



CELLFORGE: AGENTIC DESIGN OF VIRTUAL CELL MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Virtual cell modeling aims to predict cellular responses to diverse perturbations but faces challenges from biological complexity, multimodal data heterogeneity, and the need for interdisciplinary expertise. We introduce CELLFORGE, a multi-agent *framework* that autonomously designs and synthesizes neural network architectures tailored to specific single-cell datasets and perturbation tasks. Given raw multi-omics data and task descriptions, CELLFORGE discovers candidate architectures through collaborative reasoning among specialized agents, then generates executable implementations. Our core contribution is the *framework itself*: showing that multi-agent collaboration mechanisms—rather than manual human design or single-LLM prompting—can autonomously produce executable, high-quality computational methods. This approach goes beyond conventional hyperparameter tuning by enabling entirely new architectural **components such as** trajectory-aware encoders and perturbation diffusion modules to *emerge from agentic deliberation*. We evaluate CELLFORGE on six datasets spanning gene knockouts, drug treatments, and cytokine stimulations across multiple modalities (scRNA-seq, scATAC-seq, CITE-seq). The results demonstrate that the models generated by CELLFORGE are highly competitive with established baselines, while revealing systematic patterns of architectural innovation. CELLFORGE highlights the scientific value of multi-agent frameworks: collaboration among specialized agents enables genuine methodological innovation and executable solutions that single agents or human experts cannot achieve. This represents a paradigm shift toward autonomous *scientific method development* in computational biology. Code is available at <https://anonymous.4open.science/r/CellForge-FC93/>.

1 INTRODUCTION

Scientific discovery in computational biology increasingly demands interdisciplinary expertise spanning machine learning, statistics, and domain knowledge [20, 47, 79, 103]. While recent advances in large language models have enabled AI systems to excel at individual research tasks, from literature analysis [43] to hypothesis generation [88, 116], integrating these capabilities into complete scientific workflows remains challenging [66, 75]. This gap is particularly evident in virtual cell modeling, where researchers must design computational methods that capture complex biological mechanisms across diverse experimental conditions [10, 97].

Virtual cell modeling aims to predict how cells respond to genetic edits, chemical treatments, and environmental perturbations across multiple biological modalities [13, 98]. While foundation models like scGPT [21] and Geneformer [110] have advanced single-cell analysis, they often struggle with dataset-specific perturbation patterns and experimental nuances. Creating effective predictive models typically requires extensive manual effort to integrate domain knowledge, design appropriate architectures, and validate results empirically [74, 85]. **Perturbation datasets differ substantially in mechanisms, modalities, and sparsity regimes, creating different learning problems that benefit from dataset-specific inductive biases.**

We present CELLFORGE, an agentic *framework* designed to autonomously create computational methods for virtual cell modeling. **Its core contribution is showing that multi-agent collaboration mechanisms, including graph-structured discussions with confidence scoring and iterative refinement, enable autonomous scientific method development that outperforms single-agent prompting.** Rather

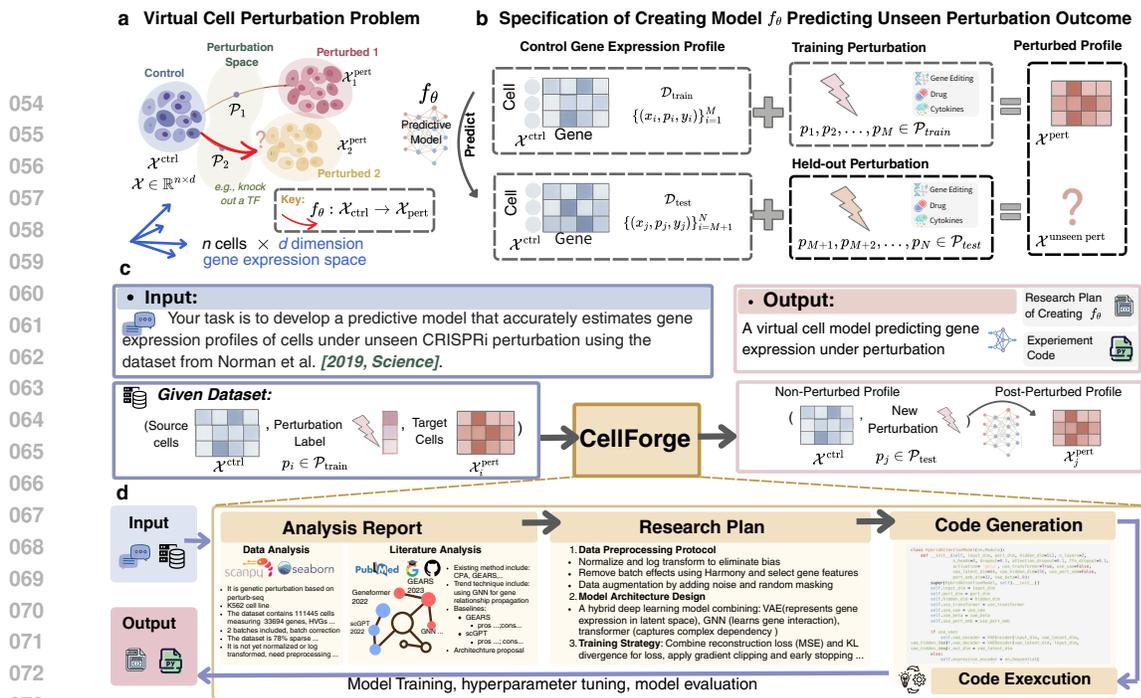


Figure 1: (a) Perturbation prediction learns mappings from control cell states to post-perturbation states in high-dimensional expression space. (b) Models train on control-perturbed cell pairs across modalities (scRNA-seq, scATAC-seq, CITE-seq) to predict responses to unseen perturbations. (c) CELLFORGE receives datasets and task descriptions, autonomously designing models for predicting expression under novel perturbations ($p_i \in \mathcal{P}_{\text{test}}$). (d) System workflow.

than selecting from predefined pipelines, CELLFORGE generates novel deep learning architectures through emergent collaborative reasoning among specialized agents. Architectures such as trajectory-aware encoders with perturbation diffusion modules serve as evidence of the framework’s creative synthesis, but the primary novelty is the framework itself.

CELLFORGE operates through three modules that address distinct research challenges: Task Analysis agents profile datasets and extract design principles from literature through alternating breadth-first and depth-first retrieval; Design agents engage in graph-structured discussions where specialized experts collaboratively propose, critique, and refine architectures until convergence; and Experiment Execution agents translate research plans into production-ready code with automated debugging and iterative refinement. This structured collaboration enables discovery of optimized models that individual agents cannot achieve, while also revealing systematic patterns in architectural innovation.

We evaluate CELLFORGE on single-cell perturbation prediction across six datasets and benchmark against a comprehensive set of baselines. Our results show that models produced by CELLFORGE are highly competitive and frequently surpass established baselines, though performance varies across runs due to the stochastic nature of automated design. Through systematic analysis, we highlight scenarios where multi-agent collaboration yields advantages, examine the kinds of architectural innovations that emerge from the framework, and delineate the boundaries of current automated design performance. Overall, CELLFORGE positions multi-agent frameworks as a general approach for automating scientific method development in computational biology, moving beyond isolated task execution toward end-to-end autonomous research workflows.

2 RELATED WORK

Single-Cell Perturbation Analysis Virtual cell modeling, the computational simulation of cellular responses to perturbations, represents a fundamental challenge in systems biology [13, 98]. Existing approaches span diverse paradigms: early methods [22, 104] treat genes independently; deep generative models [41, 72, 74] model perturbations as latent space transformations; network-based methods [7, 90, 96] incorporate gene regulatory knowledge; optimal transport approaches [12, 51] align control and perturbed cell distributions; and transformers [21, 38, 110] leverage large-scale pretraining. This diversity creates a vast design space where architecture selection remains highly context-dependent, typically requiring extensive domain expertise.

