</p> <pre>< im_start >user Chromosome Number: 17 Network Definition of the Pathway: Actin(monomeric) // PFN1* // Actin(filamentous) Genes in the Pathway: ACTB; actin beta ACTG1; actin gamma 1 PFN1; profilin 1 Given this context, what is the biological effect of this PFN1 allele, specifically what disease does this contribute to? < im_end >< im_start >assistant</pre>	<p>Step 1: The variant KEGG_800 represents a C>G substitution at position 4945969 on chromosome 17, occurring in the PFN1 gene which encodes profilin-1 protein.</p> <p>Step 2: This nucleotide change... potentially affecting a functionally important domain.</p> <p>Step 10: The progressive loss of motor neurons results in muscle weakness, atrophy, and eventually paralysis - the hallmark clinical features of ALS.</p> <p>The C>G mutation in the PFN1 gene... This cytoskeletal dysfunction contributes to axonal transport defects and motor neuron degeneration, ultimately leading to familial Amyotrophic Lateral Sclerosis (ALS).</p>	<pre><think> ... Step 3: This mutation likely alters the protein structure or function of profilin-1, potentially affecting its ability to bind actin monomers or regulate actin polymerization. : Step 10: This specific mutation contributes to ALS through a cytoskeletal/osteostasis pathway... : This C>G mutation in the PFN1 gene likely disrupts profilin-1 function, impairing its ability to regulate actin dynamics by binding monomeric actin and promoting polymerization. The resulting cytoskeletal dysfunction affects motor neuron axonal transport and cellular homeostasis, contributing to the development of familial amyotrophic lateral sclerosis (ALS). </think> Answer: amyotrophic lateral sclerosis (ALS)< im_end ></pre>

Figure 3: Case Study of BIOREASON's Output

5.5 Case Study

To illustrate BIOREASON's reasoning capabilities, consider its analysis of a PFN1 allele on chromosome 17 within the pathway *Actin(monomeric) // PFN1* // Actin(filamentous)*. BIOREASON correctly predicted Amyotrophic Lateral Sclerosis (ALS) as the resultant disease. Significantly, the model generated a plausible 10-step mechanistic rationale, initiating by identifying a specific C>G substitution in the PFN1 gene. Its reasoning then connected this variant to profilin-1 dysfunction, impaired actin dynamics critical for cytoskeletal integrity, subsequent disruption of axonal transport in motor neurons, and finally, the motor neuron degeneration characteristic of ALS. This example highlights BIOREASON's ability to not only make accurate predictions but also to articulate a step-by-step, biologically coherent pathway from a genomic variant to a complex disease phenotype.

6 Discussion

BIOREASON successfully integrates DNA foundation models with large language models, enabling direct LLM reasoning on genomic data. This overcomes key limitations of opaque DNA models and the inability of LLMs to natively process DNA sequences, resulting in enhanced multi-step biological reasoning and superior predictive performance over single-modality approaches.

A core strength of BIOREASON is its interpretable reasoning. By processing contextualized DNA embeddings within the LLM, cultivated through supervised fine-tuning, the system provides not only accurate predictions but also articulates its decision-making via step-by-step mechanistic explanations formatted with '<think>' tokens. This transparency is crucial, allowing researchers to scrutinize the model's logic and translate computational outputs into testable scientific hypotheses.

The broader impact of this work lies in its potential to accelerate biological discovery. BIOREASON offers a robust tool for gaining deeper, mechanistic insights from genomic data, aiding in understanding complex disease pathways and the formulation of novel research questions. The development and application of benchmarks focused on multi-step reasoning, as utilized in this study, will further propel the advancement of AI systems capable of sophisticated biological understanding.

Limitations. Despite its strengths, BIOREASON has several limitations. Reliance on curated datasets such as KEGG introduces potential biases and limits coverage of less-characterized regions. The computational cost of encoding long DNA sequences and applying reinforcement learning (GRPO) raises training and inference time, reducing scalability to whole-genome or real-time settings. DNA sequences were truncated to 2048 tokens due to hardware limits, potentially omitting distal context. Finally, lack of robust uncertainty quantification limits reliability in high-stakes decisions.

Future Work. Future work will focus on expanding BIOREASON's scope and applicability. Key directions include incorporating orthologous sequences to enhance data diversity and model generalizability, and adapting the core framework to other biological modalities such as RNA and protein sequences, thereby addressing a broader range of research questions. Additionally, BIOREASON's improved variant effect prediction capabilities can be leveraged for impactful applications in genome-wide association studies (GWAS) and clinical mutation interpretation.

7 Conclusion

BIOREASON advances computational biology by seamlessly integrating high-capacity DNA sequence encoders with the flexible reasoning of large language models, yielding a unified framework that excels at both mechanistic pathway inference and variant pathogenicity prediction. Across KEGG-derived reasoning tasks and VEP benchmarks, our DNA-LLM hybrids consistently outperform models restricted to a single modality while generating transparent, stepwise explanations that facilitate expert validation. This tight multimodal fusion, further refined through reinforcement learning, not only boosts accuracy but also opens new avenues for interpretable genomic analysis. Future efforts will focus on scaling model size and data, designing leaner architectures, and leveraging modalities such as protein and RNA to broaden BIOREASON 's utility in medicine and biological discovery.

Acknowledgments

We would like to thank Parsa Idehpour for his foundational contributions to the design and implementation of the multimodal GRPO infrastructure, which was essential for large-scale experimentation. We are deeply grateful to Guillaume Fillion for his thoughtful guidance on the selection of biologically meaningful benchmarks. We also thank Ronald Xie for his valuable insights into model architecture. We appreciate Arman Sayed-Ahmadi for stimulating discussions that helped shape the project's future research directions. Finally, we extend our gratitude to Clem Delangue for his support and encouragement in making our models publicly accessible through the Hugging Face platform.

