
LANGPERT: LLM-DRIVEN CONTEXTUAL SYNTHESIS FOR UNSEEN PERTURBATION PREDICTION

Anonymous authors

Paper under double-blind review

ABSTRACT

Predicting cellular responses to previously unseen genetic perturbations remains a fundamental challenge in computational biology, with broad applications in understanding gene function, disease mechanisms, and therapeutic development. Despite advances in computational approaches, developing models that generalise effectively to novel perturbations continues to be difficult. Large Language Models (LLMs) have shown promise in biological applications by synthesizing scientific knowledge, but their direct application to high-dimensional gene expression data has been impractical due to numerical limitations. We propose LangPert, a novel hybrid framework that leverages LLMs to guide a downstream k-nearest neighbors (kNN) aggregator, combining biological reasoning with efficient numerical inference. We demonstrate that LangPert achieves state-of-the-art performance on single-gene perturbation prediction tasks across multiple datasets.

1 INTRODUCTION

Understanding cellular responses to perturbations, in particular gene knockouts, is a cornerstone in deciphering complex biological systems. By systematically altering cellular components via genetic perturbations, researchers can observe cellular behavior changes in genome-wide gene expression vectors, thereby uncovering the genetic mechanisms underlying health and disease. Mapping out even the single-gene-perturbation landscape in a particular cell line requires significant experimental effort. This challenge scales with the number of cell types under investigation, and grows exponentially when considering combinatorial perturbations. This motivates the development of computational approaches that could reduce the need for exhaustive experimental testing by predicting the results of unseen perturbations, particularly important when trying to deconvolute the multicellular functional impact of the thousands of genetic variants associated with complex metabolic disorders like Type 2 Diabetes and Obesity.

The ability to computationally predict the results of unseen genetic perturbations would dramatically accelerate biological discovery while reducing experimental costs. This challenge has recently attracted significant attention, with researchers developing various approaches to leverage prior biological knowledge. These range from transformer-based foundation models, such as scGPT (Cui et al., 2024) and scFoundation (Hao et al., 2024), pre-trained on large-scale cell atlases to methods that explicitly incorporate structured knowledge like gene-gene relationships and ontologies (Roohani et al., 2023). However, despite the sophistication of these approaches, recent studies have revealed a surprising finding: seemingly simple baselines, such as predicting the mean expression response, often outperform more complex deep learning methods (Ahlmann-Eltze et al., 2025; Kernfeld et al., 2024; Wong et al., 2025).

These findings motivate the exploration of alternative approaches that can better leverage biological knowledge while maintaining the ability to handle high-dimensional gene expression data. Particularly promising are methods that can incorporate the vast amount of unstructured biological knowledge present in the scientific literature, which contains detailed and interrelated (if unstructured) information about gene functions, interactions, and regulatory mechanisms that could potentially inform and enhance current perturbation prediction methods.

Large Language Models (LLMs) have recently demonstrated remarkable success in scientific applications, particularly in assisting with data analysis, literature mining, and complex reasoning tasks (Guo et al., 2025; Gao et al., 2024). Their ability to synthesize knowledge from vast scientific corpora and

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

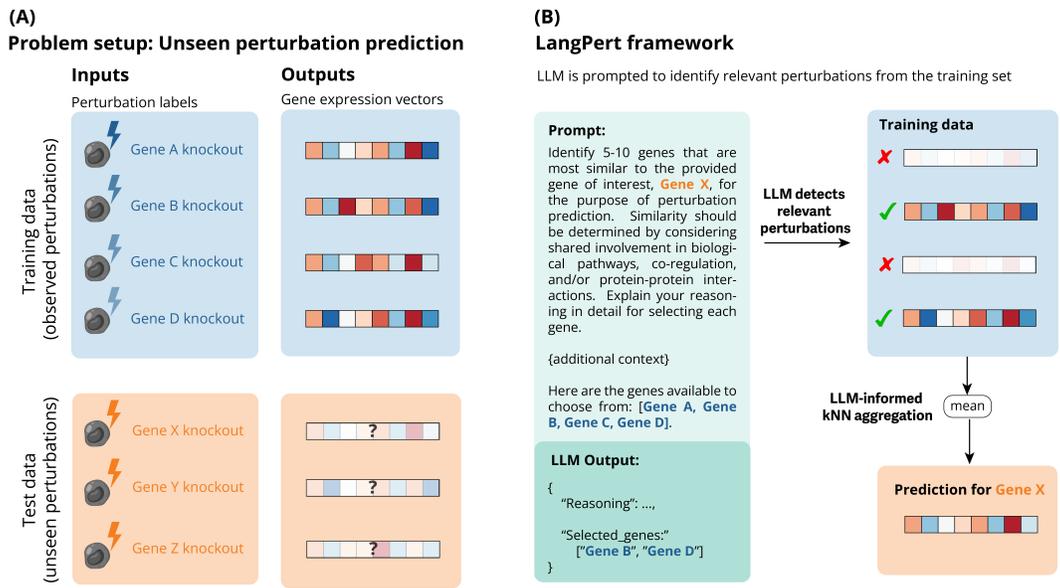


Figure 1: **(A) The task of unseen perturbation outcome prediction illustrated.** The training set consists of pairs $\{(\mathbf{x}_n, \mathbf{y}_n)\}$, where inputs \mathbf{x}_n are discrete perturbation labels and outputs \mathbf{y}_n are high-dimensional gene expression vectors. At test time, the goal is to predict outcomes \mathbf{y}_* corresponding to unseen perturbation labels \mathbf{x}_* . **(B) Our proposed LLM-based LangPert framework.** Instead of using the LLM to directly predict high-dimensional \mathbf{y}_* , the LLM is tasked with finding a relevant small subset from the training perturbation labels $\{(\mathbf{x}_n)\}$ for every \mathbf{x}_* from the held out test set. We use the LLM output to aggregate the corresponding subset of training set expression vectors, effectively resulting in an LLM-informed contextual nearest neighbour prediction.

perform step-by-step reasoning has shown promise for various scientific domains. However, their ability to handle high-dimensional numerical data remains limited due to tokenization constraints and fundamental challenges in numerical computation (Gambardella et al., 2024; Johnson and Hyland-Wood, 2024) with particular difficulties in generating precise continuous values. This limitation is particularly acute in transcriptional response modeling, where predictions are high-dimensional and must capture complex, noisy patterns across thousands of genes.

To leverage the impressive knowledge synthesis capabilities of LLMs while overcoming their limitations with high-dimensional data, we propose a novel framework for predicting cellular outcomes to unseen genetic perturbations. Instead of directly interacting with expression data, we engineer a system where the LLM guides a downstream k-nearest neighbors (kNN) aggregator, combining contextual biological insights with efficient numerical computation, as illustrated in Figure 1. We demonstrate that this hybrid LLM-kNN framework achieves state-of-the-art performance on single-gene perturbation prediction tasks.

2 BACKGROUND

2.1 EXISTING METHODS FOR PERTURBATION DATA

VAE-based methods Variational Autoencoders (VAEs) (Kingma and Welling, 2014; Rezende et al., 2014) have been widely employed as a (conditional) generative model for single-cell perturbation data. Various adaptations, including the Compositional Perturbational Autoencoder (Lotfollahi et al., 2023), sVAE+ (Lopez et al., 2023) and SAMS-VAE (Bereket and Karaletsos, 2023), all capture perturbation effects in the latent space. However, such VAE-based approaches fundamentally lack a mechanism to generalize to *unseen* single-gene perturbations, as they are inherently limited to modeling only the perturbations observed during training. This is because these methods effectively represent perturbations as distinct categorical conditions—while they learn representations of all training-set perturbations, the respective representations for test-set perturbations are absent. The

one-hot or categorical encoding of perturbations does not provide a natural way to infer relationships between observed and unseen conditions, preventing these models from extrapolating beyond the training set.