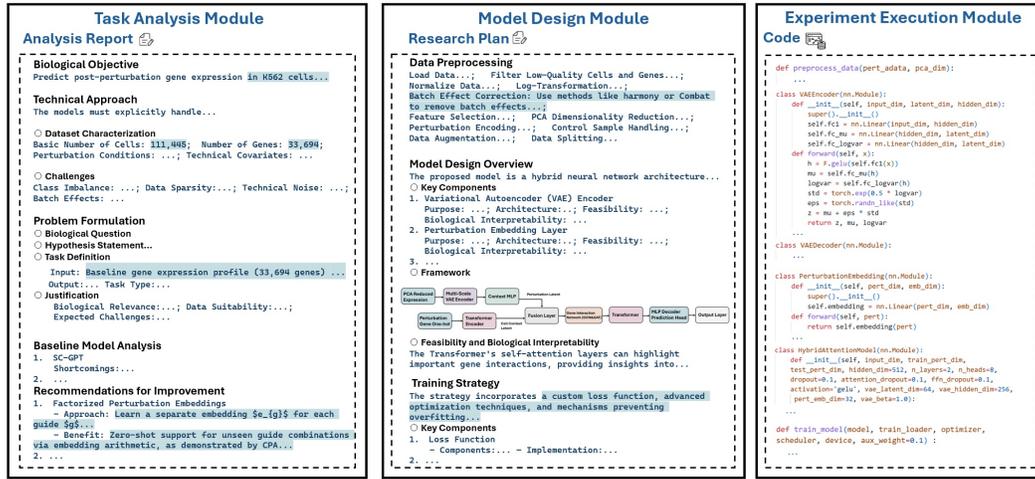


Figure 2: Multi-agent collaboration generates scientific research artifacts. Task Analysis produces dataset characterization and literature-grounded insights, Design Module synthesizes novel methodological approaches through structured agent discussions, and Experiment Execution demonstrates code generation capability.

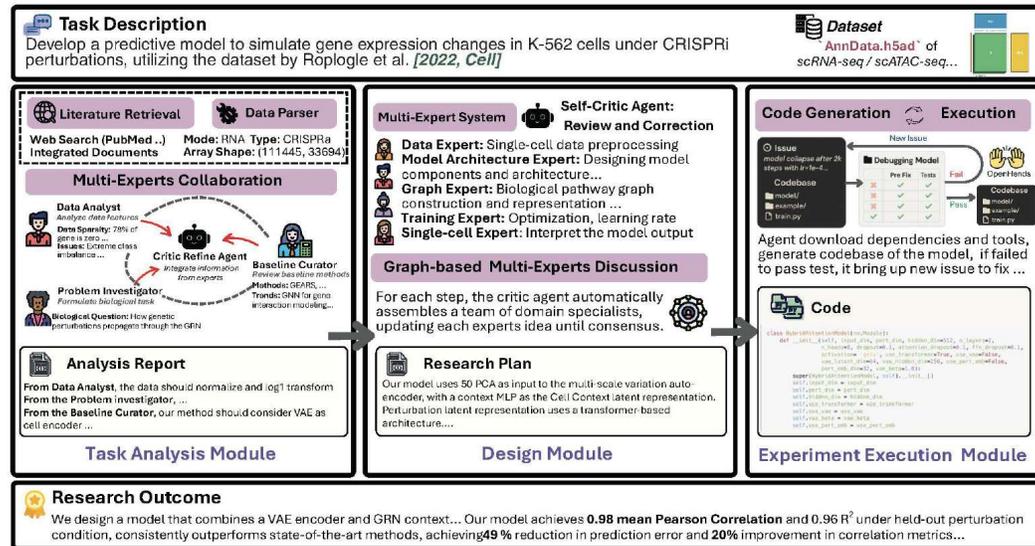


Figure 3: The CELFORGE architecture and workflow.

AI Agents in Biomedical Research Agentic systems are transforming biomedical discovery across the research pipeline [25, 39, 45]. Current approaches fall into three categories: reasoning-focused agents that interpret biological phenomena but lack code generation capabilities; experimental design agents that optimize wet-lab protocols rather than computational models; and workflow automation systems [46, 54, 112] constrained by predefined toolsets. While these advances demonstrate agentic potential, none address autonomous computational model design the systematic creation of novel architectures tailored to specific biological datasets and research objectives. This gap represents a critical opportunity for agentic frameworks that can autonomously navigate the complex design space of computational biology methods.

3 PRELIMINARY AND BACKGROUND

Notations. Let $X \in \mathbb{R}^{n \times d}$ denote the matrix of single-cell profiles, where n is the number of cells and d is the feature dimensionality (typically $\sim 20,000$ genes for scRNA-seq, chromatin peaks for scATAC-seq, or protein markers for CITE-seq). For virtual cell modeling, we are given dataset $\mathcal{D} = \{(x_i, p_i, y_i)\}_{i=1}^N$ and task description S , where $x_i \in \mathbb{R}^d$ is the pre-perturbation profile, $p_i \in \mathcal{P}$ is the applied perturbation (gene knockout, drug, cytokine), and $y_i \in \mathbb{R}^d$ is the post-perturbation profile. We partition \mathcal{D} into training $\mathcal{D}_{\text{train}}$ and test $\mathcal{D}_{\text{test}}$ sets, where test perturbations $p_i \in \mathcal{P}_{\text{test}} \subset \mathcal{P}$ are held-out during training to evaluate generalization to novel perturbations and cellular contexts.

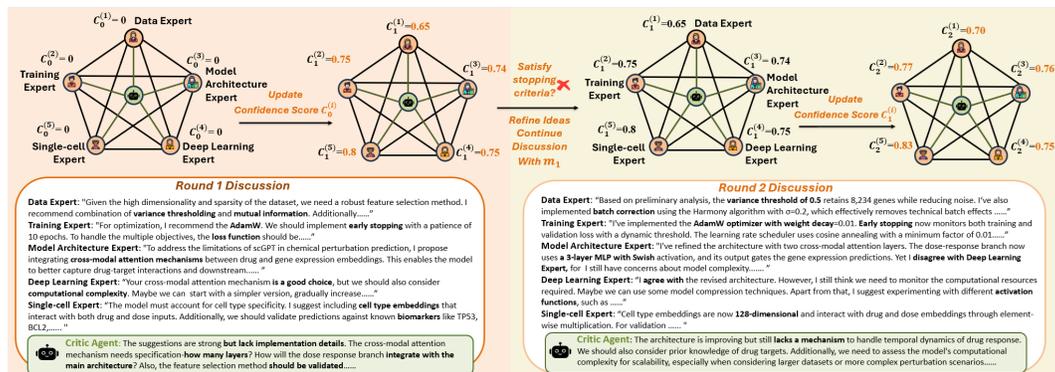


Figure 4: **The Graph-based discussion architecture and workflow.** This is an example of two rounds of discussion from the beginning. After each round, confidence scores are updated, and the agentic system will judge if the current state satisfies the stopping criteria. If not, each expert will refine their ideas based on the critic agent’s suggestions and other experts’ viewpoints. This graph-based critic refinement continues until reaching the termination state. The figure includes an example formula for computing each experts’ confidence score per round, based on a weighted combination of historical scores, peer evaluations, and critic agent’s assessments. Complete multi-rounds of discussions are presented in Appendix G.2.

Problem Formulation. We formalize perturbation prediction as learning mapping function $f_\theta : \mathbb{R}^d \times \mathcal{P} \rightarrow \mathbb{R}^{d'}$ that generalizes to unseen perturbations and cell states, where θ represents trainable parameters. To capture cell-state structure, we incorporate learnable encoders $g_\phi : \mathbb{R}^d \rightarrow \mathbb{R}^h$ producing latent embeddings $z_i = g_\phi(x_i)$ that preserve geometric relationships between control and perturbed states. Importantly, CELLFORGE learns perturbation-response mappings *de novo* for each dataset without importing pretrained representations, capturing dataset-specific perturbation signatures and experimental nuances.

Evaluation. We assess $f_\theta(x_i, p_i)$ for all $(x_i, p_i) \in \mathcal{X}_{\text{test}} \times \mathcal{P}_{\text{test}}$ using mean squared error, Pearson correlation, and perturbation consistency metrics adapted from [9, 96] to ensure biological significance (detailed in Appendix E).

4 METHOD

CELLFORGE profiles each study’s perturbations, modalities, and data characteristics to design tailored architectures through knowledge-guided analysis rather than exhaustive search (Figure 3). Task Analysis modules diagnose problems while Design modules synthesize solutions through structured multi-agent collaboration. Agents’ output maintains reasoning traceability throughout the discovery process (detailed are provided in Appendix D.3, configurations are detailed in Appendix D.2, communication protocols are outlined in Appendix H, prompts are included in Appendix R, and detailed outputs and logs are available at https://anonymous.4open.science/r/CellForge_egs-1E54/.