References

- [1] J. Amberger, C. A. Bocchini, A. F. Scott, and A. Hamosh. Mckusick’s online mendelian inheritance in man (omim®). *Nucleic Acids Research*, 37:D793, 2008. ISSN 03051048. doi: 10.1093/NAR/GKN665. URL <https://pmc.ncbi.nlm.nih.gov/articles/PMC2686440/>.
- [2] Anthropic. Claude 3.7 sonnet, February 2025. URL <https://www.anthropic.com/news/claude-3-7-sonnet>. Accessed: 2025-05-13.
- [3] Z. Avsec, V. Agarwal, D. Visentin, J. R. Ledsam, A. Grabska-Barwinska, K. R. Taylor, Y. Assael, J. Jumper, P. Kohli, and D. R. Kelley. Effective gene expression prediction from sequence by integrating long-range interactions. *Nature Methods*, 18(10):1196–1203, Oct. 2021. ISSN 1548-7105. doi: 10.1038/s41592-021-01252-x. URL <https://www.nature.com/articles/s41592-021-01252-x>.
- [4] G. Benegas, C. Ye, C. Albors, J. C. Li, and Y. S. Song. Genomic language models: Opportunities and challenges. *ArXiv*, page arXiv:2407.11435v2, 9 2024. ISSN 2331-8422. URL <https://pmc.ncbi.nlm.nih.gov/articles/PMC11275703/http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC11275703>.
- [5] G. Benegas, C. Albors, A. J. Aw, C. Ye, and Y. S. Song. A dna language model based on multispecies alignment predicts the effects of genome-wide variants. *Nature Biotechnology*, pages 1–6, 1 2025. ISSN 15461696. doi: 10.1038/S41587-024-02511-W;SUBJMETA=114,1305,208,631;KWRD=GENETICS,MACHINE+LEARNING. URL <https://www.nature.com/articles/s41587-024-02511-w>.
- [6] G. Brixì, M. G. Durrant, J. Ku, M. Poli, G. Brockman, D. Chang, G. A. Gonzalez, S. H. King, D. B. Li, A. T. Merchant, M. Naghipourfar, E. Nguyen, C. Ricci-Tam, D. W. Romero, G. Sun, A. Taghibakshi, A. Vorontsov, B. Yang, M. Deng, L. Gorton, N. Nguyen, N. K. Wang, E. Adams, S. A. Baccus, S. Dillmann, S. Ermon, D. Guo, R. Ilango, K. Janik, A. X. Lu, R. Mehta, M. R. Mofrad, M. Y. Ng, J. Pannu, C. Ré, J. C. Schmok, J. S. John, J. Sullivan, K. Zhu, G. Zynda, D. Balsam, P. Collison, A. B. Costa, T. Hernandez-Boussard, E. Ho, M.-Y. Liu, T. McGrath, K. Powell, D. P. Burke, H. Goodarzi, P. D. Hsu, and B. L. Hie. Genome modeling and design across all domains of life with *evo 2*. *bioRxiv*, 2025. doi: 10.1101/2025.02.18.638918. URL <https://www.biorxiv.org/content/early/2025/02/21/2025.02.18.638918>.
- [7] W. Brown. Granular format rewards for eliciting mathematical reasoning capabilities in small language models. <https://gist.github.com/willccbb/4676755236bb08cab5f4e54a0475d6fb>. GitHub Gist.
- [8] S. Chen, L. C. Francioli, J. K. Goodrich, R. L. Collins, M. Kanai, Q. Wang, J. Alféldi, N. A. Watts, C. Vittal, L. D. Gauthier, T. Poterba, M. W. Wilson, Y. Tarasova, W. Phu, R. Grant, M. T. Yohannes, Z. Koenig, Y. Farjoun, E. Banks, S. Donnelly, S. Gabriel, N. Gupta, S. Ferreira, C. Tolonen, S. Novod, L. Bergelson, D. Roazen, V. Ruano-Rubio, M. Covarrubias, C. Llanwarne, N. Petrillo, G. Wade, T. Jeandet, R. Munshi, K. Tibbetts, M. Abreu, C. A. A. Salinas, T. Ahmad, C. M. Albert, D. Ardissino, I. M. Armean, E. G. Atkinson, G. Atzmon, J. Barnard, S. M. Baxter, L. Beaugerie, E. J. Benjamin, D. Benjamin, M. Boehnke, L. L. Bonnycastle, E. P. Bottinger, D. W. Bowden, M. J. Bown, H. Brand, S. Brant, T. Brookings, S. Bryant, S. E. Calvo, H. Campos, J. C. Chambers, J. C. Chan, K. R. Chao, S. Chapman, D. I. Chasman, R. Chisholm, J. Cho, R. Chowdhury, M. K. Chung, W. K. Chung, K. Cibulskis, B. Cohen, K. M. Connolly, A. Correa, B. B. Cummings, D. Dabelea, J. Danesh, D. Darbar, P. Darnowsky, J. Denny, R. Duggirala, J. Dupuis, P. T. Ellinor, R. Elosua, J. Emery, E. England, J. Erdmann, T. Esko, E. Evangelista, D. Fatkin, J. Florez, A. Franke, J. Fu, M. Färkkilä, K. Garimella, J. Gentry, G. Getz, D. C. Glahn, B. Glaser, S. J. Glatt, D. Goldstein, C. Gonzalez, L. Groop, S. Gudmundsson, A. Haessly, C. Haiman, I. Hall, C. L. Hanis, M. Harms, M. Hiltunen, M. M. Holi, C. M. Hultman, C. Jalas, M. Kallela, D. Kaplan, J. Kaprio, S. Kathiresan, E. E. Kenny, B. J. Kim, Y. J. Kim, D. King, G. Kirov, J. Kooner, S. Koskinen, H. M. Krumholz, S. Kugathasan, S. H. Kwak, M. Laakso, N. Lake, T. Langsford, K. M. Laricchia, T. Lehtimäki, M. Lek, E. Lipscomb, R. J. Loos, W. Lu, S. A. Lubitz, T. T. Luna, R. C. Ma, G. M. Marcus, J. Marrugat, K. M. Mattila, S. McCarroll, M. I. McCarthy, J. L. McCauley, D. McGovern, R. McPherson, J. B. Meigs, O. Melander, A. Metspalu, D. Meyers, E. V. Minikel, B. D. Mitchell, V. K. Mootha, A. Naheed, S. Nazarian, P. M. Nilsson, M. C. O’Donovan, Y. Okada, D. Ongur, L. Orozco,