GEARS Roohani et al. (2023) proposed a graph neural-network based model called GEARS for perturbation response modelling. Unlike VAE-based approaches, GEARS explicitly incorporates structured biological prior knowledge, allowing it to generalize to perturbations involving genes that have not been experimentally tested in the prediction task of interest. Specifically, GEARS incorporates information about gene-gene relationships in two ways, using a gene co-expression graph as well as a gene ontology (GO) knowledge graph.

Single-cell foundation models The success of transformer-based foundation models has spurred their adaptation to single-cell biology. Models such as GeneFormer (Theodoris et al., 2023), scGPT (Cui et al., 2024), and scFoundation (Hao et al., 2024) are pre-trained on large-scale single-cell atlases to learn gene expression patterns. While these models have shown promise in various single-cell analysis tasks, careful evaluations have questioned their fundamental capabilities compared to simpler approaches (Boiarsky et al., 2024; Kedzierska et al., 2023). Particularly in the context of perturbation prediction, where models are fine-tuned to predict responses to unseen genetic interventions, these sophisticated approaches often fail to outperform simple baselines such as mean prediction (Ahlmann-Eltze et al., 2025; Kernfeld et al., 2024).

LLM-informed gene embeddings LLMs have been widely applied across scientific domains, including biology (Lee et al., 2020). In the context of gene-level biological knowledge, recent methods such as GenePT (Chen and Zou, 2024) have taken a novel approach: instead of training foundation models on gene expression data, they leverage LLMs’ understanding of scientific literature to generate gene embeddings. These embeddings, derived from NCBI text descriptions of genes, have shown promising results in observational single-cell analysis tasks. Building on this idea, Märtens et al. (2024) extended the approach to interventional settings, developing a GP+LLM model that combines a Gaussian Process with literature-derived embeddings as well as protein language model embeddings to predict perturbation outcomes, demonstrating that LLMs can effectively encode biologically relevant prior knowledge.

2.2 LEVERAGING LLMs FOR BIOLOGICAL KNOWLEDGE

Large Language Models have emerged as powerful tools for synthesizing biological knowledge from scientific literature, offering new approaches to understanding gene functions and relationships. While methods like GP+LLM (Märtens et al., 2024) have shown promise by leveraging LLM-derived embeddings in predictive models, these embeddings are inherently *static*, as they are extracted from a fixed body of literature (e.g., NCBI abstracts in (Chen and Zou, 2024)) as condensed into an LLM at a specific training data corpus cutoff-date. This limits their adaptability when reasoning about unseen perturbations, where context-dependent interactions may play a crucial role.

Recently, Wu et al. (2024) introduced PerturbQA, a benchmark that represents perturbation data as “(perturbation, gene, outcome)” triplets where outcome is a binary variable. They consider two types of outcomes: whether a particular gene is differentially expressed, and for differentially expressed genes, whether they are up- or down-regulated. This classification setup makes the data more amenable for use with LLMs. In contrast, our work aims to predict high-dimensional gene expression vectors directly, an approach that presents different technical challenges for LLM implementation.

This direct application of LLMs to high-dimensional perturbation modeling is constrained by tokenization limits and numerical precision issues. These challenges motivate our hybrid approach, which integrates LLM-driven biological reasoning with computational models designed for handling high-dimensional data – a paradigm we explore in LangPert.

3 LANGPERT: A HYBRID LLM-KNN FRAMEWORK

We propose LangPert, a framework that leverages LLMs’ ability to reason about biological mechanisms while ameliorating their limitations in handling high-dimensional data. Instead of using

LLMs to generate static embeddings or to directly predict expression values, LangPert employs an LLM to identify biologically relevant training examples that can inform predictions for unseen perturbations. These LLM-selected examples then guide a downstream aggregation function (here we employ a k-nearest neighbors (kNN) scheme) that performs the actual numerical computations in the high-dimensional expression space.

Problem formulation As illustrated in Figure 1, the task of unseen perturbation response prediction presents a challenging supervised learning challenge, where given pairs $\{(\mathbf{x}_n, \mathbf{y}_n)\}$, inputs \mathbf{x}_n are discrete perturbation labels and outputs $\mathbf{y}_n \in \mathbb{R}^D$ are high-dimensional numeric readouts corresponding to those labels. The goal is to predict responses \mathbf{y}_* for test inputs \mathbf{x}_* which are distinct from those in the training set. This makes the task fundamentally different from standard supervised learning, as it requires extrapolation to entirely new perturbations, which would fall outside of the one-hot representational space of the training data perturbation categories.

Naive application of LLMs One approach to utilise LLMs for this problem is via in-context learning (ICL), where $(\mathbf{x}_n, \mathbf{y}_n)$ pairs are given to the LLM as part of the input prompt, alongside with new inputs \mathbf{x}_* . However, this approach is problematic due to the high dimensionality of gene expression vectors. Therefore, we propose a strategy to remedy this challenge inherent in a naive / brute force application of LLMs in this domain.

LangPert adaptation for high-dimensional outcomes Specifically, we propose to only show the model training inputs $\{\mathbf{x}_n\}$ alongside with a test input \mathbf{x}_* . As these inputs correspond to perturbation labels – for genetic perturbations these would be gene names – LLMs have demonstrated strong capabilities in reasoning about biological relationships and identifying functionally related genes through their training on scientific literature (Hu et al., 2025). We leverage these established capabilities, using the LLM’s comprehensive knowledge of biological systems and its proven few-shot learning abilities (Brown et al., 2020) to identify genes from the training set that are functionally relevant for the prediction target \mathbf{x}_* . This biological reasoning can be further enhanced through relevant context included in the prompt by a human user.

Contextual aggregation Given the LLM output of a subset of gene perturbation labels relevant to the input prompt, the final output of an unseen perturbation prediction pipeline can be made via an aggregation / reduction of the gene expression vectors of those relevant genes. That is, given a relevant subset \mathcal{G} where $\{(\mathbf{x}_n, \mathbf{y}_n)\}$ for $n \in \mathcal{G}$, we propose to make the prediction $\frac{1}{|\mathcal{G}|} \sum_{n \in \mathcal{G}} \mathbf{y}_n$ for the unseen output. This can be interpreted as a k nearest neighbour predictor, where the relevant neighbours are identified by the LLM and then averaged. In principle, different aggregation techniques can be adapted here, for example weighting the inputs, using nonparametric (e.g. median) reductions, or even passing the subset data $\{(\mathbf{x}_n, \mathbf{y}_n)\}$ to a small tabular prediction framework. In this work, we found it sufficient to use simple averaging to achieve SOTA results, but future work investigating more principled aggregation/reduction strategies may yield even stronger predictions.

In summary, we have proposed a hybrid LLM-kNN framework, where predictions take the following form

$$\mathbf{y}_* = \frac{1}{\sum_n w_n} \sum_n w_n \mathbf{y}_n, \text{ where } w_n = \text{LLM}(\mathbf{x}_*, \{\mathbf{x}_n\}, \text{context}) \in \{0, 1\}$$

where the LLM sees the test perturbation label \mathbf{x}_* , all training labels $\{\mathbf{x}_n\}$ and potentially additional information presented in the prompt. The number of chosen training perturbations, i.e. the number of nearest neighbours $k := \sum_n w_n$ can either be specified in the prompt or remain unspecified, giving the LLM flexibility to choose.