4.1 TASK ANALYSIS MODULE

The Task Analysis module autonomously characterizes datasets and discovers architectural design principles through three sequential components: **(1) Data Parser.** Extracts experimental metadata across modalities (scRNA-seq, scATAC-seq, CITE-seq), including perturbation types, gene features, and cellular contexts. The component standardizes heterogeneous experimental information and generates foundational statistics without human intervention (Appendix S.1). **(2) Literature Retrieval.** Combines static knowledge (46 single-cell perturbation related articles, listed in Appendix O) with dynamic PubMed searches using alternating breadth-first and depth-first strategies. Starting with query $Q^{(0)}$, the system employs Sentence-BERT embeddings where BFS layers retrieve diverse concepts ($\mathcal{N}_t = \text{TopK}(Q^{(t)})$) and DFS layers follow promising paths. Document relevance scoring via $\text{Score}(Q, d) = \frac{e(Q) \cdot e(d)}{\|e(Q)\| \|e(d)\|}$ guides systematic exploration of architectural design spaces (detailed algorithms in Appendix G.1). **(3) Multi-Experts Collaboration.** Specialized agents (Dataset Analyst, Problem Investigator, Baseline Assessor) combine retrieved insights to propose architectural innovations beyond established paradigms. Instead of recombining known modules, they extract design principles and synthesize new components such as trajectory-aware encoders for temporal dynamics or perturbation diffusion modules for combinatorial interventions. The Baseline Assessor

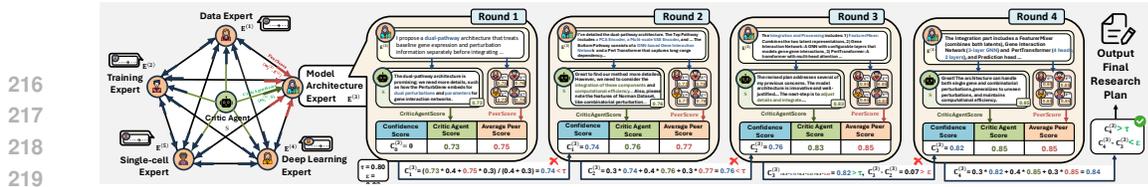


Figure 5: **Confidence Score Update in Graph-based Expert Discussion.** This figure illustrates an example of how a domain experts confidence score evolves during iterative rounds of discussion in the Graph-based Expert Discussion framework. While this example focuses on the Model Architecture Expert, the same confidence updating process applies to all participating experts in the graph, each iteratively refining their proposals and adjusting their confidence based on multi-agent evaluations.

grounds proposals in theoretical analysis across diverse deep learning paradigms, supporting principled innovation. This process yields dataset-tailored modules that emerge from creative integration of biological insights and computational methods, rather than systematic enumeration.

4.2 DESIGN MODULE

The Design module implements scientific creativity through graph-based multi-agent collaboration, generating integrated research plans encompassing preprocessing strategies, architectural designs, and implementation details. The core innovation is the autonomous discovery of optimized architecture rather than the hyperparameter tuning of fixed designs. This architectural discovery process tailoring neural components to dataset-specific biological characteristics constitutes the primary source of CELLFORGE’s performance advantages in perturbation prediction.

Multi-Expert Critic System. We construct a panel of domain experts through role-play prompting: each expert is instantiated from similar dedicated prompt templates that encode its specialty while using the same underlying LLM. See Appendix R.3 for the full templates. For each task, the system dynamically selects a subset of domain experts $E^{(k)}$ (e.g., Data Expert, Single-Cell Expert, Deep Learning Expert) based on task requirements, along with a permanent critic agent S . These agents form an undirected collaboration graph $G^{(k)} = (S, E^{(k)})$, where each expert node maintains a confidence score $c_t^{(i)}$ that evolves through discussion rounds, where t is the discussion round and i represents different domain experts.

Graph-based Discussion. The framework runs up to $T_{\max} = 10$ rounds of graph-based message passing, where experts propose architectural solutions. In each round t every expert $E^{(i)}$ proposes an architectural candidate $m_t^{(i)}$. After all proposals are submitted, a *critic agent* S reviews every $m_t^{(i)}$, summarizes strengths and weaknesses, and assigns a score.

At the end of round t the value is updated by both the critic agent and peer experts. Specifically, the confidence score $c_t^{(i)}$ for expert i at round t is computed as: $c_t^{(i)} = \lambda_1 \cdot c_{t-1}^{(i)} + \lambda_2 \cdot \text{CriticAgentScore}(m_t^{(i)}, S) + \lambda_3 \cdot \frac{1}{k-1} \sum_{j \neq i} \text{PeerScore}(m_t^{(i)}, E^{(j)})$, where $c_{t-1}^{(i)}$ represents the historical confidence, $\text{CriticAgentScore}(m_t^{(i)}, S)$ evaluates the scientific rigor and feasibility of proposal $m_t^{(i)}$ by the critic agent S , $\text{PeerScore}(m_t^{(i)}, E^{(j)})$ captures the evaluation from peer expert j , k is the total number of participating experts, and $(\lambda_1, \lambda_2, \lambda_3) = (0.3, 0.4, 0.3)$ are empirically determined weights with $\lambda_1 + \lambda_2 + \lambda_3 = 1$. The discussion ends when all experts’ confidence scores exceed the threshold $\tau = 0.8$ with minimal variance ($\max_{i,j} |c_{t^*}^{(i)} - c_{t^*}^{(j)}| < \epsilon$, $\epsilon = 0.03$), where t^* represents the final round when the discussion ends, i and j represent domain experts.

If this condition is not met, otherwise it stops at the round limit T_{\max} to balance computational cost, inference time, and token consumption. Before reaching the ending criteria, experts refine their proposals using historical context and proceed to the next round. This process ensures convergence toward scientifically valid and technically feasible model designs with explicit reasoning chains throughout several rounds of discussion. Further information on expert selection and discussion construction is in Appendix D.4, detailed algorithm and mathematical formulation are presented in Appendix G.2, and hyperparameter configuration is presented in Appendix D.7.

4.3 EXPERIMENT EXECUTION MODULE

The Experiment Execution module turns high-level research plans into fully tested, empirically validated results:

Table 1: Post-perturbation gene expression prediction results on datasets where multiple existing baseline models are available. The reported metrics for *CellForge-Models* are shown as mean \pm standard deviation across three automatically designed models. **Ranking markers are determined by best-case bounds:** for metrics where lower is better, we use mean $-$ std; for metrics where higher is better, we use mean $+$ std. All baseline methods are reproduced on the corresponding dataset under the unseen perturbation setting.