- M. J. Owen, C. Palmer, N. D. Palmer, A. Palotie, K. S. Park, C. Pato, A. E. Pulver, D. Rader, N. Rahman, A. Reiner, A. M. Remes, D. Rhodes, S. Rich, J. D. Rioux, S. Ripatti, D. M. Roden, J. I. Rotter, N. Sahakian, D. Saleheen, V. Salomaa, A. Saltzman, N. J. Samani, K. E. Samocha, A. Sanchis-Juan, J. Scharf, M. Schleicher, H. Schunkert, S. Schönherr, E. G. Seaby, S. H. Shah, M. Shand, T. Sharpe, M. B. Shoemaker, T. Shyong, E. K. Silverman, M. Singer-Berk, P. Sklar, J. T. Smith, J. G. Smith, H. Soyninen, H. Sokol, R. G. Son, J. Soto, T. Spector, C. Stevens, N. O. Stitzel, P. F. Sullivan, J. Suvisaari, E. S. Tai, K. D. Taylor, Y. Y. Teo, M. Tsuang, T. Tuomi, D. Turner, T. Tusie-Luna, E. Vartiainen, M. Vawter, L. Wang, A. Wang, J. S. Ware, H. Watkins, R. K. Weersma, B. Weisburd, M. Wessman, N. Whiffin, J. G. Wilson, R. J. Xavier, A. O'Donnell-Luria, M. Solomonson, C. Seed, A. R. Martin, M. E. Talkowski, H. L. Rehm, M. J. Daly, G. Tiao, B. M. Neale, D. G. MacArthur, and K. J. Karczewski. A genomic mutational constraint map using variation in 76,156 human genomes. *Nature* 2023 625:7993, 625:92–100, 12 2023. ISSN 1476-4687. doi: 10.1038/s41586-023-06045-0. URL <https://www.nature.com/articles/s41586-023-06045-0>.
- [9] M. E. Consens, B. Li, A. R. Poetsch, and S. Gilbert. Genomic language models could transform medicine but not yet. *npj Digital Medicine*, 8:1–4, 12 2025. ISSN 23986352. doi: 10.1038/S41746-025-01603-4;SUBJMETA=1538,692,700;KWRD=HEALTH+CARE,HEALTH+POLICY. URL <https://www.nature.com/articles/s41746-025-01603-4>.
- [10] H. Dalla-Torre, L. Gonzalez, J. Mendoza-Revilla, N. L. Carranza, A. H. Grzywaczewski, F. Oteri, C. Dallago, E. Trop, B. P. de Almeida, H. Sirelkhatim, G. Richard, M. Skwark, K. Beguir, M. Lopez, and T. Pierrot. Nucleotide transformer: building and evaluating robust foundation models for human genomics. *Nature Methods*, 22:287–297, 2 2024. ISSN 15487105. doi: 10.1038/S41592-024-02523-Z;SUBJMETA=114,1305,1647,208,212,631,794;KWRD=GENOMICS,MACHINE+LEARNING,SOFTWARE. URL <https://www.nature.com/articles/s41592-024-02523-z>.
- [11] H. Dalla-Torre, L. Gonzalez, J. Mendoza-Revilla, N. Lopez Carranza, A. H. Grzywaczewski, F. Oteri, C. Dallago, E. Trop, B. P. De Almeida, H. Sirelkhatim, G. Richard, M. Skwark, K. Beguir, M. Lopez, and T. Pierrot. Nucleotide transformer: building and evaluating robust foundation models for human genomics. *Nature Methods*, 22(2):287–297, Feb. 2025. ISSN 1548-7091, 1548-7105. doi: 10.1038/s41592-024-02523-z. URL <https://www.nature.com/articles/s41592-024-02523-z>.
- [12] DeepSeek-AI, D. Guo, D. Yang, H. Zhang, J. Song, R. Zhang, R. Xu, Q. Zhu, S. Ma, P. Wang, X. Bi, X. Zhang, X. Yu, Y. Wu, Z. F. Wu, Z. Gou, Z. Shao, Z. Li, Z. Gao, A. Liu, B. Xue, B. Wang, B. Wu, B. Feng, C. Lu, C. Zhao, C. Deng, C. Zhang, C. Ruan, D. Dai, D. Chen, D. Ji, E. Li, F. Lin, F. Dai, F. Luo, G. Hao, G. Chen, G. Li, H. Zhang, H. Bao, H. Xu, H. Wang, H. Ding, H. Xin, H. Gao, H. Qu, H. Li, J. Guo, J. Li, J. Wang, J. Chen, J. Yuan, J. Qiu, J. Li, J. L. Cai, J. Ni, J. Liang, J. Chen, K. Dong, K. Hu, K. Gao, K. Guan, K. Huang, K. Yu, L. Wang, L. Zhang, L. Zhao, L. Wang, L. Zhang, L. Xu, L. Xia, M. Zhang, M. Zhang, M. Tang, M. Li, M. Wang, M. Li, N. Tian, P. Huang, P. Zhang, Q. Wang, Q. Chen, Q. Du, R. Ge, R. Zhang, R. Pan, R. Wang, R. J. Chen, R. L. Jin, R. Chen, S. Lu, S. Zhou, S. Chen, S. Ye, S. Wang, S. Yu, S. Zhou, S. Pan, S. S. Li, S. Zhou, S. Wu, S. Ye, T. Yun, T. Pei, T. Sun, T. Wang, W. Zeng, W. Zhao, W. Liu, W. Liang, W. Gao, W. Yu, W. Zhang, W. L. Xiao, W. An, X. Liu, X. Wang, X. Chen, X. Nie, X. Cheng, X. Liu, X. Xie, X. Liu, X. Yang, X. Li, X. Su, X. Lin, X. Q. Li, X. Jin, X. Shen, X. Chen, X. Sun, X. Wang, X. Song, X. Zhou, X. Wang, X. Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Y. Zhang, Y. Xu, Y. Li, Y. Zhao, Y. Sun, Y. Wang, Y. Yu, Y. Zhang, Y. Shi, Y. Xiong, Y. He, Y. Piao, Y. Wang, Y. Tan, Y. Ma, Y. Liu, Y. Guo, Y. Ou, Y. Wang, Y. Gong, Y. Zou, Y. He, Y. Xiong, Y. Luo, Y. You, Y. Liu, Y. Zhou, Y. X. Zhu, Y. Xu, Y. Huang, Y. Li, Y. Zheng, Y. Zhu, Y. Ma, Y. Tang, Y. Zha, Y. Yan, Z. Z. Ren, Z. Ren, Z. Sha, Z. Fu, Z. Xu, Z. Xie, Z. Zhang, Z. Hao, Z. Ma, Z. Yan, Z. Wu, Z. Gu, Z. Zhu, Z. Liu, Z. Li, Z. Xie, Z. Song, Z. Pan, Z. Huang, Z. Xu, Z. Zhang, and Z. Zhang. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. 1 2025. URL <https://arxiv.org/pdf/2501.12948>.
- [13] A. Fallahpour, V. Gureghian, G. J. Filion, A. B. Lindner, and A. Pandi. Codontransformer: A multispecies codon optimizer using context-aware neural networks. *Nature Communications*, 16(1), Apr 2025. doi: 10.1038/s41467-025-58588-7.
- [14] A. Fallahpour, J. Ma, A. Munim, H. Lyu, and B. Wang. Medrax: Medical reasoning agent for chest x-ray, 2025. URL <https://arxiv.org/abs/2502.02673>.