4 RESULTS

4.1 EXPERIMENTAL SETUP

Datasets For evaluation, we consider data from large-scale Perturb-seq screens across two cell lines: the leukemia cell line (K562) and the retinal pigment epithelial (RPE1) cell line from (Replogle

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

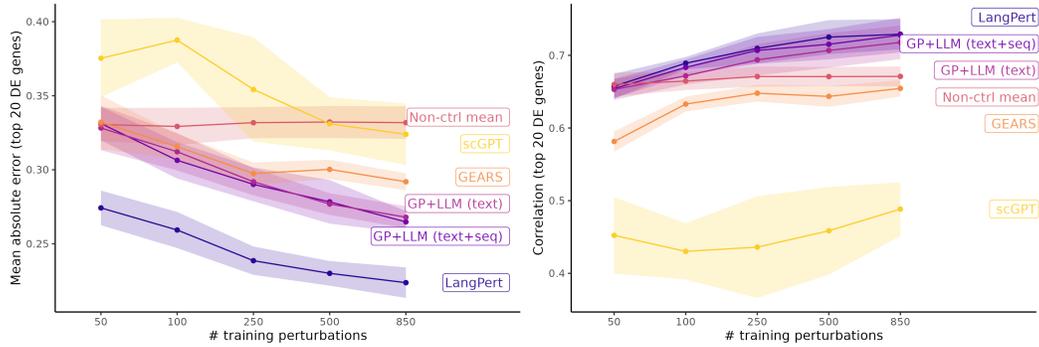


Figure 2: Performance comparison across models (scGPT, GEARS, non-control mean baseline, GP+LLM, and LangPert) evaluated using mean absolute error (MAE, lower is better) and Pearson correlation between predicted and observed differences from control cells (higher is better). Shaded regions indicate ± 1.96 standard errors across data splits. Model performance is shown as a function of training set size, varying from 50 to 850 perturbations.

et al., 2022). We use the version of the data curated by Roohani et al. (2023)¹, with a total of 1092 perturbations in the K562 cell line and a total of 1543 perturbations in the RPE1 cell line.

Experimental details and metrics In all our evaluations, we assess performance in a 5-fold cross-validation setting, so in the end, all metrics are calculated on the entire set of 1092 perturbations in K562 cell line, and 1543 perturbations in RPE1. Predictions on every held-out cross-validation fold are made independently to avoid any data leakage. In experiments where we consider a gradually increasing number of training perturbations (e.g. along x -axis in Figure 2), for every cross-validation split we repeatedly downsample the training set.

Following practices from literature, we quantify perturbation prediction performance relative to control cells, i.e. using Pearson correlation, mean absolute error (MAE), and mean squared error (MSE) on the differences $\Delta_n := \mathbf{y}_n - \mathbf{y}_{\text{control}}$. Following Roohani et al. (2023), we calculate both metrics across the top 20 differentially expressed genes relative to control cells, resulting in gene sets that are specific for every perturbation.

Baselines As discussed in the Background section 2.1, we consider existing methods which have a capability to generalise to unseen single-gene perturbations. Specifically, we consider the graph neural network approach GEARS, fine-tuning a single-cell foundation model scGPT, and two versions of the GP+LLM model (one using NCBI text embeddings as input, the other combining text embeddings and protein sequence embeddings). We also consider a non-control mean baseline that has been shown to be surprisingly effective (Kernfeld et al., 2024; Märtens et al., 2024).

Choice of LLM A crucial component of LangPert is its LLM engine, meaning the choice of LLM can significantly impact its behavior and performance. For all comparisons in Section 4.2, we use Claude 3.5 Sonnet. Later, in Section 4.4, we examine how performance varies across different LLMs, ranging from small open-weight LLMs with 8 billion parameters to frontier models.

4.2 PERFORMANCE COMPARISON

Results on K562 cell line Figure 2 shows the performance metrics (MAE and correlation) on the K562 cell line dataset across a varying number of training perturbations. The ordering of existing baselines is aligned with what has been reported in literature: the fine-tuned scGPT is the lowest performing model, followed by GEARS. GEARS outperforms the non-control mean baseline according to the MAE metric, but slightly underperforms in correlation. The GP+LLM models outperform both scGPT and GEARS.

¹Available in <https://github.com/snap-stanford/GEARS>

Table 1: **Results on K562 cell line dataset:** Performance comparison of different models at different training data sizes. Values shown as mean \pm standard error.

Model	100 training perturbations			850 training perturbations		
	MAE \downarrow	MSE \downarrow	Correlation \uparrow	MAE \downarrow	MSE \downarrow	Correlation \uparrow
scGPT	0.388 \pm 0.016	0.234 \pm 0.014	0.430 \pm 0.021	0.324 \pm 0.008	0.195 \pm 0.012	0.488 \pm 0.016
GEARS	0.316 \pm 0.004	0.164 \pm 0.004	0.633 \pm 0.005	0.292 \pm 0.002	0.147 \pm 0.003	0.655 \pm 0.007
Non-ctrl mean	0.329 \pm 0.006	0.175 \pm 0.008	0.665 \pm 0.006	0.332 \pm 0.006	0.176 \pm 0.007	0.671 \pm 0.007
GP+LLM (text)	0.312 \pm 0.006	0.162 \pm 0.008	0.672 \pm 0.006	0.268 \pm 0.004	0.132 \pm 0.006	0.718 \pm 0.012
GP+LLM (text+seq)	0.306 \pm 0.006	0.158 \pm 0.008	0.683 \pm 0.006	0.265 \pm 0.004	0.130 \pm 0.005	0.728 \pm 0.012
LangPert	0.259 \pm 0.006	0.132 \pm 0.005	0.689 \pm 0.004	0.224 \pm 0.005	0.108 \pm 0.005	0.731 \pm 0.011

Table 2: **Results on RPE1 cell line dataset:** Performance comparison of different models at different training data sizes. Values shown as mean \pm standard error.

Model	50 training perturbations			100 training perturbations		
	MAE \downarrow	MSE \downarrow	Correlation \uparrow	MAE \downarrow	MSE \downarrow	Correlation \uparrow
scGPT	0.452 \pm 0.008	0.346 \pm 0.016	0.627 \pm 0.009	0.451 \pm 0.009	0.354 \pm 0.006	0.642 \pm 0.003
GEARS	0.471 \pm 0.013	0.343 \pm 0.019	0.670 \pm 0.007	0.435 \pm 0.022	0.307 \pm 0.026	0.706 \pm 0.005
Non-ctrl mean	0.427 \pm 0.010	0.303 \pm 0.015	0.737 \pm 0.006	0.430 \pm 0.006	0.305 \pm 0.013	0.738 \pm 0.005
GP+LLM (text)	0.414 \pm 0.011	0.289 \pm 0.014	0.721 \pm 0.002	0.408 \pm 0.005	0.282 \pm 0.009	0.725 \pm 0.002
GP+LLM (text+seq)	0.417 \pm 0.013	0.294 \pm 0.016	0.717 \pm 0.004	0.403 \pm 0.005	0.277 \pm 0.009	0.723 \pm 0.004
LangPert	0.368 \pm 0.008	0.249 \pm 0.012	0.726 \pm 0.005	0.361 \pm 0.008	0.239 \pm 0.013	0.737 \pm 0.008
Model	250 training perturbations			1170 training perturbations		
	MAE \downarrow	MSE \downarrow	Correlation \uparrow	MAE \downarrow	MSE \downarrow	Correlation \uparrow
scGPT	0.436 \pm 0.024	0.329 \pm 0.030	0.641 \pm 0.009	0.449 \pm 0.038	0.349 \pm 0.040	0.603 \pm 0.040
GEARS	0.430 \pm 0.011	0.291 \pm 0.015	0.720 \pm 0.007	0.405 \pm 0.015	0.266 \pm 0.015	0.710 \pm 0.011
Non-ctrl mean	0.432 \pm 0.005	0.306 \pm 0.011	0.741 \pm 0.004	0.434 \pm 0.005	0.308 \pm 0.011	0.743 \pm 0.005
GP+LLM (text)	0.400 \pm 0.004	0.271 \pm 0.010	0.735 \pm 0.003	0.371 \pm 0.004	0.238 \pm 0.009	0.759 \pm 0.005
GP+LLM (text+seq)	0.397 \pm 0.003	0.269 \pm 0.009	0.732 \pm 0.003	0.364 \pm 0.004	0.233 \pm 0.009	0.760 \pm 0.005
LangPert	0.344 \pm 0.005	0.218 \pm 0.010	0.753 \pm 0.009	0.318 \pm 0.004	0.192 \pm 0.006	0.772 \pm 0.005

Our proposed LangPert significantly outperforms all existing models according to the MAE and MSE metrics (see Table 1 for numerical values), and also achieves a slightly higher correlation value. For example, in the scenario with 850 training perturbations, LangPert achieves MAE of 0.224(\pm 0.005) which is a substantial improvement over the second best GP+LLM’s 0.265(\pm 0.004).