MODEL	$MSE \downarrow$	$PCC \uparrow$	$R^2 \uparrow$	$MSE_{DE} \downarrow$	$PCC_{DE} \uparrow$	$R^2_{DE} \uparrow$
<i>Gene Knock Out Perturbation – scRNAseq Dataset (Adamson et al. [2])</i>						
Unperturbed	0.9840	0.0001	-0.0127	3.7865	0.0012	-4.2437
Random Forest	0.3053	0.2063	0.0504	0.5923	0.2632	0.1653
Linear Regression	0.5803	0.0026	0.0435	0.6995	0.0257	0.1074
CPA [73]	0.0067 ³	0.9833 ³	0.9845²	0.1447 ³	0.9024	0.8896
scGen [72]	0.0082	0.9805	0.9611	0.1301 ²	0.8994	0.7263
CondOT [11]	0.0062	0.9608	0.9740	0.1997	0.9341 ²	0.9002 ²
Biolord [86]	0.0044 ²	0.7799	0.9844 ³	0.1256 ²	0.9097	0.9276²
scGPT [21]	0.0100	0.9861	0.9649	0.2562	0.9088	0.7911
CellForge-Models	0.0051 \pm 0.0063 ¹	0.9883 \pm 0.0459¹	0.9761 \pm 0.0803 ¹	0.2013 \pm 0.0444 ¹	0.9474 \pm 0.0601¹	0.8912 \pm 0.0518 ¹
<i>Gene Knock Out Perturbation – scRNAseq Dataset (Norman et al. [82])</i>						
Unperturbed	0.9251	0.0000	-0.1738	5.1214	-0.0021	-4.2047
Random Forest	0.4059	0.1625	0.0623	0.6817	0.1428	0.0498
Linear Regression	0.4989	0.0244	0.0314	0.7331	0.0265	0.0238
CPA [73]	0.0051 ³	0.9779	0.9603	0.3400	0.5754	0.4555
scGen [72]	0.0053	0.9221	0.9521	0.3877	0.5605	0.3220
CondOT [11]	0.0420	0.9847 ²	0.9619	0.2791 ³	0.8022	0.7470
Biolord [86]	0.0027 ²	0.4374	0.9830 ²	0.2450 ²	0.4646	0.8112 ²
scGPT [21]	0.0076	0.9823	0.9536	0.5318	0.8630 ²	0.5652
CellForge-Models	0.0034 \pm 0.0023 ¹	0.9846 \pm 0.0418¹	0.9609 \pm 0.0081 ¹	0.1736 \pm 0.0677 ¹	0.8109 \pm 0.0133 ¹	0.5975 \pm 0.0539 ¹
<i>Drug Perturbation – scRNA-seq Dataset (Srivatsan et al. [106])</i>						
Unperturbed	0.8919	0.0002	-2.4282	9.3326	0.0077	-6.8585
Random Forest	0.5289	0.0527	0.0986	0.6138	0.0245	0.0817
Linear Regression	0.6703	0.0711	0.2826	0.5625	0.0763	0.0421
ChemCPA [41]	0.0847	0.7221	0.6930	0.1035 ³	0.8053	0.7412
scGen [72]	0.0579	0.7871	0.7334	0.1263	0.6575	0.5610
CondOT [11]	0.0499	0.8674	0.6531	0.0933 ²	0.8341	0.4378
Biolord [86]	0.0011 ²	0.9658	0.9287	0.0162 ²	0.9283²	0.8236
CellFlow [59]	0.0003 ¹	0.9906¹	0.9813¹	0.0045 ¹	0.7918	0.9794¹
CellForge-Models	0.0053 \pm 0.0290 ³	0.8664 \pm 0.1332 ³	0.8317 \pm 0.0740 ³	0.0080 \pm 0.0835 ³	0.9278 \pm 0.1001 ¹	0.7887 \pm 0.0548 ³

(1) *Code Generation & Self-Debugging.* The Code Generator converts the selected architecture into production-ready scripts and notebooks with complete dependency management. If a syntax or runtime error occurs, the agent receives the traceback via the OpenHands event stream, analyses the failure, patches the code, and re-executes it, repeating until unit tests pass or a rollback to the last stable state is triggered (see Appendix J for a breakdown of resolved error types).

(2) *Training Orchestration.* An automated scheduler launches training with best-practice safeguards: early stopping, cross-validation, adaptive learning-rate schedules, and checkpointing. When the Validation Agent detects under- or over-fitting, it initiates lightweight hyper-parameter tuning (e.g. adjusting regularisation strength or training epochs) to restore convergence. **A brief human semantic check is performed before training to ensure that the generated code corresponds to the correct perturbation-prediction objective. This step does not involve any human design or modification of the architectures, but solely a lightweight semantic check to ensure that the agent-generated code conforms to the intended perturbation-prediction objective.** (3) *Validation, Refinement & Output Assurance.* After each training cycle, the Validation Agent scores checkpoints on MSE, PCC, and R^2 , identifies failure modes, and feeds structured critiques back to the generator. Because the task outputs numerical gene-expression matrices which are always well-formed the focus is on accuracy rather than structural validity.

5 MAIN RESULTS

5.1 EVALUATION SETUP

We evaluate the models designed and implemented by CELLFORGE in various types of perturbation from scPerturb [85], including gene knockouts, drug treatments, and cytokine stimulation in multiple modalities (scRNA-seq, scATAC-seq, CITE-seq). Each dataset represents distinct biological challenges: The Adamson [2] and Norman [82] datasets capture CRISPR gene knockouts in different cell lines, providing fundamental test cases for genetic perturbation. The Srivatsan [106] dataset assesses the prediction of cellular responses to chemical compounds. To rigorously test generalization to **unseen perturbations**, we held out entire perturbation types during training, ensuring that test perturbations (P_{test}) were never observed. We then systematically surveyed prior work and reproduced all baselines applicable under this setting.

Table 2: DEG Recovery Performance Across Benchmark Datasets

Dataset	DEG Recall	ROC-AUC	PR-AUC
<i>Gene Knock Out Perturbation – scRNAseq Datasets</i>			
Adamson et al. [2]	0.695 \pm 0.08	0.652 \pm 0.06	0.285 \pm 0.08
Norman et al. [82]	0.779 \pm 0.13	0.704 \pm 0.05	0.375 \pm 0.07
<i>Drug Perturbation – scRNA-seq Dataset</i>			
Srivatsan et al. [106]	0.689 \pm 0.20	0.646 \pm 0.06	0.182 \pm 0.02

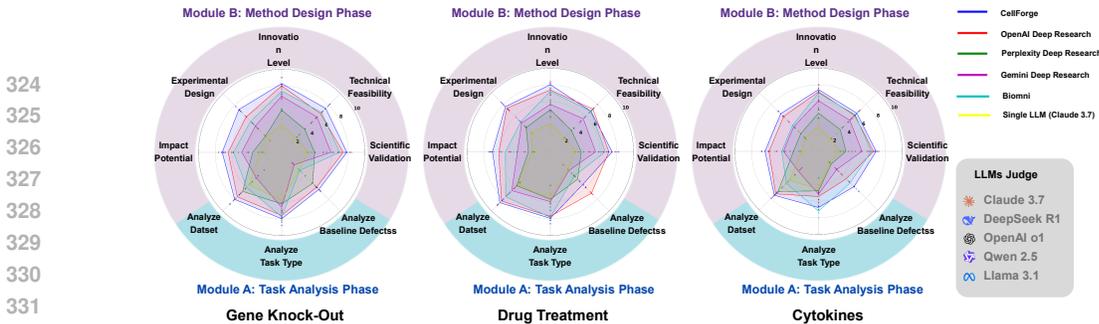


Figure 6: We manually prompt four different DeepResearch variants, Biomni and Single LLM (Claude 3.7) to generate research plans, which were then evaluated by five independent LLMs across eight dimensions, with scores ranging from 1 to 10. Detailed prompts, outputs, and scores are provided in Appendix P. For other modalities where prior perturbation-response methods are scarce or unavailable, including scATAC-seq (Liscovitch et al. [67]) and scCITE-seq (Papalexi et al. [84]), we position our experiments as exploratory applications rather than direct performance comparisons. The results for these datasets, illustrating CELLFORGEs ability to generalize across modalities and design custom architectures, are provided in Appendix B.

5.2 PREDICTIVE PERFORMANCE

Table 1 evaluates the prediction accuracy of models designed by CELLFORGE across diverse perturbation datasets. We report results under two complementary perspectives: overall predictive fidelity and biological relevance. **Overall fidelity** is measured via mean squared error (MSE_{\downarrow}), where lower values indicate predictions closer to actual gene expression; Pearson correlation coefficient (PCC_{\uparrow}), which quantifies how well predicted expression patterns correlate with actual patterns; and coefficient of determination (R^2_{\uparrow}), which measures the proportion of expression variance explained by the model. **Biological relevance** is assessed through metrics on differentially expressed genes (DE), i.e., those showing significant expression changes after perturbation and thus biologically most meaningful. For each dataset, we select the top 20 DE genes based on ground truth perturbation responses and compute the same metrics restricted to this subset (MSE_{DE} , PCC_{DE} , R^2_{DE}).

Across **gene knockout** datasets, CELLFORGE-designed models are competitive with, and in some metrics surpass, strong baselines such as CPA, CondOT, Biolord, and scGPT. On the Adamson dataset, the best-performing CELLFORGE model achieves near-zero MSE and the highest PCC, rivaling Biolord and CPA. On the Norman dataset, it again ranks among the top methods, though performance varies across runs. For **drug perturbations**, the results are more mixed: CellFlow and Biolord remain the strongest overall, while CELLFORGE ranges from near state-of-the-art to substantially weaker depending on the instantiated architecture. In its best-performing configurations, it attains PCC_{DE} close to Biolord and CellFlow; in others, predictive accuracy declines noticeably, reflecting variability introduced by the automated design process.