- [15] H. Feng, L. Wu, B. Zhao, C. Huff, J. Zhang, J. Wu, L. Lin, P. Wei, C. Wu, P. W. pwei, and A. Professor. Benchmarking dna foundation models for genomic sequence classification running title: Dna foundation models benchmarking. doi: 10.1101/2024.08.16.608288. URL <https://doi.org/10.1101/2024.08.16.608288>.
- [16] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen. Lora: Low-rank adaptation of large language models, 2021. URL <https://arxiv.org/abs/2106.09685>.
- [17] E. Huckvale and H. N. Moseley. kegg pull: a software package for the restful access and pulling from the kyoto encyclopedia of gene and genomes. *BMC Bioinformatics*, 24:1–17, 12 2023. ISSN 14712105. doi: 10.1186/S12859-023-05208-0/TABLES/12. URL <https://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-023-05208-0http://creativecommons.org/publicdomain/zero/1.0/>.
- [18] A.-M. Istrate, F. Milletari, F. Castrotrorres, J. M. Tomczak, M. Torkar, D. Li, and T. Karaletsos. rbio1 - training scientific reasoning llms with biological world models as soft verifiers. *bioRxiv*, 2025. doi: 10.1101/2025.08.18.670981. URL <https://www.biorxiv.org/content/early/2025/08/21/2025.08.18.670981>.
- [19] Q. Jin, Y. Yang, Q. Chen, and Z. Lu. Genegpt: augmenting large language models with domain tools for improved access to biomedical information. *Bioinformatics*, 40, 2 2024. ISSN 13674811. doi: 10.1093/BIOINFORMATICS/BTAE075. URL <https://dx.doi.org/10.1093/bioinformatics/btae075>.
- [20] M. Kanehisa, M. Furumichi, Y. Sato, Y. Matsuura, and M. Ishiguro-Watanabe. Kegg: biological systems database as a model of the real world. *Nucleic Acids Research*, 53:D672–D677, 1 2025. ISSN 0305-1048. doi: 10.1093/NAR/GKAE909. URL <https://dx.doi.org/10.1093/nar/gkae909>.
- [21] J. Kans. Entrez direct: E-utilities on the unix command line - entrez programming utilities help - ncbi bookshelf, 4 2013. URL <https://www.ncbi.nlm.nih.gov/books/NBK179288/>.
- [22] S. H. King, C. L. Driscoll, D. B. Li, D. Guo, A. T. Merchant, G. Brix, M. E. Wilkinson, and B. L. Hie. Generative design of novel bacteriophages with genome language models. *bioRxiv*, 2025. doi: 10.1101/2025.09.12.675911. URL <https://www.biorxiv.org/content/early/2025/09/17/2025.09.12.675911>.
- [23] W. Kwon, Z. Li, S. Zhuang, Y. Sheng, L. Zheng, C. H. Yu, J. E. Gonzalez, H. Zhang, and I. Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.
- [24] M. J. Landrum, J. M. Lee, G. R. Riley, W. Jang, W. S. Rubinstein, D. M. Church, and D. R. Maglott. Clinvar: Public archive of relationships among sequence variation and human phenotype. *Nucleic Acids Research*, 42, 1 2014. ISSN 03051048. doi: 10.1093/NAR/GKT1113,. URL <https://pubmed.ncbi.nlm.nih.gov/24234437/>.
- [25] J. Lee, W. Yoon, S. Kim, D. Kim, S. Kim, C. H. So, and J. Kang. Biobert: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36: 1234–1240, 2 2020. ISSN 1367-4803. doi: 10.1093/BIOINFORMATICS/BTZ682. URL <https://dx.doi.org/10.1093/bioinformatics/btz682>.
- [26] Q. Li, Z. Hu, Y. Wang, L. Li, Y. Fan, I. King, G. Jia, S. Wang, L. Song, and Y. Li. Progress and opportunities of foundation models in bioinformatics. *Briefings in Bioinformatics*, 25:548, 9 2024. ISSN 14774054. doi: 10.1093/BIB/BBAE548. URL <https://dx.doi.org/10.1093/bib/bbae548>.
- [27] Z. Liu, C. Chen, W. Li, P. Qi, T. Pang, C. Du, W. S. Lee, and M. Lin. Understanding r1-zero-like training: A critical perspective. In *Conference on Language Modeling (COLM)*, 2025.
- [28] Z. Ma, J. Zhang, Z. Liu, J. Zhang, J. Tan, M. Shu, J. C. Niebles, S. Heinecke, H. Wang, C. Xiong, R. Krishna, and S. Savarese. Latte: Learning to think with vision specialists, 2025. URL <https://arxiv.org/abs/2412.05479>.