Results on RPE1 cell line We conducted a similar experiment on the RPE1 cell line, with results summarized in Table 2. The ranking of methods remains consistent with previous benchmarks, with LangPert achieving the best results in MAE and MSE metrics. For the correlation metric, at the smallest sample size (50 training perturbations), the non-control mean achieves the highest correlation (0.737 \pm 0.006), slightly surpassing LangPert (0.726 \pm 0.005). However, as the sample size increases, LangPert outperforms all baselines. At the largest sample size (1170 perturbations), LangPert achieves an MAE of 0.318 \pm 0.004 (compared to the second-best LLM+GP at 0.364 \pm 0.004) and a correlation of 0.772 \pm 0.005 (vs. LLM+GP’s 0.760 \pm 0.005).

Overall, LangPert sets a new state-of-the-art performance on both the K562 and RPE1 benchmarks.

4.3 LANGPERT’S CONTEXT-DEPENDENT REASONING

To understand the source of LangPert’s performance gains, we analysed the model’s reasoning traces to examine how it selects relevant genes for each perturbation. A key challenge in kNN-based approaches is determining which similarity criterion to use: genes can be related through shared pathways, protein complexes, subcellular localization, or functional processes. Our analysis revealed that rather than applying a single similarity metric uniformly, LangPert adaptively identifies the most relevant organising principle for each gene based on its specific biological context. The model selects different organising principles based on each gene’s functional role:

- **MTOR** \rightarrow **Signaling pathway hierarchy:** selects RPTOR (direct complex partner), EIF4E/EIF4G1 (downstream effectors), PPP2R1A (feedback regulator)
- **EIF3E** \rightarrow **Protein complex membership:** selects other eIF3 subunits (EIF3A, EIF3B, EIF3I, EIF3F) within the translation initiation complex

- **PSMD11** → **Functional machinery**: selects other proteasome subunits (PSMD1, PSMC4, PSMC6, PSMD3, PSMA3) required for protein degradation

This context-dependent reasoning allows LangPert to capture synergistic interactions within biological systems by adaptively choosing the right biological lens for each gene—pathway hierarchies for regulators, protein complexes for structural components, or functional machinery for metabolic genes. This contrasts with static database methods like GEARS, which apply Gene Ontology uniformly across genes, rather than selecting the most relevant similarity framework for each perturbation.

4.4 COMPARATIVE PERFORMANCE OF LLM VARIANTS

The LLM component in the LangPert framework has two main elements: the choice of the LLM itself and the prompting strategy. In this section, we evaluate how different LLM variants in the LangPert framework affect its performance.

We compare a range of models, from small open LLMs such as Llama 3.1 8b (Grattafiori et al., 2024) and TxGemma-9b (a Gemma variant pre-trained by its creators for therapeutic applications) (Wang et al., 2025), to a medium-sized Llama 3.3 70b model and the larger DeepSeek R1 model (DeepSeek-AI et al., 2025). We also evaluate frontier models Claude Sonnet 3.5 and OpenAI o1.

The LLM comparison is illustrated in Figure 3, where the blue shaded area represents the performance coverage across the LLM variants. Here we show two metrics, correlation and MAE, with the non-control mean baseline dividing the plot into four quadrants. Methods that outperform the mean in both metrics appear in the green area. Existing non-LLM baselines are shown in red (scGPT, GEARS, GP+LLM). While all methods outperform the mean baseline in MAE, not all of them do in correlation.

Among LLM variants, the smaller TxGemma-9b and Llama 8b models show the lowest performance (achieving MAE values of 0.315 and 0.284, respectively), yet still outperform scGPT. Next in the performance hierarchy, LangPert using either DeepSeek-R1 or Llama 3.3 70b models (with MAE values of 0.254 and 0.247, respectively) outperforms GEARS. Finally, the frontier models Claude

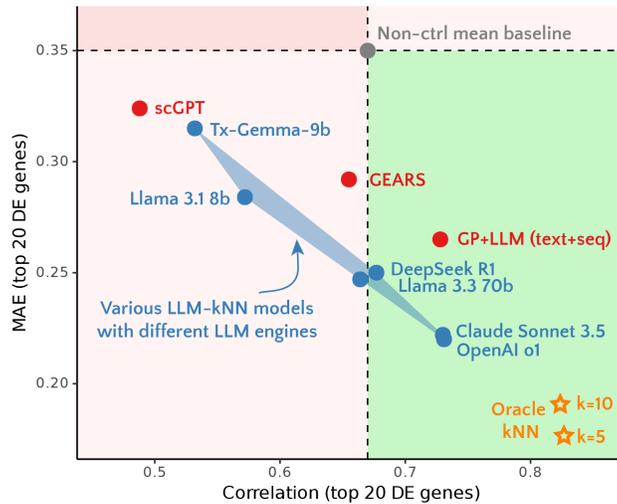


Figure 3: Performance comparison of various versions of LangPert across different LLM backends (blue shaded area) to existing methods (in red), according to two metrics, correlation (x -axis) and MAE (y -axis) on the full K562 cell line dataset. The non-control mean baseline divides the area into four quadrants: Ideally, strong performance lies in the green shaded area (i.e. a smaller MAE and higher correlation). Among all LLMs, the small open models (Llama 8b and TxGemma-9b) are the least performant, and frontier models (Claude Sonnet 3.5 and OpenAI o1) are best performing. For reference, we have added an Oracle kNN method (in orange, for $k = 5$ and $k = 10$) to indicate an approximate upper bound on the performance of a kNN-based approach which has access to the ground truth.

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

Claude Sonnet 3.5 output for query gene SMG5:

Reasoning: SMG5 is a critical component of the nonsense-mediated mRNA decay (NMD) pathway, which degrades mRNAs containing premature termination codons. In analyzing potential functional similarities, I focused on several key aspects:

1. Core NMD machinery components and regulators that would disrupt similar RNA surveillance mechanisms
2. Genes involved in mRNA quality control and processing that functionally overlap with NMD
3. Translation termination factors and related proteins, as SMG5 interacts with release factors
4. Proteins involved in mRNA degradation pathways

In K562 cells specifically: (1) The NMD pathway is particularly important due to high transcriptional activity in these cancer cells, (2) Many aberrant transcripts need surveillance due to genomic instability, (3) Translation regulation is critical for K562 survival. Key predicted similar phenotypes would come from disrupting: (1) Core NMD factors, (2) General mRNA decay machinery, (3) Translation termination components, (4) Related quality control pathways.

Inferred kNN gene set: [UPF2, UPF1, DIS3, EXOSC10, RBM8A, RNPS1]

OpenAI o1 output for query gene SMG5:

Reasoning: SMG5 encodes a core component of the nonsense-mediated mRNA decay (NMD) pathway, where it helps recruit phosphatase activity to UPF1, leading to the degradation of aberrant transcripts. In K562 cells (a leukemic line with high transcription/translation activity), knocking out other critical NMD factors or key Exon Junction Complex (EJC) components should yield phenotypes most similar to SMG5 loss, because these are all essential for detecting and degrading faulty mRNAs. First among these is UPF1, the central NMD helicase that physically interacts with SMG5. UPF2 is another core NMD scaffold protein downstream of UPF1. RBM8A (Y14) and MAGOH are core EJC members required for coupling exon junction marking to transcript decay. RNPS1 also associates with EJC/NMD surveillance. Finally, DDX6, while more broadly involved in mRNA decapping and turnover, is still closely tied to mRNA quality-control outcomes and thus expected to share certain knockout effects with SMG5 in K562.