Taken together, these results show that while CELLFORGE can autonomously generate models that match or exceed hand-designed baselines in some settings, outcomes vary across instantiations. This variability underscores both the promise and the limitations of automated design: CELLFORGE can discover highly competitive models, but careful evaluation across multiple runs remains essential to ensure robustness. As large language models introduce inherent randomness, we provide detailed variability analysis in Appendix J and cross-model comparisons in Appendix K.

5.3 BIOLOGICAL VALIDATION

Beyond overall expression-level accuracy, we assess whether CELLFORGE produces biologically meaningful predictions at multiple levels of resolution. At the gene level, we evaluate recovery of differentially expressed genes (DEGs), a critical signal of perturbation response. Following STAMP [30], we measure DEG recall (sensitivity to true DEGs), ROC-AUC (discriminative power), and PR-AUC (precision under class imbalance). As shown in Table 2, CELLFORGE consistently achieves DEG recall above 0.68 with ROC-AUC values above 0.65 across datasets. On the Norman dataset, performance is relatively stronger, reaching 0.779 recall, indicating effective prioritization of biologically meaningful genes despite the imbalance.

At the pathway and cellular-structure level, enrichment analysis using KEGG annotations shows that the models recover perturbation-relevant pathways: NF-B and p53 signaling for genetic perturbations, autophagy and Wnt pathways for cytokine responses, and coordinated RNAprotein pathways in multimodal CITE-seq data. Complementing this, UMAP visualization demonstrates that predicted

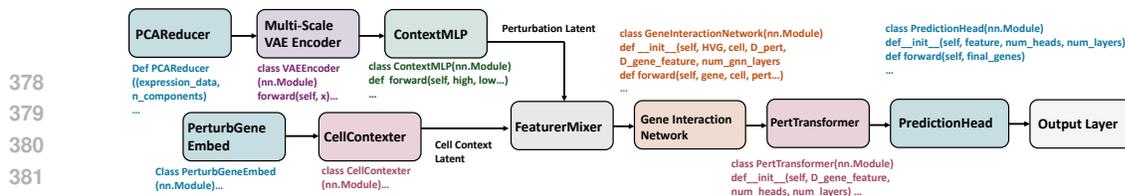


Figure 7: An example diagram of the model framework designed by CELLFORGE on the scRNAseq gene knockout perturbation prediction task ([82].)

cellular states preserve the manifold structure across perturbation types (Appendix M.1), indicating that the models not only capture gene-level perturbation effects but also maintain coherent global organization of cellular states.

5.4 LLM-AS-A-JUDGE AND HUMAN EVALUATION FOR TASK ANALYSIS MODULE AND DESIGN MODULE

To assess the scientific validity of research plans generated by CELLFORGE, we employ a multi-perspective evaluation framework combining automated LLM assessment with independent human expert review. All evaluations are conducted in a randomized, blinded manner.

Our evaluation protocol employs five independent LLM judges from different model families (Claude 3.7, o1, DeepSeek-R1, Qwen-plus, LLaMA3.1) to minimize model-specific biases. Each judge evaluates research plans across eight scientific dimensions: scientific validity, technical feasibility, experimental design quality, biological relevance, innovation level, impact potential, resource efficiency, and methodological rigor (detailed criteria in Appendix I). This methodology follows established LLM-as-judge practices [36, 64].

Three domain experts with extensive single-cell biology experience conducted independent blinded evaluations using identical criteria, each spending approximately 10 hours on assessment. Figure 6 shows that CELLFORGE consistently outperforms DeepResearch variants across scientific validity, innovation level, and experimental design.

Critically, both LLM-assigned scores and agent-generated confidence scores demonstrate strong correlation with human expert evaluations (Pearson $r = 0.83$, $p < 0.01$, Figure 19). This correlation validates that our evaluation framework captures genuine scientific merit rather than stylistic preferences, as domain experts with years of perturbation biology experience would not correlate highly with LLMs based solely on presentation quality.

To address potential concerns about LLM-as-judge evaluation reliability, we conducted comprehensive inter-judge consistency analysis and style-robustness testing. We computed Krippendorff’s and Kendall’s W concordance coefficients across all evaluation dimensions. The results demonstrate strong inter-judge agreement with an average Krippendorff’s of 0.844 ± 0.056 and average inter-judge correlation of 0.925 ± 0.026, indicating high reliability of our evaluation methodology. Detailed evaluation results are presented in P.5.

5.5 NOVEL ARCHITECTURAL DISCOVERIES

CellForge automatically discovers new models that can surpass hand-designed models. An example of a CELLFORGE-designed model framework is shown in Figure 7. Importantly, these architectures were *not* pre-specified or hard-wired: no rule in the code dictates which modules should be chosen. Instead, they emerge from literature retrieval and multi-agent debate, often yielding hybrid or entirely new designs that move beyond simple recombination of known templates (Appendix L). This demonstrates that the system is capable of autonomously internalizing domain knowledge and translating it into new model designs, rather than merely searching or permuting known templates. Because the design process is conditioned on data characteristics and retrieved knowledge rather than fixed heuristics, the pipeline generalizes naturally to new modalities and perturbation types without requiring manual re-engineering.

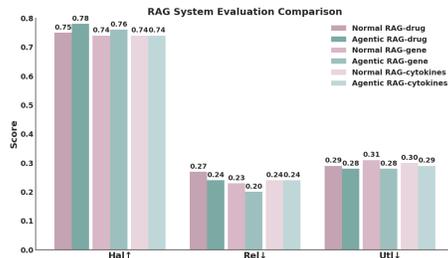


Figure 8: The performance of CellForge’s RAG compared to standard RAG methods. (1) hal: hallucination detection, (2) rel: context relevance, (3) utl: context utilization. Results are stratified by perturbation type (Drug, Cytokine, Gene). Detailed evaluation methods are stated in Appendix F.

CELLFORGE automatically discovered architectures that move beyond standard hand-crafted baselines such as VAE-MLP stacks or plain Transformer encoders. For instance, on the cytokine perturbation dataset *SchiebingerLander2019*, which contains temporal scRNA-seq profiles, CELLFORGE produced a model with three key components. The `TrajectoryAwareEncoder` separates shared versus condition-specific latent dimensions while incorporating temporal embeddings, capturing both global developmental trajectories and cytokine-specific effects. The `PerturbationDiffusionModule` introduces perturbation-conditioned latent diffusion dynamics to represent non-linear, combinatorial interactions. Finally, the `GraphRegularizedDecoder` integrates gene-gene co-regulatory constraints to ensure biologically coherent predictions. Detailed presentation of the novel model components is illustrated in Appendix L.

5.6 RETRIEVAL EFFECTIVENESS

To test whether CELLFORGE effectively retrieves and integrates literature, we evaluated it on RAG-Bench [27] with the PubMedQA dataset [52]. Metrics include hallucination detection (AUROC \uparrow), context relevance (RMSE \downarrow), and context utilization (RMSE \downarrow). As shown in Figure 8, the largest gain occurs in context utilization for gene-perturbation tasks, with consistent performance across perturbation types. This robustness is driven by two key designs: **graph-based expert discussion**, which fuses reasoning paths, and **retrieval-augmented task analysis**, which grounds design in literature while adapting to dataset statistics.

5.7 COMPONENT CONTRIBUTIONS

To disentangle the contributions of individual modules, we performed ablations varying only internal components while holding datasets, tasks, and LLM interface constant. Table 3 shows that adding **Agentic Retrieval** substantially improves performance over the basic version (e.g., Adamson dataset PCC from 0.0087 to 0.5643), while **Graph-Based Discussion** provides complementary gains (to 0.5310). Their combination yields synergistic improvements far exceeding either alone (PCC up to 0.9883), highlighting how knowledge-guided collaborative reasoning drives effective discovery. Similar patterns hold across drug and cytokine perturbations, underscoring the generality of these mechanisms. Notably, our graph-based discussion protocol consistently outperforms alternative approaches: Round-Robin sequential protocols achieve lower performance (Adamson PCC 0.9456 vs. 0.9883), while Moderator-centered evaluation without peer interaction shows further degradation (PCC 0.9123), demonstrating that collaborative expert reasoning is essential for optimal results.