- [29] F. I. Marin, F. Teufel, M. Horlacher, D. Madsen, D. Pultz, O. Winther, and W. Boomsma. Bend: Benchmarking dna language models on biologically meaningful tasks. *12th International Conference on Learning Representations, ICLR 2024*, 11 2023. URL <https://arxiv.org/pdf/2311.12570>.
- [30] S. M. Narayanan, J. D. Braza, R.-R. Griffiths, A. Bou, G. Wellawatte, M. C. Ramos, L. Mitchener, S. G. Rodrigues, and A. D. White. Training a scientific reasoning model for chemistry, 2025. URL <https://arxiv.org/abs/2506.17238>.
- [31] E. Nguyen, M. Poli, M. Faizi, A. W. Thomas, C. B. Sykes, M. Wornow, A. Patel, C. Rabideau, S. Massaroli, Y. Bengio, S. Ermon, S. A. Baccus, and C. Ré. Hyenadna: Long-range genomic sequence modeling at single nucleotide resolution. *ArXiv*, 6 2023. ISSN 2331-8422. URL <https://arxiv.org/pdf/2306.15794>.
- [32] OpenAI, :, A. Hurst, A. Lerer, A. P. Goucher, A. Perelman, A. Ramesh, A. Clark, A. Ostrow, A. Welihinda, A. Hayes, A. Radford, A. Mađry, A. Baker-Whitcomb, A. Beutel, A. Borzunov, A. Carney, A. Chow, A. Kirillov, A. Nichol, A. Paino, A. Renzin, A. T. Passos, A. Kirillov, A. Christakis, A. Conneau, A. Kamali, A. Jabri, A. Moyer, A. Tam, A. Crookes, A. Tootoochian, A. Tootoonchian, A. Kumar, A. Vallone, A. Karpathy, A. Braunstein, A. Cann, A. Codispoti, A. Galu, A. Kondrich, A. Tulloch, A. Mishchenko, A. Baek, A. Jiang, A. Pelisse, A. Woodford, A. Gosalia, A. Dhar, A. Pantuliano, A. Nayak, A. Oliver, B. Zoph, B. Ghorbani, B. Leimberger, B. Rossen, B. Sokolowsky, B. Wang, B. Zweig, B. Hoover, B. Samic, B. McGrew, B. Spero, B. Giertler, B. Cheng, B. Lightcap, B. Walkin, B. Quinn, B. Guarraci, B. Hsu, B. Kellogg, B. Eastman, C. Lugaresi, C. Wainwright, C. Bassin, C. Hudson, C. Chu, C. Nelson, C. Li, C. J. Shern, C. Conger, C. Barette, C. Voss, C. Ding, C. Lu, C. Zhang, C. Beaumont, C. Hal-lacy, C. Koch, C. Gibson, C. Kim, C. Choi, C. McLeavey, C. Hesse, C. Fischer, C. Winter, C. Czarnecki, C. Jarvis, C. Wei, C. Koumouzelis, D. Sherburn, D. Kappler, D. Levin, D. Levy, D. Carr, D. Farhi, D. Mely, D. Robinson, D. Sasaki, D. Jin, D. Valladares, D. Tsipras, D. Li, D. P. Nguyen, D. Findlay, E. Oiwoh, E. Wong, E. Asdar, E. Proehl, E. Yang, E. Antonow, E. Kramer, E. Peterson, E. Sigler, E. Wallace, E. Brevdo, E. Mays, F. Khorasani, F. P. Such, F. Raso, F. Zhang, F. von Lohmann, F. Sulit, G. Goh, G. Oden, G. Salmon, G. Starace, G. Brockman, H. Salman, H. Bao, H. Hu, H. Wong, H. Wang, H. Schmidt, H. Whitney, H. Jun, H. Kirchner, H. P. de Oliveira Pinto, H. Ren, H. Chang, H. W. Chung, I. Kivlichan, I. O’Connell, I. O’Connell, I. Osband, I. Silber, I. Sohl, I. Okuyucu, I. Lan, I. Kostrikov, I. Sutskever, I. Kanitscheider, I. Gulrajani, J. Coxon, J. Menick, J. Pachocki, J. Aung, J. Betker, J. Crooks, J. Lennon, J. Kiros, J. Leike, J. Park, J. Kwon, J. Phang, J. Teplitz, J. Wei, J. Wolfe, J. Chen, J. Harris, J. Varavva, J. G. Lee, J. Shieh, J. Lin, J. Yu, J. Weng, J. Tang, J. Yu, J. Jang, J. Q. Candela, J. Beutler, J. Landers, J. Parish, J. Heidecke, J. Schulman, J. Lachman, J. McKay, J. Uesato, J. Ward, J. W. Kim, J. Huizinga, J. Sitkin, J. Kraaijeveld, J. Gross, J. Kaplan, J. Snyder, J. Achiam, J. Jiao, J. Lee, J. Zhuang, J. Harriman, K. Fricke, K. Hayashi, K. Singhal, K. Shi, K. Karthik, K. Wood, K. Rimbach, K. Hsu, K. Nguyen, K. Gu-Lemberg, K. Button, K. Liu, K. Howe, K. Muthukumar, K. Luther, L. Ahmad, L. Kai, L. Itow, L. Workman, L. Pathak, L. Chen, L. Jing, L. Guy, L. Fedus, L. Zhou, L. Mamitsuka, L. Weng, L. McCallum, L. Held, L. Ouyang, L. Feuvrier, L. Zhang, L. Kondraciuk, L. Kaiser, L. Hewitt, L. Metz, L. Doshi, M. Aflak, M. Simens, M. Boyd, M. Thompson, M. Dukhan, M. Chen, M. Gray, M. Hudnall, M. Zhang, M. Aljube-h, M. Litwin, M. Zeng, M. Johnson, M. Shetty, M. Gupta, M. Shah, M. Yatbaz, M. J. Yang, M. Zhong, M. Glaese, M. Chen, M. Janner, M. Lampe, M. Petrov, M. Wu, M. Wang, M. Fradin, M. Pokrass, M. Castro, M. O. T. de Castro, M. Pavlov, M. Brundage, M. Wang, M. Khan, M. Murati, M. Bavarian, M. Lin, M. Yesildal, N. Soto, N. Gimelshein, N. Cone, N. Staudacher, N. Summers, N. LaFontaine, N. Chowdhury, N. Ryder, N. Stathas, N. Turley, N. Tezak, N. Felix, N. Kudige, N. Keskar, N. Deutsch, N. Bundick, N. Puckett, O. Nachum, O. Okelola, O. Boiko, O. Murk, O. Jaffe, O. Watkins, O. Godement, O. Campbell-Moore, P. Chao, P. McMillan, P. Belov, P. Su, P. Bak, P. Bakkum, P. Deng, P. Dolan, P. Hoeschele, P. Welinder, P. Tillet, P. Pronin, P. Tillet, P. Dhariwal, Q. Yuan, R. Dias, R. Lim, R. Arora, R. Troll, R. Lin, R. G. Lopes, R. Puri, R. Miyara, R. Leike, R. Gaubert, R. Zamani, R. Wang, R. Donnelly, R. Honsby, R. Smith, R. Sahai, R. Ramchandani, R. Huet, R. Carmichael, R. Zellers, R. Chen, R. Chen, R. Nigmatullin, R. Cheu, S. Jain, S. Altman, S. Schoenholz, S. Toizer, S. Miserendino, S. Agarwal, S. Culver, S. Ethersmith, S. Gray, S. Grove, S. Metzger, S. Hermani, S. Jain, S. Zhao, S. Wu, S. Jomoto, S. Wu, Shuaiqi, Xia, S. Phene, S. Papay, S. Narayanan, S. Coffey, S. Lee, S. Hall, S. Balaji, T. Broda, T. Stramer, T. Xu, T. Gogineni, T. Christianson, T. Sanders,