Inferred kNN gene set: [UPF1, UPF2, RBM8A, MAGOH, RNPS1, DDX6]

Figure 4: Comparison of LLM-based inference for identifying genes with similar perturbation outcomes to SMG5 in K562 cells. Both Claude and OpenAI o1 identified key nonsense-mediated mRNA decay (NMD) pathway components and associated factors.

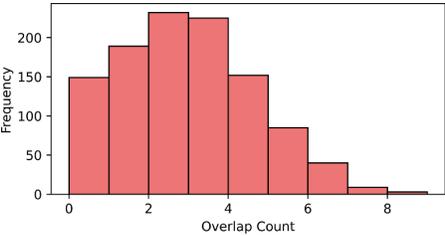


Figure 5: Overlap in inferred gene sets between the two frontier models: Claude Sonnet 3.5 and OpenAI o1.

432 Sonnet 3.5 and OpenAI o1 outperform all other approaches. Interestingly, both frontier models
433 achieve near-identical performance on this task.

434 A natural question is whether future improvements in LLM capabilities could lead to further per-
435 formance gains in our framework. To establish an approximate upper bound on what LLM-kNN
436 approaches can achieve, we implemented an Oracle kNN model that has access to the ground truth
437 (unseen) perturbation outcome. We added the performance of this Oracle model with $k = 5$ and
438 $k = 10$ to Figure 3 as a reference point (“Oracle kNN” in orange).

439 Given the near-identical performance of the two frontier models, we investigated whether their
440 inferred kNN gene sets were also identical. Figure 5 shows the overlap distribution between the two
441 models, revealing considerable variation in gene set selection despite resulting in similarly predictive
442 outcomes. This suggests that different models identify distinct but equally relevant gene relationships,
443 potentially corresponding to alternative biological hypotheses with similar predictive power.

444 In Figure 4, we illustrate this phenomenon with outputs for a selected query perturbation, SMG5,
445 comparing the two frontier models. Both Claude Sonnet 3.5 and OpenAI o1 successfully identified
446 key nonsense-mediated mRNA decay (NMD) pathway components and associated factors. Consistent
447 with our overlap analysis, they inferred partially overlapping but not identical gene sets. These
448 observations suggest promising directions for future work, including exploring ensembling strategies
449 or methods to integrate insights from multiple LLMs to enhance both interpretability and robustness.

451 5 DISCUSSION

452 We introduce LangPert, a novel hybrid framework that leverages Large Language Models to guide
453 k-nearest neighbour predictions for unseen genetic perturbations. LangPert achieves state-of-the-art
454 performance across multiple datasets and metrics, particularly as measured by mean absolute error
455 of predictions across the entire spectrum of low-to-high data regimes. Our analysis reveals that
456 this performance stems from LangPert’s ability to adaptively select different biological similarity
457 frameworks based on each gene’s functional context.

458 LangPert’s key innovation lies in its ability to harness LLMs’ biological reasoning capabilities
459 without being constrained by their numerical limitations. By using LLMs to identify relevant training
460 examples and employing kNN for aggregation, our approach effectively bridges the gap between
461 knowledge-driven and data-driven methodologies. This hybrid strategy outperforms existing methods
462 such as GEARS, scGPT, and GP+LLM.

463 The superior performance of LangPert has significant implications for perturbation biology. Improved
464 predictive models could dramatically reduce the need for exhaustive experimental testing, accelerating
465 biological discovery and potentially informing more efficient experimental design strategies. Fur-
466 thermore, LangPert’s flexible framework allows for the incorporation of different LLMs, prompting
467 strategies, and aggregation methods, suggesting ample room for further optimization and adaptation
468 to various biological contexts.

469 Despite its promising results, LangPert is not without limitations. Its performance depends on the
470 quality and up-to-date nature of the LLM’s knowledge, and potential biases in LLM training data
471 could influence predictions. Future work should explore more sophisticated aggregation methods
472 beyond averaging, incorporate uncertainty quantification, and extend the approach to multi-gene
473 perturbations or other types of biological interventions. Additionally, investigating the impact of
474 different LLMs as well as prompting strategies could further enhance the model’s capabilities.

475 The success of LangPert in integrating LLMs with traditional machine learning techniques for
476 high-dimensional biological data suggests potential applications beyond perturbation biology. This
477 approach could be adapted to other scientific domains characterized by high-dimensional outcomes
478 and rich contextual knowledge, such as multi-objective molecular property prediction.

479 In conclusion, LangPert represents a significant advance in our ability to predict cellular responses to
480 unseen genetic perturbations. By effectively combining the strengths of LLMs and traditional machine
481 learning approaches, it opens new avenues for accelerating biological discovery and deepening our
482 understanding of complex cellular systems. As we continue to refine and expand this approach,
483 we anticipate its impact to grow, potentially transforming how we approach predictive modeling in
484 biology and beyond.

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

REFERENCES

- C. Ahlmann-Eltze, W. Huber, and S. Anders. Deep learning-based predictions of gene perturbation effects do not yet outperform simple linear baselines, Feb. 2025. URL <https://www.biorxiv.org/content/10.1101/2024.09.16.613342v5>. Pages: 2024.09.16.613342 Section: New Results.
- M. Bereket and T. Karaletsos. Modelling Cellular Perturbations with the Sparse Additive Mechanism Shift Variational Autoencoder. Nov. 2023. URL <https://openreview.net/forum?id=DzaCE00jGV>.
- R. Boiarsky, N. M. Singh, A. Buendia, A. P. Amini, G. Getz, and D. Sontag. Deeper evaluation of a single-cell foundation model. *Nature Machine Intelligence*, 6(12):1443–1446, Dec. 2024. ISSN 2522-5839. doi: 10.1038/s42256-024-00949-w. URL <https://www.nature.com/articles/s42256-024-00949-w>. Publisher: Nature Publishing Group.
- T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei. Language Models are Few-Shot Learners, July 2020. URL <http://arxiv.org/abs/2005.14165>. arXiv:2005.14165 [cs].
- Y. Chen and J. Zou. Simple and effective embedding model for single-cell biology built from chatgpt. *Nature Biomedical Engineering*, pages 1–11, 2024.
- H. Cui, C. Wang, H. Maan, K. Pang, F. Luo, N. Duan, and B. Wang. scgpt: toward building a foundation model for single-cell multi-omics using generative ai. *Nature Methods*, pages 1–11, 2024.
- DeepSeek-AI, D. Guo, D. Yang, H. Zhang, J. Song, A. Liu, 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. 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, Jan. 2025. URL <http://arxiv.org/abs/2501.12948>. arXiv:2501.12948 [cs].
- A. Gambardella, Y. Iwasawa, and Y. Matsuo. Language Models Do Hard Arithmetic Tasks Easily and Hardly Do Easy Arithmetic Tasks, June 2024. URL <http://arxiv.org/abs/2406.02356>. arXiv:2406.02356 [cs].
- S. Gao, A. Fang, Y. Huang, V. Giunchiglia, A. Noori, J. R. Schwarz, Y. Ektefaie, J. Kondic, and M. Zitnik. Empowering biomedical discovery with AI agents. *Cell*, 187(22):6125–6151, Oct. 2024. ISSN 0092-8674, 1097-4172. doi: 10.1016/j.cell.2024.09.022. URL [https://www.cell.com/cell/abstract/S0092-8674\(24\)01070-5](https://www.cell.com/cell/abstract/S0092-8674(24)01070-5). Publisher: Elsevier.
- A. Grattafiori, A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Vaughan, A. Yang, A. Fan, A. Goyal, A. Hartshorn, A. Yang, A. Mitra, A. Srivankumar, A. Korenev, A. Hinsvark, A. Rao, A. Zhang, A. Rodriguez, A. Gregerson, A. Spataru,