5.8 COSTS AND FAILURE CASES

We also report practical considerations for reproducibility and deployment. Training on two NVIDIA H20 GPUs with a 16-core CPU, 150 GB RAM, and 2 TB SSD typically converges within 38 hours for models of 1030M parameters. Token usage analysis across 50+ experiments shows an average input/output ratio of 60K/300K, with per-request costs averaging \$5.18 (details in Appendix I). Code execution succeeds in roughly 80% of runs; most failures arise from tensor operation errors or invalid configurations, with rarer cases due to hallucinated code or data access issues. The agentic pipeline mitigates many of these through iterative error recovery, further improving robustness (Appendix J).

Multi-agent frameworks face criticism for computational overhead relative to single-LLM approaches. We systematically evaluate CELLFORGE against single-LLM baselines across token consumption, execution time, costs, and success rates, provided in Appendix I.4.

Table 3: Ablation study on the impact of key framework components on designed models’ performance. Detailed settings are listed in D.6.

MODEL	MSE \downarrow	PCC \uparrow	R ² \uparrow	MSE _{DE} \downarrow	PCC _{DE} \uparrow	R ² _{DE} \uparrow
<i>Gene Knock Out Perturbation (Adamson Dataset [2])</i>						
CELLFORGE (Basic Version without RAG, etc.)	0.4776	0.0087	0.0410	0.6061	0.0940	0.1280
Normal RAG	0.2442	0.1008	0.1119	0.3997	0.3354	0.3667
Agentic Retrieval	0.1267	0.5643	0.5431	0.1152	0.5922	0.6067
Graph-Based Discussion	0.2751	0.5310	0.5874	0.2792	0.6540	0.5311
Normal RAG & Graph-Based Discussion	0.0909	0.8951	0.8658	0.3416	0.8547	0.6770
Agentic Retrieval & Graph-Based Discussion	0.0051	0.9883	0.9761	0.2013	0.9474	0.8912
Agentic Retrieval & Round-Robin Discussion	0.0123	0.9456	0.9234	0.1567	0.9123	0.8345
Agentic Retrieval & Moderator-centered Discussion	0.0156	0.9123	0.8876	0.1789	0.8967	0.8123
<i>Drug Perturbation (Srivatsan Dataset [106])</i>						
CELLFORGE (Basic Version without RAG, etc.)	0.5760	0.0298	0.0475	0.6409	0.0992	0.1039
Normal RAG	0.2572	0.1584	0.1038	0.3022	0.3472	0.2901
Agentic Retrieval	0.1309	0.3437	0.4350	0.1210	0.3836	0.4169
Graph-Based Discussion	0.1670	0.4193	0.3764	0.1325	0.4266	0.3865
Normal RAG & Graph-Based Discussion	0.0995	0.6512	0.5933	0.0985	0.6784	0.7548
Agentic Retrieval & Graph-Based Discussion	0.0053	0.9881	0.9665	0.0080	0.9953	0.9802
Agentic Retrieval & Round-Robin Discussion	0.0145	0.9567	0.9234	0.0234	0.9456	0.9123
Agentic Retrieval & Moderator-centered Discussion	0.0189	0.9234	0.8967	0.0345	0.9123	0.8789

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

6 CONCLUSION

CELLFORGE is an autonomous multi-agent system that designs and implements model architectures for single-cell perturbation prediction without human intervention. By combining agentic retrieval with graph-based collaborative reasoning, it integrates computational, biological, and statistical expertise to adaptively improve across datasets and modalities. This work demonstrates that knowledge-grounded agentic frameworks can transcend manual or template-based design, yielding architectures that are both computationally effective and biologically meaningful. More broadly, our results highlight agent-based systems as a paradigm for automating scientific model development, enabling scalable exploration of modeling strategies in complex domains such as single-cell biology. Future directions include extending to new modalities, improving robustness, and generalizing to other areas of computational biology.

ETHICS STATEMENT

Our work uses only publicly available single-cell perturbation datasets [2, 67, 82, 84, 85, 106] under their respective licenses. No personally identifiable or sensitive data is involved. We emphasize that while CELLFORGE automates model design, any downstream use in biomedical applications should be carefully reviewed for ethical compliance, especially regarding clinical translation and potential misuse.

REPRODUCIBILITY STATEMENT

We detail datasets, evaluation metrics, and baseline implementations in the main text and Appendix.

USAGE OF LANGUAGE MODELS

We utilized a large language model (LLM) to aid in the preparation of this manuscript. Its use was limited to editorial tasks, including proofreading for typographical errors, correcting grammar, and improving the clarity and readability of the text.

REFERENCES

- 540
541
542 [1] Taghrid Abdelaal, Lennart Michielsen, Dries Cats, Daan Hoogduin, Hailiang Mei, Marcel J T
543 Reinders, and Ahmed Mahfouz. A comparison of automatic cell identification methods for
544 single-cell RNA sequencing data. *Genome Biology*, 20(1):1–17, 2019.
- 545
546 [2] Britt Adamson, Thomas M Norman, Marco Jost, Min Y Cho, James K Nuñez, Yuwen Chen,
547 Jacqueline E Villalta, Luke A Gilbert, Max A Horlbeck, Marco Y Hein, et al. A multiplexed
548 single-cell CRISPR screening platform enables systematic dissection of the unfolded protein
549 response. *Cell*, 167(7):1867–1882, 2016.
- 550
551 [3] Abhinav K. Adduri, Dhruv Gautam, Beatrice Bevilacqua, Alishba Imran, Rohan Shah, Mohsen
552 Naghypourfar, Noam Teyssier, Rajesh Ilango, Sanjay Nagaraj, Mingze Dong, Chiara Ricci-
553 Tam, Christopher Carpenter, Vishvak Subramanyam, Aidan Winters, Sravya Tirukkovular,
554 Jeremy Sullivan, Brian S. Plosky, Basak Eraslan, Nicholas D. Youngblut, Jure Leskovec,
555 Luke A. Gilbert, Silvana Konermann, Patrick D. Hsu, Alexander Dobin, Dave P. Burke,
556 Hani Goodarzi, and Yusuf H. Roohani. Predicting cellular responses to perturbation across
557 diverse contexts with state. *bioRxiv*, 2025. doi: 10.1101/2025.06.26.661135. URL <https://www.biorxiv.org/content/10.1101/2025.06.26.661135v2>.
- 558
559 [4] Constantin Ahlmann-Eltze, Wolfgang Huber, and Simon Anders. Deep learning-based predic-
560 tions of gene perturbation effects do not yet outperform simple linear baselines. *BioRxiv*, pp.
561 2024–09, 2024.
- 562
563 [5] Tal Ashuach, Mariano I. Gabitto, Rohan V. Koodli, Giuseppe Antonio Saldi, Michael I. Jordan,
564 and Nir Yosef. Multivi: deep generative model for the integration of multimodal data. *Nature
565 Methods*, 20(8):1222–1231, Aug 2023. doi: 10.1038/s41592-023-01909-9.
- 566
567 [6] Jinheon Baek, Sujay Kumar Jauhar, Silviu Cucerzan, and Sung Ju Hwang. ResearchAgent:
568 Iterative research idea generation over scientific literature with large language models. *arXiv
569 preprint arXiv:2404.07738*, 2024.
- 570
571 [7] Ding Bai, Caleb N Ellington, Shentong Mo, Le Song, and Eric P Xing. AttentionPert:
572 accurately modeling multiplexed genetic perturbations with multi-scale effects. *Bioinformatics*,
573 40(Supplement_1):i453–i461, 2024.
- 574
575 [8] Joeran Beel, Min-Yen Kan, and Moritz Baumgart. Evaluating sakana’s AI scientist for
576 autonomous research: Wishful thinking or an emerging reality towards’ artificial research
577 intelligence’(ARI)? *arXiv preprint arXiv:2502.14297*, 2025.
- 578
579 [9] Ihab Bendidi, Shawn Whitfield, Kian Kenyon-Dean, Hanene Ben Yedder, Yassir El Mesbahi,
580 Emmanuel Noutahi, and Alisandra K. Denton. Benchmarking transcriptomics foundation
581 models for perturbation analysis: one PCA still rules them all, 11 2024. URL [http://
582 arxiv.org/abs/2410.13956](http://arxiv.org/abs/2410.13956).
- 583
584 [10] Daniil A. Boiko, Robert MacKnight, and Gabe Gomes. Autonomous chemical research with
585 large language models. *Nature*, 623:760–768, 2023. doi: 10.1038/s41586-023-06792-0.
- 586
587 [11] Charlotte Bunne, Andreas Krause, and Marco Cuturi. Supervised training of conditional
588 monge maps. *Advances in Neural Information Processing Systems*, 35:6859–6872, 2022.
- 589
590 [12] Charlotte Bunne, Stefan G Stark, Gabriele Gut, Jacobo Sarabia Del Castillo, Mitch Levesque,
591 Kjong-Van Lehmann, Lucas Pelkmans, Andreas Krause, and Gunnar Rätsch. Learning single-
592 cell perturbation responses using neural optimal transport. *Nature Methods*, 20(11):1759–1768,
593 2023.
- [13] Charlotte Bunne, Yusuf Roohani, Yanay Rosen, Ankit Gupta, Xikun Zhang, Marcel Roed,
Theo Alexandrov, Mohammed AlQuraishi, Patricia Brennan, Daniel B Burkhardt, et al. How
to build the virtual cell with artificial intelligence: Priorities and opportunities. *Cell*, 187(25):
7045–7063, 2024.