- T. Patwardhan, T. Cunninghamman, T. Degry, T. Dimson, T. Raoux, T. Shadwell, T. Zheng, T. Underwood, T. Markov, T. Sherbakov, T. Rubin, T. Stasi, T. Kaftan, T. Heywood, T. Peterson, T. Walters, T. Eloundou, V. Qi, V. Moeller, V. Monaco, V. Kuo, V. Fomenko, W. Chang, W. Zheng, W. Zhou, W. Manassra, W. Sheu, W. Zaremba, Y. Patil, Y. Qian, Y. Kim, Y. Cheng, Y. Zhang, Y. He, Y. Zhang, Y. Jin, Y. Dai, and Y. Malkov. Gpt-4o system card. 10 2024. URL <https://arxiv.org/pdf/2410.21276>.
- [33] M. Poli, J. Wang, S. Massaroli, J. Quesnelle, R. Carlow, E. Nguyen, and A. Thomas. Striped-Hyena: Moving Beyond Transformers with Hybrid Signal Processing Models, 12 2023. URL <https://github.com/togethercomputer/stripedhyena>.
- [34] Qwen, :, A. Yang, B. Yang, B. Zhang, B. Hui, B. Zheng, B. Yu, C. Li, D. Liu, F. Huang, H. Wei, H. Lin, J. Yang, J. Tu, J. Zhang, J. Yang, J. Yang, J. Zhou, J. Lin, K. Dang, K. Lu, K. Bao, K. Yang, L. Yu, M. Li, M. Xue, P. Zhang, Q. Zhu, R. Men, R. Lin, T. Li, T. Tang, T. Xia, X. Ren, X. Ren, Y. Fan, Y. Su, Y. Zhang, Y. Wan, Y. Liu, Z. Cui, Z. Zhang, and Z. Qiu. Qwen2.5 technical report. 12 2024. URL <https://arxiv.org/pdf/2412.15115>.
- [35] G. Richard, B. P. de Almeida, H. Dalla-Torre, C. Blum, L. Hexemer, P. Pandey, S. Laurent, M. Lopez, A. Laterre, M. Lang, U. Şahin, K. Beguir, and T. Pierrot. Chatnt: A multimodal conversational agent for dna, rna and protein tasks. *bioRxiv*, 2024. doi: 10.1101/2024.04.30.591835. URL <https://www.biorxiv.org/content/early/2024/05/02/2024.04.30.591835>.
- [36] Z. Shao, P. Wang, Q. Zhu, R. Xu, J. Song, X. Bi, H. Zhang, M. Zhang, Y. K. Li, Y. Wu, and D. Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models, 2024. URL <https://arxiv.org/abs/2402.03300>.
- [37] S. T. Sherry, M. H. Ward, M. Kholodov, J. Baker, L. Phan, E. M. Smigielski, and K. Sirotkin. dbsnp: the ncbi database of genetic variation. *Nucleic Acids Research*, 29:308–311, 1 2001. ISSN 0305-1048. doi: 10.1093/NAR/29.1.308. URL <https://dx.doi.org/10.1093/nar/29.1.308>.
- [38] Z. Sondka, N. B. Dhir, D. Carvalho-Silva, S. Jupe, Madhumita, K. McLaren, M. Starkey, S. Ward, J. Wilding, M. Ahmed, J. Argasinska, D. Beare, M. S. Chawla, S. Duke, I. Fasanella, A. G. Neogi, S. Haller, B. Hetenyi, L. Hodges, A. Holmes, R. Lyne, T. Maurel, S. Nair, H. Pedro, A. Sangrador-Vegas, H. Schuilenburg, Z. Sheard, S. Y. Yong, and J. Teague. Cosmic: a curated database of somatic variants and clinical data for cancer. *Nucleic Acids Research*, 52:D1210–D1217, 1 2024. ISSN 0305-1048. doi: 10.1093/NAR/GKAD986. URL <https://dx.doi.org/10.1093/nar/gkad986>.
- [39] J. Su, Y. Lu, S. Pan, A. Murtadha, B. Wen, and Y. Liu. Roformer: Enhanced transformer with rotary position embedding, 2023. URL <https://arxiv.org/abs/2104.09864>.
- [40] E. Wang, S. Schmidgall, P. F. Jaeger, F. Zhang, R. Pilgrim, Y. Matias, J. Barral, D. Fleet, and S. Azizi. Txgemma: Efficient and agentic llms for therapeutics.
- [41] A. Yang, B. Yang, B. Hui, B. Zhang, B. Yu, C. Zhou, C. Li, C. Li, D. Liu, F. Huang, G. Dong, H. Wei, H. Lin, J. Tang, J. Wang, J. Yang, J. Tu, J. Zhang, J. Ma, J. Xu, J. Zhou, J. Bai, J. He, J. Lin, K. Dang, K. Lu, K. Chen, K. Yang, M. Li, M. Xue, N. Ni, P. Zhang, P. Wang, R. Peng, R. Men, R. Gao, R. Lin, S. Wang, S. Bai, S. Tan, T. Zhu, T. Li, T. Liu, W. Ge, X. Deng, X. Zhou, X. Ren, X. Zhang, X. Wei, X. Ren, Y. Fan, Y. Yao, Y. Zhang, Y. Wan, Y. Chu, Y. Liu, Z. Cui, Z. Zhang, and Z. Fan. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2024.
- [42] A. Yang, B. Yang, B. Zhang, B. Hui, B. Zheng, B. Yu, C. Li, D. Liu, F. Huang, H. Wei, H. Lin, J. Yang, J. Tu, J. Zhang, J. Yang, J. Yang, J. Zhou, J. Lin, K. Dang, K. Lu, K. Bao, K. Yang, L. Yu, M. Li, M. Xue, P. Zhang, Q. Zhu, R. Men, R. Lin, T. Li, T. Xia, X. Ren, X. Ren, Y. Fan, Y. Su, Y. Zhang, Y. Wan, Y. Liu, Z. Cui, Z. Zhang, and Z. Qiu. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- [43] Q. Zhang, K. Ding, T. Lyv, X. Wang, Q. Yin, Y. Zhang, J. Yu, Y. Wang, X. Li, Z. Xiang, K. Feng, X. Zhuang, Z. Wang, M. Qin, M. Zhang, J. Zhang, J. Cui, T. Huang, P. Yan, R. Xu, H. Chen, X. Li, X. Fan, H. Xing, and H. Chen. Scientific large language models: A survey on biological and chemical domains. *A Survey on Biological and Chemical Domains*, 1:90, 1 2024. doi: 10.1145/nnnnnnn.nnnnnnn. URL <https://arxiv.org/pdf/2401.14656v2>.

- [44] J. Zhou and O. G. Troyanskaya. Predicting effects of noncoding variants with deep learning-based sequence model. *Nature Methods*, 12(10):931–934, Oct. 2015. ISSN 1548-7091, 1548-7105. doi: 10.1038/nmeth.3547. URL <https://www.nature.com/articles/nmeth.3547>.
- [45] Z. Zhou, Y. Ji, W. Li, P. Dutta, R. V. Davuluri, and H. Liu. Dnabert-2: Efficient foundation model and benchmark for multi-species genome. *12th International Conference on Learning Representations, ICLR 2024*, 6 2023. URL <https://arxiv.org/pdf/2306.15006>.

A Training Details

A.1 Hyper-Parameters

All experiments share the following base settings unless otherwise noted. Please look at our [GitHub](#) for all other details and training scripts.

Optimizer & regularization (SFT).

- Optimizer: AdamW
- Learning rate: 5×10^{-5}
- Weight decay: 1×10^{-2}
- Gradient accumulation: 8 steps
- Random seed: 23
- Devices: 1

LoRA adapters (SFT).

- Rank: 32, Alpha: 64, Dropout: 0.05

Optimizer & regularization (GRPO).

- Optimizer: AdamW
- Learning rate: 1×10^{-5}
- Weight decay: 1×10^{-2}
- Gradient accumulation: 4 steps
- Learning Rate Scheduler: Cosine, 0.03 warmup ratio
- Random seed: 23

LoRA adapters (GRPO).

- Rank: 16, Alpha: 32, Dropout: 0.00

GRPO Parameters.

- Number of generations: 8
- Per device batch size: 8
- Steps: 1000 (7 epochs)
- Devices: 2
- Temperature (4B parameters): 0.7
- Temperature (1.7B parameters): 1
- Top p: 0.95
- Top k: 20
- Beta: 0.0
- Epsilon: 0.2

DeepSpeed & hardware.

- Strategy: deepspeed_stage_2
- CPUs per task: 8
- RAM per node: 128–256 GB
- Data-loader workers: 4

- **Task-specific settings:**
- *KEGG pathway reasoning:*
 - Batch size: 1
 - Epochs: 5
 - Max length DNA: 2048
 - Max text length: 1024 (for LLM only increases to 8192 to fit the raw DNA sequences)
- *Variant effect prediction (coding & non-SNV):*
 - Batch size: 2
 - Epochs: 3
 - Max length DNA: 2048
 - Max text length: 1024 (for LLM only increases to 8192 to fit the raw DNA sequences)

A.2 Computational Resources

We conducted experiments using multiple GPU clusters equipped with NVIDIA A100 and H100 GPUs. A100 systems were equipped with Intel Xeon Silver CPUs, featuring 16-24 CPU cores, 24-32 threads, and 188-251 GB of RAM. We used 4 A100 GPUs for reinforcement learning, while other experiments were performed on single H100 GPUs with Slurm-based orchestration and Deepspeed.