540 B. Roziere, B. Biron, B. Tang, B. Chern, C. Caucheteux, C. Nayak, C. Bi, C. Marra, C. McConnell,
541 C. Keller, C. Touret, C. Wu, C. Wong, C. C. Ferrer, C. Nikolaidis, D. Allonsius, D. Song, D. Pintz,
542 D. Livshits, D. Wyatt, D. Esiobu, D. Choudhary, D. Mahajan, D. Garcia-Olano, D. Perino, D. Hup-
543 kes, E. Lakomkin, E. AlBadawy, E. Lobanova, E. Dinan, E. M. Smith, F. Radenovic, F. Guzmán,
544 F. Zhang, G. Synnaeve, G. Lee, G. L. Anderson, G. Thattai, G. Nail, G. Mialon, G. Pang, G. Cu-
545 curell, H. Nguyen, H. Korevaar, H. Xu, H. Touvron, I. Zarov, I. A. Ibarra, I. Kloumann, I. Misra,
546 I. Evtimov, J. Zhang, J. Copet, J. Lee, J. Geffert, J. Vranes, J. Park, J. Mahadeokar, J. Shah,
547 J. v. d. Linde, J. Billock, J. Hong, J. Lee, J. Fu, J. Chi, J. Huang, J. Liu, J. Wang, J. Yu, J. Bitton,
548 J. Spisak, J. Park, J. Rocca, J. Johnstun, J. Saxe, J. Jia, K. V. Alwala, K. Prasad, K. Upasani,
549 K. Plawiak, K. Li, K. Heafield, K. Stone, K. El-Arini, K. Iyer, K. Malik, K. Chiu, K. Bhalla,
550 K. Lakhotia, L. Rantala-Yearly, L. v. d. Maaten, L. Chen, L. Tan, L. Jenkins, L. Martin, L. Madaan,
551 L. Malo, L. Blecher, L. Landzaat, L. d. Oliveira, M. Muzzi, M. Pasupuleti, M. Singh, M. Paluri,
552 M. Kardas, M. Tsimpoukelli, M. Oldham, M. Rita, M. Pavlova, M. Kambadur, M. Lewis, M. Si,
553 M. K. Singh, M. Hassan, N. Goyal, N. Torabi, N. Bashlykov, N. Bogoychev, N. Chatterji, N. Zhang,
554 O. Duchenne, O. Çelebi, P. Alrassy, P. Zhang, P. Li, P. Vasic, P. Weng, P. Bhargava, P. Dubal, P. Kr-
555 ishnan, P. S. Koura, P. Xu, Q. He, Q. Dong, R. Srinivasan, R. Ganapathy, R. Calderer, R. S. Cabral,
556 R. Stojnic, R. Raileanu, R. Maheswari, R. Girdhar, R. Patel, R. Sauvestre, R. Polidoro, R. Sumbaly,
557 R. Taylor, R. Silva, R. Hou, R. Wang, S. Hosseini, S. Chennabasappa, S. Singh, S. Bell, S. S. Kim,
558 S. Edunov, S. Nie, S. Narang, S. Rapparth, S. Shen, S. Wan, S. Bhosale, S. Zhang, S. Vandenhende,
559 S. Batra, S. Whitman, S. Sootla, S. Collot, S. Gururangan, S. Borodinsky, T. Herman, T. Fowler,
560 T. Sheasha, T. Georgiou, T. Scialom, T. Speckbacher, T. Mihaylov, T. Xiao, U. Karn, V. Goswami,
561 V. Gupta, V. Ramanathan, V. Kerkez, V. Gonguet, V. Do, V. Vogeti, V. Albiero, V. Petrovic, W. Chu,
562 W. Xiong, W. Fu, W. Meers, X. Martinet, X. Wang, X. Wang, X. E. Tan, X. Xia, X. Xie, X. Jia,
563 X. Wang, Y. Goldschlag, Y. Gaur, Y. Babaei, Y. Wen, Y. Song, Y. Zhang, Y. Li, Y. Mao, Z. D.
564 Coudert, Z. Yan, Z. Chen, Z. Papakipos, A. Singh, A. Srivastava, A. Jain, A. Kelsey, A. Shajnfeld,
565 A. Gangidi, A. Victoria, A. Goldstand, A. Menon, A. Sharma, A. Boesenberg, A. Baevski, A. Fein-
566 stein, A. Kallet, A. Sangani, A. Teo, A. Yunus, A. Lupu, A. Alvarado, A. Caples, A. Gu, A. Ho,
567 A. Poulton, A. Ryan, A. Ramchandani, A. Dong, A. Franco, A. Goyal, A. Saraf, A. Chowdhury,
568 A. Gabriel, A. Bharambe, A. Eisenman, A. Yazdan, B. James, B. Maurer, B. Leonhardi, B. Huang,
569 B. Loyd, B. D. Paola, B. Paranjape, B. Liu, B. Wu, B. Ni, B. Hancock, B. Wasti, B. Spence,
570 B. Stojkovic, B. Gamido, B. Montalvo, C. Parker, C. Burton, C. Mejia, C. Liu, C. Wang, C. Kim,
571 C. Zhou, C. Hu, C.-H. Chu, C. Cai, C. Tindal, C. Feichtenhofer, C. Gao, D. Civin, D. Beaty,
572 D. Kreymer, D. Li, D. Adkins, D. Xu, D. Testuggine, D. David, D. Parikh, D. Liskovich, D. Foss,
573 D. Wang, D. Le, D. Holland, E. Dowling, E. Jamil, E. Montgomery, E. Presani, E. Hahn, E. Wood,
574 E.-T. Le, E. Brinkman, E. Arcaute, E. Dunbar, E. Smothers, F. Sun, F. Kreuk, F. Tian, F. Kokkinos,
575 F. Ozgenel, F. Caggioni, F. Kanayet, F. Seide, G. M. Florez, G. Schwarz, G. Badeer, G. Swee,
576 G. Halpern, G. Herman, G. Sizov, Guangyi, Zhang, G. Lakshminarayanan, H. Inan, H. Shojanazeri,
577 H. Zou, H. Wang, H. Zha, H. Habeeb, H. Rudolph, H. Suk, H. Aspegren, H. Goldman, H. Zhan,
578 I. Damlaj, I. Molybog, I. Tufanov, I. Leontiadis, I.-E. Veliche, I. Gat, J. Weissman, J. Geboski,
579 J. Kohli, J. Lam, J. Asher, J.-B. Gaya, J. Marcus, J. Tang, J. Chan, J. Zhen, J. Reizenstein, J. Teboul,
580 J. Zhong, J. Jin, J. Yang, J. Cummings, J. Carvill, J. Shepard, J. McPhie, J. Torres, J. Ginsburg,
581 J. Wang, K. Wu, K. H. U, K. Saxena, K. Khandelwal, K. Zand, K. Matosich, K. Veeraraghavan,
582 K. Michelena, K. Li, K. Jagadeesh, K. Huang, K. Chawla, K. Huang, L. Chen, L. Garg, L. A.
583 L. Silva, L. Bell, L. Zhang, L. Guo, L. Yu, L. Moshkovich, L. Wehrstedt, M. Khabsa, M. Avalani,
584 M. Bhatt, M. Mankus, M. Hasson, M. Lennie, M. Reso, M. Groshev, M. Naumov, M. Lathi,
585 M. Keneally, M. Liu, M. L. Seltzer, M. Valko, M. Restrepo, M. Patel, M. Vyatskov, M. Samvelyan,
586 M. Clark, M. Macey, M. Wang, M. J. Hermoso, M. Metanat, M. Rastegari, M. Bansal, N. San-
587 thanam, N. Parks, N. White, N. Bawa, N. Singhal, N. Egebo, N. Usunier, N. Mehta, N. P. Laptev,
588 N. Dong, N. Cheng, O. Chernoguz, O. Hart, O. Salpekar, O. Kalinli, P. Kent, P. Parekh, P. Saab,
589 P. Balaji, P. Rittner, P. Bontrager, P. Roux, P. Dollar, P. Zvyagina, P. Ratanchandani, P. Yuvraj,
590 Q. Liang, R. Alao, R. Rodriguez, R. Ayub, R. Murthy, R. Nayani, R. Mitra, R. Parthasarathy,
591 R. Li, R. Hogan, R. Battey, R. Wang, R. Howes, R. Rinott, S. Mehta, S. Siby, S. J. Bondu,
592 S. Datta, S. Chugh, S. Hunt, S. Dhillon, S. Sidorov, S. Pan, S. Mahajan, S. Verma, S. Yamamoto,
593 S. Ramaswamy, S. Lindsay, S. Lindsay, S. Feng, S. Lin, S. C. Zha, S. Patil, S. Shankar, S. Zhang,
S. Zhang, S. Wang, S. Agarwal, S. Sajuyigbe, S. Chintala, S. Max, S. Chen, S. Kehoe, S. Satterfield,
S. Govindaprasad, S. Gupta, S. Deng, S. Cho, S. Virk, S. Subramanian, S. Choudhury, S. Goldman,
T. Remez, T. Glaser, T. Best, T. Koehler, T. Robinson, T. Li, T. Zhang, T. Matthews, T. Chou,
T. Shaked, V. Vontimitta, V. Ajayi, V. Montanez, V. Mohan, V. S. Kumar, V. Mangla, V. Ionescu,
V. Poenaru, V. T. Mihailescu, V. Ivanov, W. Li, W. Wang, W. Jiang, W. Bouaziz, W. Constable,