- 594 [14] Daniel B. Burkhardt, Jay S. Stanley, Alexander Tong, Ana Luisa Perdigoto, Scott A. Gigante,
595 Kevan C. Herold, Guy Wolf, Antonio J. Giraldez, David van Dijk, and Smita Krishnaswamy.
596 Quantifying the effect of experimental perturbations at single-cell resolution. *Nature Biotech-*
597 *nology*, 39(5):619–629, May 2021. ISSN 1546-1696. doi: 10.1038/s41587-020-00803-5.
598 URL <https://www.nature.com/articles/s41587-020-00803-5>.
- 599 [15] Zhi-Jie Cao and Ge Gao. Multi-omics single-cell data integration and regulatory inference
600 with graph-linked embedding. *Nature Biotechnology*, 40(10):1458–1466, 2022.
- 601 [16] ZhiJie Cao and Ge Gao. Multiomics singlecell data integration and regulatory inference
602 with graphlinked embedding. *Nature Biotechnology*, 40(10):1458–1466, May 2022. doi:
603 10.1038/s41587-022-01284-4.
- 604 [17] Bo Chen, Xingyi Cheng, Pan Li, Yangli-ao Geng, Jing Gong, Shen Li, Zhilei Bei, Xu Tan,
605 Boyan Wang, Xin Zeng, et al. xtrimopglm: unified 100b-scale pre-trained transformer for
606 deciphering the language of protein. *arXiv preprint arXiv:2401.06199*, 2024.
- 607 [18] Tingting Chen, Srinivas Anumasa, Beibei Lin, Vedant Shah, Anirudh Goyal, and Dianbo
608 Liu. Auto-Bench: An automated benchmark for scientific discovery in LLMs. *arXiv preprint*
609 *arXiv:2502.15224*, 2025.
- 610 [19] Yiqun Chen and James Zou. Simple and effective embedding model for single-cell biology
611 built from ChatGPT. *Nature Biomedical Engineering*, 9(4):483–493, April 2025. ISSN
612 2157-846X. doi: 10.1038/s41551-024-01284-6. URL [https://www.nature.com/](https://www.nature.com/articles/s41551-024-01284-6)
613 [articles/s41551-024-01284-6](https://www.nature.com/articles/s41551-024-01284-6).
- 614 [20] Ziru Chen, Shijie Chen, Yuting Ning, Qianheng Zhang, Boshi Wang, Botao Yu, Yifei Li,
615 Zeyi Liao, Chen Wei, Zitong Lu, et al. ScienceAgentBench: Toward rigorous assessment of
616 language agents for data-driven scientific discovery. *arXiv preprint arXiv:2410.05080*, 2024.
- 617 [21] Haotian Cui, Chloe Wang, Hassaan Maan, Kuan Pang, Fengning Luo, Nan Duan, and Bo Wang.
618 scGPT: toward building a foundation model for single-cell multi-omics using generative ai.
619 *Nature Methods*, 21(8):1470–1480, 08 2024. ISSN 1548-7091.
- 620 [22] Atray Dixit, Oren Parnas, Biyu Li, Jenny Chen, Charles P Fulco, Livnat Jerby-Arnon, Ne-
621 manja D Marjanovic, Danielle Dionne, Tyler Burks, Raktima Raychowdhury, et al. Perturb-Seq:
622 dissecting molecular circuits with scalable single-cell RNA profiling of pooled genetic screens.
623 *Cell*, 167(7):1853–1866, 2016.
- 624 [23] Mingze Dong, Bao Wang, Jessica Wei, Antonio H de O. Fonseca, Curtis J Perry, Alexander
625 Frey, Ferial Ouerghi, Ellen F Foxman, Jeffrey J Ishizuka, Rahul M Dhodapkar, et al. Causal
626 identification of single-cell experimental perturbation effects with CINEMA-OT. *Nature*
627 *methods*, 20(11):1769–1779, 2023.
- 628 [24] Steffen Eger, Yong Cao, Jennifer D’Souza, Andreas Geiger, Christian Greisinger, Stephanie
629 Gross, Yufang Hou, Brigitte Krenn, Anne Lauscher, Yizhi Li, et al. Transforming science with
630 large language models: A survey on AI-assisted scientific discovery, experimentation, content
631 generation, and evaluation. *arXiv preprint arXiv:2502.05151*, 2025.
- 632 [25] Adibvafa Fallahpour, Andrew Magnuson, Purav Gupta, Shihao Ma, Jack Naimer, Arnav Shah,
633 Haonan Duan, Omar Ibrahim, Hani Goodarzi, Chris J Maddison, et al. Bioreason: Incentivizing
634 multimodal biological reasoning within a dna-llm model. *arXiv preprint arXiv:2505.23579*,
635 2025.
- 636 [26] Rochelle V Flores, Shicong Wang, et al. Deep learning tackles single-cell analysis—a survey
637 of deep learning for scRNA-seq analysis. *Briefings in Bioinformatics*, 23(5):bbac327, 2022.
- 638 [27] Robert Friel, Masha Belyi, and Atindriyo Sanyal. RAGBench: Explainable benchmark for
639 retrieval-augmented generation systems, 2025. URL [http://arxiv.org/abs/2407.](http://arxiv.org/abs/2407.11005)
640 [11005](http://arxiv.org/abs/2407.11005).
- 641 [28] Xi Fu, Shentong Mo, Alejandro Buendia, Anouchka P Laurent, Anqi Shao, Maria del Mar
642 Alvarez-Torres, Tianji Yu, Jimin Tan, Jiayu Su, Romella Sagatelian, et al. A foundation model
643 of transcription across human cell types. *Nature*, 637(8047):965–973, 2025.