A.3 Reward Details

We use a deterministic composite reward function adapted from [7, 12] to guide reinforcement learning with GRPO, emphasizing correctness and strict adherence to the reasoning format. Each completion is parsed using an XML-aware extractor that isolates the final answer following the last `</think>` tag.

For each output i , the total reward r_i is defined as the sum of the following components:

- **Correctness Reward.** Rewards +2.0 if the extracted final answer matches or contains the ground-truth answer (case-insensitive substring match), and 0.0 otherwise.
- **Conciseness Reward.** Rewards +0.5 if the extracted final answer contains ten or fewer words, encouraging brevity in final responses and preventing loopholes around the correctness reward.
- **Format Reward.** Rewards +0.5 if the completion strictly follows the required reasoning structure of a single `<think>` block followed by a newline and final answer, with properly closed tags.

Each reward component is computed independently and summed per sample, yielding a total reward in $[0, 2.5]$. Rewards are non-differentiable and used solely within GRPO to compute group-normalized advantages. No learned critic or value function is used.

A.4 GRPO Details

Formally, given an input query X_{LLM} , GRPO samples a set of G outputs $\{o_1, \dots, o_G\}$ from the current policy $\pi_{\theta_{\text{old}}}$. Each candidate output o_i comprises a reasoning trace and a final response. Outputs are evaluated using a composite domain-specific reward function $r(q, o_i)$, incorporating the rewards from Appendix A.3.

Dr. GRPO normalizes these rewards into an advantage using the average and standard deviation:

$$A_i = r_i - \text{mean}(\{r_1, \dots, r_G\}) \quad (3)$$

The policy parameters θ are then optimized by maximizing the clipped surrogate objective:

$$\begin{aligned} \mathcal{J}_{\text{GRPO}}(\theta) = & \mathbb{E}[X_{\text{LLM}} \sim P(Q), \{o_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(O|X_{\text{LLM}})] \\ & \frac{1}{G} \sum_{i=1}^G \left(\min \left(\frac{\pi_{\theta}(o_i|X_{\text{LLM}})}{\pi_{\theta_{\text{old}}}(o_i|X_{\text{LLM}})} A_i, \text{clip} \left(\frac{\pi_{\theta}(o_i|X_{\text{LLM}})}{\pi_{\theta_{\text{old}}}(o_i|X_{\text{LLM}})}, 1 - \epsilon, 1 + \epsilon \right) A_i \right) - \beta \mathbb{D}_{\text{KL}}(\pi_{\theta} || \pi_{\text{ref}}) \right) \end{aligned} \quad (4)$$

with hyperparameters ϵ and β .

A.5 GRPO Learning Curve

Figure 4 illustrates the reward progression during GRPO training across three model configurations. Several key patterns emerge from the training dynamics:

Rapid initial learning. All models exhibit steep reward increases in the first 100-200 steps, demonstrating that GRPO efficiently guides the policy toward correct reasoning patterns early in training.

Model size effects. The 4B parameter models reach stable high reward values significantly faster than the 1.7B model, with the larger models stabilizing around step 400 compared to step 800 for the smaller model. The 4B models also demonstrate lower variance in the later training stages, suggesting that increased model capacity provides more robust optimization under GRPO.

Near-optimal performance. All models eventually reach rewards approaching the theoretical maximum of 2.5, indicating successful acquisition of both correct answer generation and adherence to the required reasoning format. The plateau behavior after stabilization suggests stable policy optimization without catastrophic forgetting.

Architecture-agnostic learning. Both NT (Nucleotide Transformer) and Evo2 DNA encoders paired with Qwen3-4B backbones follow nearly identical learning curves, indicating that the GRPO training procedure generalizes effectively across different DNA foundation model architectures.

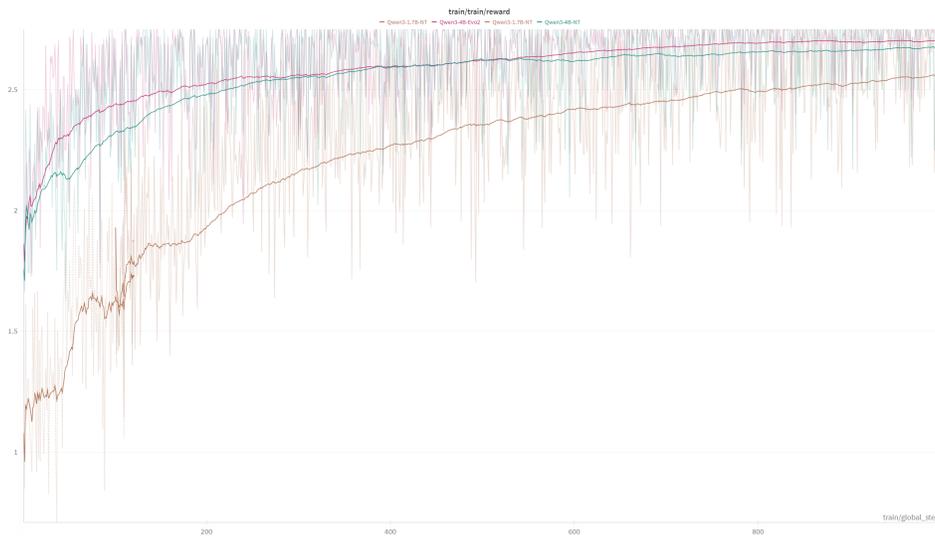


Figure 4: Mean reward progression during GRPO training across different BioReason model configurations over 1000 training steps. Shaded regions represent per-step reward variance across the batch of 8 generations per query.

B Generalization to Supervised DNA Models

To demonstrate that BIOREASON’s architecture generalizes beyond unsupervised DNA foundation models, we conducted an additional experiment integrating a supervised DNA model, Enformer [3], for chromatin accessibility prediction. This experiment addresses two key questions:

- (1) Can BIOREASON leverage supervised DNA encoders trained on specific genomic tasks?
- (2) Does the framework extend to prediction tasks beyond variant-to-disease reasoning?

While our primary experiments utilized unsupervised DNA foundation models that learn general-purpose representations from diverse genomic data, supervised models like Enformer offer complementary strengths. Enformer is explicitly trained to predict genomic annotations including chromatin accessibility, outperforming DeepSEA on DNase-seq variant effect prediction [3]. By integrating Enformer into the BIOREASON framework, we evaluate whether our architecture can extract additional value from task-specific pretrained representations.

B.1 Experimental Setup

We utilized a public dataset derived from the DeepSEA paper [44] containing 60,000 data points for long-range genomics. The human genome was segmented into 200bp bins, with each bin representing a classification instance. The task was formulated as multi-label binary classification to predict chromatin accessibility states ("open" vs "closed" for transcription factor binding) across 20 different DNase-seq tracks encompassing diverse cell types.