594 X. Tang, X. Wu, X. Wang, X. Wu, X. Gao, Y. Kleinman, Y. Chen, Y. Hu, Y. Jia, Y. Qi, Y. Li,
595 Y. Zhang, Y. Zhang, Y. Adi, Y. Nam, Yu, Wang, Y. Zhao, Y. Hao, Y. Qian, Y. Li, Y. He, Z. Rait,
596 Z. DeVito, Z. Rosnbrick, Z. Wen, Z. Yang, Z. Zhao, and Z. Ma. The Llama 3 Herd of Models, Nov.
597 2024. URL <http://arxiv.org/abs/2407.21783>. arXiv:2407.21783 [cs].

598
599 D. Guo, D. Yang, H. Zhang, J. Song, R. Zhang, R. Xu, Q. Zhu, S. Ma, P. Wang, X. Bi, et al.
600 Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint*
601 *arXiv:2501.12948*, 2025.

602 M. Hao, J. Gong, X. Zeng, C. Liu, Y. Guo, X. Cheng, T. Wang, J. Ma, X. Zhang, and
603 L. Song. Large-scale foundation model on single-cell transcriptomics. *Nature Meth-*
604 *ods*, pages 1–11, 2024. URL [https://idp.nature.com/authorize/casa?](https://idp.nature.com/authorize/casa?redirect_uri=https://www.nature.com/articles/s41592-024-02305-7&casa_token=pRbl500caJkAAAAA:FZirFybc1gGR1gHqgRiUyMsPm_-3B8fMt+hLux3Lc7AjTWfp_MSZUTYOfd4TsubiDWqvwphBXOUuBdKa3Q)
605 [redirect_uri=https://www.nature.com/articles/s41592-024-02305-7&](https://idp.nature.com/authorize/casa?redirect_uri=https://www.nature.com/articles/s41592-024-02305-7&casa_token=pRbl500caJkAAAAA:FZirFybc1gGR1gHqgRiUyMsPm_-3B8fMt+hLux3Lc7AjTWfp_MSZUTYOfd4TsubiDWqvwphBXOUuBdKa3Q)
606 [casa_token=pRbl500caJkAAAAA:FZirFybc1gGR1gHqgRiUyMsPm_](https://idp.nature.com/authorize/casa?redirect_uri=https://www.nature.com/articles/s41592-024-02305-7&casa_token=pRbl500caJkAAAAA:FZirFybc1gGR1gHqgRiUyMsPm_-3B8fMt+hLux3Lc7AjTWfp_MSZUTYOfd4TsubiDWqvwphBXOUuBdKa3Q)
607 [-3B8fMt+hLux3Lc7AjTWfp_MSZUTYOfd4TsubiDWqvwphBXOUuBdKa3Q](https://idp.nature.com/authorize/casa?redirect_uri=https://www.nature.com/articles/s41592-024-02305-7&casa_token=pRbl500caJkAAAAA:FZirFybc1gGR1gHqgRiUyMsPm_-3B8fMt+hLux3Lc7AjTWfp_MSZUTYOfd4TsubiDWqvwphBXOUuBdKa3Q). Publisher:
608 Nature Publishing Group US New York.

609 M. Hu, S. Alkhairy, I. Lee, R. T. Pillich, D. Fong, K. Smith, R. Bachelder, T. Ideker, and D. Pratt.
610 Evaluation of large language models for discovery of gene set function. *Nature Methods*, 22
611 (1):82–91, Jan. 2025. ISSN 1548-7105. doi: 10.1038/s41592-024-02525-x. URL [https:](https://www.nature.com/articles/s41592-024-02525-x)
612 [//www.nature.com/articles/s41592-024-02525-x](https://www.nature.com/articles/s41592-024-02525-x). Publisher: Nature Publishing
613 Group.

614 S. Johnson and D. Hyland-Wood. A Primer on Large Language Models and their Limitations, Dec.
615 2024. URL <http://arxiv.org/abs/2412.04503>. arXiv:2412.04503 [cs].

616
617 K. Z. Kedzierska, L. Crawford, A. P. Amini, and A. X. Lu. Assessing the limits of zero-shot
618 foundation models in single-cell biology, Nov. 2023. URL [https://www.biorxiv.org/](https://www.biorxiv.org/content/10.1101/2023.10.16.561085v2)
619 [content/10.1101/2023.10.16.561085v2](https://www.biorxiv.org/content/10.1101/2023.10.16.561085v2). Pages: 2023.10.16.561085 Section: New
620 Results.

621 E. Kernfeld, Y. Yang, J. S. Weinstock, A. Battle, and P. Cahan. A systematic comparison of
622 computational methods for expression forecasting, Oct. 2024. URL [https://www.biorxiv.](https://www.biorxiv.org/content/10.1101/2023.07.28.551039v2)
623 [org/content/10.1101/2023.07.28.551039v2](https://www.biorxiv.org/content/10.1101/2023.07.28.551039v2). Pages: 2023.07.28.551039 Section:
624 New Results.

625 D. P. Kingma and M. Welling. Auto-encoding variational bayes. *Proceedings of the International*
626 *Conference on Learning Representations (ICLR)*, 2014.

627
628 J. Lee, W. Yoon, S. Kim, D. Kim, S. Kim, C. H. So, and J. Kang. Biobert: a pre-trained biomedical
629 language representation model for biomedical text mining. *Bioinformatics*, 36(4):1234–1240,
630 2020.

631 R. Lopez, N. Tagasovska, S. Ra, K. Cho, J. Pritchard, and A. Regev. Learning Causal Repre-
632 sentations of Single Cells via Sparse Mechanism Shift Modeling. Mar. 2023. URL [https:](https://openreview.net/forum?id=IOWJsPJ2xGd)
633 [//openreview.net/forum?id=IOWJsPJ2xGd](https://openreview.net/forum?id=IOWJsPJ2xGd).

634
635 M. Lotfollahi, A. Klimovskaia Susmelj, C. De Donno, L. Hetzel, Y. Ji, I. L. Ibarra, S. R. Srivatsan,
636 M. Naghipourfar, R. M. Daza, B. Martin, J. Shendure, J. L. McFaline-Figueroa, P. Boyeau, F. A.
637 Wolf, N. Yakubova, S. Günemann, C. Trapnell, D. Lopez-Paz, and F. J. Theis. Predicting
638 cellular responses to complex perturbations in high-throughput screens. *Molecular Systems*
639 *Biology*, 19(6):e11517, June 2023. ISSN 1744-4292. doi: 10.15252/msb.202211517. URL
640 <https://www.embopress.org/doi/full/10.15252/msb.202211517>. Publisher:
641 John Wiley & Sons, Ltd.

642
643 K. Märtens, R. Donovan-Maiye, and J. Ferkinghoff-Borg. Enhancing generative perturbation models
644 with llm-informed gene embeddings. In *ICLR 2024 Workshop on Machine Learning for Genomics*
Explorations, 2024.

645
646 J. M. Replogle, R. A. Saunders, A. N. Pogson, J. A. Hussmann, A. Lenail, A. Guna, L. Mascibroda,
647 E. J. Wagner, K. Adelman, and G. Lithwick-Yanai. Mapping information-rich genotype-phenotype
landscapes with genome-scale Perturb-seq. *Cell*, 185(14):2559–2575, 2022. URL [https:](https://www.cell.com/cell/pdf/S0092-8674(22)00597-9.pdf)
[//www.cell.com/cell/pdf/S0092-8674\(22\)00597-9.pdf](https://www.cell.com/cell/pdf/S0092-8674(22)00597-9.pdf). Publisher: Elsevier.

648 D. J. Rezende, S. Mohamed, and D. Wierstra. Stochastic backpropagation and approximate inference
649 in deep generative models. *arXiv preprint arXiv:1401.4082*, 2014.
650

651 Y. Roohani, K. Huang, and J. Leskovec. Predicting transcriptional outcomes of novel multi-
652 gene perturbations with GEARS. *Nature Biotechnology*, pages 1–9, Aug. 2023. ISSN 1546-
653 1696. doi: 10.1038/s41587-023-01905-6. URL [https://www.nature.com/articles/
654 s41587-023-01905-6](https://www.nature.com/articles/s41587-023-01905-6). Publisher: Nature Publishing Group.

655 C. V. Theodoris, L. Xiao, A. Chopra, M. D. Chaffin, Z. R. Al Sayed, M. C. Hill, H. Mantineo, E. M.
656 Brydon, Z. Zeng, and X. S. Liu. Transfer learning enables predictions in network biology. *Nature*,
657 618(7965):616–624, 2023. URL [https://idp.nature.com/authorize/casa?
658 redirect_uri=https://www.nature.com/articles/s41586-023-06139-9&
659 casa_token=T6Kd54XSnhwAAAAA:QHLXOg37bjq_ki4B_rC-zG033wKzhCnWlP5_
660 RI1kq00de9-AGVMDZWj4-_4KC8mKZyhkDkIs8z8Jnlh-Bg](https://idp.nature.com/authorize/casa?redirect_uri=https://www.nature.com/articles/s41586-023-06139-9&casa_token=T6Kd54XSnhwAAAAA:QHLXOg37bjq_ki4B_rC-zG033wKzhCnWlP5_RI1kq00de9-AGVMDZWj4-_4KC8mKZyhkDkIs8z8Jnlh-Bg). Publisher: Nature
661 Publishing Group UK London.

662 E. Wang, S. Schmidgall, P. F. Jaeger, F. Zhang, R. Pilgrim, Y. Matias, J. Barral, D. Fleet, and
663 S. Azizi. TxGemma: Efficient and Agentic LLMs for Therapeutics, Apr. 2025. URL [http:
664 //arxiv.org/abs/2504.06196](http://arxiv.org/abs/2504.06196). arXiv:2504.06196 [cs].
665

666 D. R. Wong, A. S. Hill, and R. Moccia. Simple controls exceed best deep learning algorithms
667 and reveal foundation model effectiveness for predicting genetic perturbations, Jan. 2025. URL
668 <https://www.biorxiv.org/content/10.1101/2025.01.06.631555v2>. Pages:
669 2025.01.06.631555 Section: New Results.

670 M. Wu, R. Littman, J. Levine, L. Qiu, T. Biancalani, D. Richmond, and J.-C. Huetter. Contextualizing
671 biological perturbation experiments through language. Oct. 2024. URL [https://openreview.
672 net/forum?id=5WEpbilssv](https://openreview.net/forum?id=5WEpbilssv).
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

702 SUPPLEMENTARY MATERIAL

703
704 A PROMPTS

705
706 The prompts below aim to identify approximately 5-10 genes from a provided list that closely resemble
707 the specified gene based on shared involvement in specific biological pathways, co-regulation, or
708 protein-protein interactions in the context of a given cell line (K562 or RPE1).
709

710 **Prompt for K562 cell line**

711 Instruction: Analyze the gene {gene} and identify 5-10 most similar
712 genes from the provided list. Rank them by similarity (most similar
713 first).

714 Consider similarity based on:

- 715 • Shared biological pathways and functions
- 716 • Protein-protein interactions and complex formation
- 717 • Co-regulation patterns and co-expression relationships
- 718 • Similar effects when knocked out (predicted knockout outcomes)

718 Context: Analysis should focus on the K562 cell line (chronic myeloid
719 leukemia model). Consider cancer-relevant pathways including ribosome
720 biogenesis, transcriptional regulation, mitochondrial function, stress
721 responses, and known genetic dependencies specific to K562 survival.

722 Available genes: {list_of_genes}

723 Format your response as JSON with two parts:

- 724 1. "reasoning": Explain your analysis, discussing potential
725 connections between {gene} and relevant genes. Include both general
726 biological knowledge and cell-specific considerations where relevant.
- 727 2. "kNN": List the most similar genes in order of similarity (5-10
728 genes)

729 Example response format:

```
730 { "reasoning": "Gene X is involved in pathway Y which  
731 directly interacts with gene Z...", "kNN": ["MostSimilarGene",  
732 "SecondMostSimilar", "ThirdMostSimilar"] }
```

733 Note: Prioritize confidence over hitting exactly 5-10 genes. Provide
734 ONLY ONE JSON response.

735 **Prompt for RPE1 cell line**

736 Instruction: Analyze the gene {gene} and identify 5-10 most similar
737 genes from the provided list. Rank them by similarity (most similar
738 first).

739 Consider similarity based on:

- 740 • Shared biological pathways and functions
- 741 • Protein-protein interactions and complex formation
- 742 • Co-regulation patterns and co-expression relationships
- 743 • Similar effects when knocked out (predicted knockout outcomes)

744 Context: Analysis should focus on the RPE1 cell line (near-diploid,
745 non-transformed human retinal pigment epithelial cells). Consider
746 pathways relevant to RPE1 biology including epithelial polarity, cilia
747 signaling, cell cycle control, DNA repair, oxidative stress response,
748 and retinal metabolism.

749 Available genes: {list_of_genes}

- 750 Format your response as JSON with two parts: 1. "reasoning": Explain
751 your analysis, discussing potential connections between {gene} and
752 relevant genes. Include both general biological knowledge and
753 RPE1-specific considerations where relevant. 2. "kNN": List the most
754 similar genes in order of similarity (5-10 genes)

755 Example response format: { "reasoning": "Gene X is involved in pathway
756 Y which directly interacts with gene Z...", "kNN": ["MostSimilarGene",
757 "SecondMostSimilar", "ThirdMostSimilar"] }

758 Note: Prioritize confidence over hitting exactly 5-10 genes. Provide
759 ONLY ONE JSON response.