We chose to frame this as a classification task because DNase-seq data naturally represents binary chromatin accessibility states. While Enformer was originally trained for regression on epigenomic coverage prediction, our autoregressive architecture is optimized for discrete prediction tasks. Enformer’s demonstrated competency with DNase-seq data, as shown in its original evaluation, suggests it retains relevant knowledge for this formulation.

B.2 Model Configurations

We evaluated the following models:

- **Enformer (DNA-only):** Standalone Enformer with frozen weights and a trained attention head for classification.
- **Qwen3-1.7B and Qwen3-4B (LLM-only):** Qwen3 models fine-tuned on the task with DNA sequences treated as text.
- **Enformer-Qwen3-1.7B and Enformer-Qwen3-4B (DNA-LLM):** BIOREASON hybrids integrating Enformer embeddings with Qwen3 backbones.

All models were trained on a subset of the data and evaluated on a held-out test set following the experimental setup.

B.3 Results

Table 3 presents comprehensive results across multiple evaluation metrics. The BIOREASON hybrids significantly outperform both standalone DNA and LLM baselines. The Enformer-Qwen3-4B model achieves a Macro F1-score of 33.70%, nearly doubling the performance of standalone Enformer (17.18%). This substantial improvement demonstrates that integrating sequence embeddings with an LLM via the BIOREASON architecture provides significant performance gains even when using supervised DNA encoders. The Enformer-Qwen3-1.7B configuration similarly outperforms both component models in isolation.

These results provide two critical insights into BIOREASON’s design. First, the significant performance improvement with Enformer validates that the synergy between DNA sequence embeddings and LLM reasoning is a general architectural principle, not an artifact of using unsupervised models. This confirms the fundamental value of combining specialized sequence encoders with flexible LLM reasoning engines, regardless of the encoder’s training paradigm.

Second, these results illuminate our rationale for prioritizing unsupervised DNA foundation models in the main experiments. While supervised models like Enformer demonstrate strong performance when

Table 3: Performance comparison of f_{DNA} -only, LLM-only, and DNA-LLM (BIOREASON) models on chromatin accessibility prediction across 20 DNase-seq tracks (all metrics in %). “M” and “W” denote *Macro* and *Weighted* averages, respectively.

Model	F1-M	F1-W	Prec-M	Prec-W	Rec-M	Rec-W
[DNA] Enformer	17.18	15.86	15.61	14.85	31.07	29.92
[LLM] Qwen3 - 1.7B	13.01	12.94	20.46	20.11	9.89	9.87
[LLM] Qwen3 - 4B	18.96	18.32	24.47	23.62	16.49	16.06
[DNA-LLM] Enformer + Qwen3 - 1.7B	25.89	24.39	27.02	25.55	33.97	32.45
[DNA-LLM] Enformer + Qwen3 - 4B	33.70	33.08	34.01	33.49	40.02	39.29

integrated into BIOREASON, they are highly specialized for their training objectives, in Enformer’s case, predicting a predefined set of epigenetic marks. The standalone Enformer’s modest performance (17.18% F1) on this task suggests that its specialized embeddings, while powerful, may not be as transferable to novel downstream tasks compared to the rich, general-purpose representations learned by unsupervised models from vast, diverse genomic data.

For BIOREASON’s central goal of enabling broad, multi-step biological reasoning across diverse queries and tasks, unsupervised models provide a more robust foundation. Nevertheless, this experiment demonstrates that BIOREASON’s architecture is flexible enough to accommodate both supervised and unsupervised DNA encoders, allowing users to select the most appropriate foundation model for their specific application domain.

C Per-Disease Performance on the KEGG-Derived Reasoning Benchmark

This section reports class-wise results for the KEGG-derived biological reasoning benchmark, providing a detailed view of BioReason’s performance across individual disease categories. Each class corresponds to a distinct disease entity, with metrics averaged across multiple variant instances. The results demonstrate consistently strong performance across diverse diseases, indicating that BioReason generalizes well beyond high-frequency classes and maintains stable reasoning quality across varied mechanistic contexts.

Table 4: Per-disease performance on the KEGG-derived reasoning benchmark for NT + Qwen3-4B and Evo2 + Qwen3-4B (all metrics in %).

Disease	Freq.	NT + Qwen3-4B				Evo2 + Qwen3-4B			
		Acc.	Prec.	Rec.	F1	Acc.	Prec.	Rec.	F1
Parkinson’s disease	47	99.7	97.9	100.0	98.9	100.0	100.0	100.0	100.0
Alzheimer’s disease	40	99.7	100.0	97.5	98.7	100.0	100.0	100.0	100.0
Spinocerebellar ataxia	36	100.0	100.0	100.0	100.0	99.3	97.2	97.2	97.2
Amyotrophic lateral sclerosis	35	99.7	100.0	97.1	98.6	100.0	100.0	100.0	100.0
Melanoma	17	99.7	100.0	94.1	97.0	99.7	100.0	94.1	97.0
Prion disease	15	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0
Colorectal cancer	12	99.3	85.7	100.0	92.3	99.7	92.3	100.0	96.0
Huntington’s disease	10	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0
Gaucher disease	7	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0
Acute myeloid leukemia	7	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: Our abstract and introduction accurately reflect the contributions and scope of the paper.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: See 6.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: This paper does not include any theoretical results.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Appendix A lists the necessary information in addition to our content pages.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We provide full, anonymized source code and data alongside instructions for reproducing the main experimental results.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Appendix A lists the necessary information in addition to our content pages.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: Comprehensive statistical analyses with multiple runs were limited by the significant computational resources required for these large-scale models. To ensure fair comparison and reproducibility from our single experimental runs, all text generations from BioReason and LLM-only models were performed deterministically by setting the decoding temperature to 0, as detailed in Section 5.3. This approach provides stable point estimates for performance evaluation, with further investigation of inter-run variability deferred to future work.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.

- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: Appendix A lists the necessary information in addition to our content pages.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: The research conducted in this paper conforms to the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: The paper highlights BioReason's substantial positive societal impact, primarily through its potential to accelerate biological discovery and deepen the mechanistic understanding of diseases by generating testable hypotheses. Concurrently, Section 6 (Limitations) addresses crucial considerations for responsible development, such as mitigating potential biases arising from curated training datasets and emphasizing the critical need for robust validation of the model's interpretable outputs to ensure their cautious and effective application in scientific research.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: This research does not release nor produce any data or models that have a high risk for misuse.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: The creators and original owners of assets used in the paper have been appropriately cited.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.

- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. **New assets**

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: We release a codebase and our datasets alongside proper documentation and experiment scripts.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. **Crowdsourcing and research with human subjects**

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: This research did not utilize human subjects or crowdsourcing.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. **Institutional review board (IRB) approvals or equivalent for research with human subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: This research did not utilize human subjects or crowdsourcing.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.

- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: LLMs were not an important, original, or non-standard component of the core methods of this paper.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what should or should not be described.