- 648 [29] Xian Gao, Zongyun Zhang, Mingye Xie, Ting Liu, and Yuzhuo Fu. Graph of AI ideas:
649 Leveraging knowledge graphs and llms for AI research idea generation. *arXiv preprint*
650 *arXiv:2503.08549*, 2025.
- 651 [30] Yicheng Gao, Zhiting Wei, Kejing Dong, Ke Chen, Jingya Yang, Guohui Chuai, and Qi Liu.
652 Toward subtask decomposition-based learning and benchmarking for predicting genetic pertur-
653 bation outcomes and beyond. *Nature Computational Science*, 4(10):773–785, Sep 2024. doi:
654 10.1038/s43588-024-00698-1.
- 655 [31] Aniketh Garikaparathi, Manasi Patwardhan, Lovekesh Vig, and Arman Cohan. IRIS: Interactive
656 research ideation system for accelerating scientific discovery. *arXiv preprint arXiv:2504.16728*,
657 2025.
- 658 [32] George I. Gavriilidis, Vasileios Vasileiou, Aspasia Orfanou, Naveed Ishaque, and Fotis Pso-
659 mopoulos. A mini-review on perturbation modelling across single-cell omic modalities.
660 *Computational and Structural Biotechnology Journal*, 23:1886–1896, December 2024. ISSN
661 2001-0370. doi: 10.1016/j.csbj.2024.04.058.
- 662 [33] Adam Gayoso, Zo Steier, Romain Lopez, Jeffrey Regier, Kristopher L. Nazor, Aaron Streets,
663 and Nir Yosef. Joint probabilistic modeling of single-cell multi-omic data with totalvi. *Nature*
664 *Methods*, 18(3):272–282, Mar 2021. doi: 10.1038/s41592-020-01050-x.
- 665 [34] Alireza Ghafarollahi and Markus J Buehler. AtomAgents: Alloy design and discovery through
666 physics-aware multi-modal multi-agent artificial intelligence. *arXiv preprint arXiv:2407.10022*,
667 2024.
- 668 [35] Alireza Ghafarollahi and Markus J Buehler. Sparks: Multi-agent artificial intelligence model
669 discovers protein design principles. *arXiv preprint arXiv:2504.19017*, 2025.
- 670 [36] Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li,
671 Yinghan Shen, Shengjie Ma, Honghao Liu, et al. A survey on llm-as-a-judge. *arXiv preprint*
672 *arXiv:2411.15594*, 2024.
- 673 [37] Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V. Chawla, Olaf
674 Wiest, and Xiangliang Zhang. Large language model based multi-agents: A survey of progress
675 and challenges, 02 2024. URL <https://arxiv.org/abs/2402.01680>.
- 676 [38] Minsheng Hao, Jing Gong, Xin Zeng, Chiming Liu, Yucheng Guo, Xingyi Cheng, Taifeng
677 Wang, Jianzhu Ma, Xuegong Zhang, and Le Song. Large-scale foundation model on single-cell
678 transcriptomics. *Nature methods*, 21(8):1481–1491, 2024.
- 679 [39] Minsheng Hao, Yongju Lee, Hanchen Wang, Gabriele Scalia, and Aviv Regev. Perturboagent:
680 A self-planning agent for boosting sequential perturb-seq experiments. *bioRxiv*, pp. 2025–05,
681 2025.
- 682 [40] Yuhan Hao, Stephanie Hao, Erica Andersen-Nissen, William M. Mauck, Shiwei Zheng,
683 Andrew Butler, Maddie J. Lee, Aaron J. Wilk, Charlotte Darby, Michael Zager, Paul Hoffman,
684 Marlon Stoeckius, Efthymia Papalexli, Eleni P. Mimitou, Jaison Jain, Avi Srivastava, Tim Stuart,
685 Lamar M. Fleming, Bertrand Yeung, Angela J. Rogers, Juliana M. McElrath, Catherine A.
686 Blish, Raphael Gottardo, Peter Smibert, and Rahul Satija. Integrated analysis of multimodal
687 single-cell data. *Cell*, 184(13):3573–3587.e29, June 2021. ISSN 0092-8674, 1097-4172.
688 doi: 10.1016/j.cell.2021.04.048. URL [https://www.cell.com/cell/abstract/
689 S0092-8674\(21\)00583-3](https://www.cell.com/cell/abstract/S0092-8674(21)00583-3).
- 690 [41] Leon Hetzel, Simon Boehm, Niki Kilbertus, Stephan Gunnemann, Fabian Theis, et al. Pre-
691 dicting cellular responses to novel drug perturbations at a single-cell resolution. *Advances in*
692 *Neural Information Processing Systems*, 35:26711–26722, 2022.
- 693 [42] Lukas Heumos, Anne C Schaar, et al. Best practices for single-cell analysis across modalities.
694 *Nature Reviews Genetics*, 24(6):395–415, 2023.
- 695 [43] Chao-Chun Hsu, Erin Bransom, Jenna Sparks, Bailey Kuehl, Chenhao Tan, David Wadden,
696 Lucy Lu Wang, and Aakanksha Naik. CHIME: LLM-assisted hierarchical organization of
697 scientific studies for literature review support. *Findings of ACL 2024*, 2024.

- 702 [44] Xiang Hu, Hongyu Fu, Jinge Wang, Yifeng Wang, Zhikun Li, Renjun Xu, Yu Lu, Yaochu Jin,
703 Lili Pan, and Zhenzhong Lan. Nova: An iterative planning and search approach to enhance
704 novelty and diversity of llm generated ideas. *arXiv preprint arXiv:2410.14255*, 2024.
705
- 706 [45] Kexin Huang, Ying Jin, Ryan Li, Michael Y Li, Emmanuel Candès, and Jure Leskovec.
707 Automated hypothesis validation with agentic sequential falsifications. *arXiv preprint*
708 *arXiv:2502.09858*, 2025.
- 709 [46] Kexin Huang, Serena Zhang, Hanchen Wang, Yuanhao Qu, Yingzhou Lu, Yusuf Roohani,
710 Ryan Li, Lin Qiu, Junze Zhang, Yin Di, et al. Biomni: A general-purpose biomedical ai agent.
711 *bioRxiv*, pp. 2025–05, 2025.
712
- 713 [47] Qian Huang, Jian Vora, Percy Liang, and Jure Leskovec. MAgentBench: Evaluating language
714 agents on machine learning experimentation. In *ICML 2024*, 2024.
- 715 [48] Ana-Maria Istrate, Donghui Li, and Theofanis Karaletsos. scGenePT: Is language all you need
716 for modeling single-cell perturbations?, October 2024. URL <https://www.biorxiv.org/content/10.1101/2024.10.23.619972v1>.
717
- 718 [49] Yanrong Ji, Zhihan Zhou, Han Liu, and Ramana V. Davuluri. Dnabert: pretrained bidirectional
719 encoder representations from transformers model for dnalanguage in genome. *Bioinformatics*,
720 37(15):2112–2120, August 2021. doi: 10.1093/bioinformatics/btab083. URL <https://doi.org/10.1093/bioinformatics/btab083>.
721
- 722 [50] Yuge Ji, Mohammad Lotfollahi, F. Alexander Wolf, and Fabian J. Theis. Machine learning
723 for perturbational single-cell omics. *Cell Systems*, 12(6):522–537, June 2021. ISSN 2405-
724 4712, 2405-4720. doi: 10.1016/j.cels.2021.05.016. URL [https://www.cell.com/cell-systems/abstract/S2405-4712\(21\)00202-7](https://www.cell.com/cell-systems/abstract/S2405-4712(21)00202-7).
725
- 726 [51] Qun Jiang, Shengquan Chen, Xiaoyang Chen, and Rui Jiang. scPRAM accurately predicts
727 single-cell gene expression perturbation response based on attention mechanism. *Bioinformat-*
728 *ics*, 40(5):btae265, 2024.
729
- 730 [52] Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William W Cohen, and Xinghua Lu. PubMedQA:
731 A dataset for biomedical research question answering. *arXiv preprint arXiv:1909.06146*, 2019.
732
- 733 [53] Qiao Jin, Yifan Yang, Qingyu Chen, and Zhiyong Lu. GeneGPT: Augmenting large language
734 models with domain tools for improved access to biomedical information. *Bioinformatics*, 40
735 (2):btae075, February 2024. ISSN 1367-4811. doi: 10.1093/bioinformatics/btae075. URL
736 <https://doi.org/10.1093/bioinformatics/btae075>.
737
- 738 [54] Ruofan Jin, Zaixi Zhang, Mengdi Wang, and Le Cong. Stella: Self-evolving llm agent for
739 biomedical research. *arXiv preprint arXiv:2507.02004*, 2025.
740
- 741 [55] Liqiang Jing, Zhehui Huang, Xiaoyang Wang, Wenlin Yao, Wenhao Yu, Kaixin Ma, Hongming
742 Zhang, Xinya Du, and Dong Yu. DSBench: How far are data science agents to becoming data
743 science experts? *arXiv preprint arXiv:2409.07703*, 2024.
- 744 [56] Julia Joung, Silvana Konermann, Jonathan S. Gootenberg, Omar O. Abudayyeh, Randall J.
745 Platt, Mark D. Brigham, Neville E. Sanjana, and Feng Zhang. Genome-scale CRISPR-Cas9
746 knockout and transcriptional activation screening. *Nature Protocols*, 12(4):828–863, April
747 2017. ISSN 1750-2799. doi: 10.1038/nprot.2017.016. URL <https://www.nature.com/articles/nprot.2017.016>.
748
- 749 [57] Kenji Kamimoto, Blerta Stringa, Christy M Hoffmann, Kunal Jindal, Lilianna Solnica-Krezel,
750 and Samantha A Morris. Dissecting cell identity via network inference and in silico gene
751 perturbation. *Nature*, 614(7949):742–751, 2023.
752
- 753 [58] Kasia Z. Kedzierska, Lorin Crawford, Ava P. Amini, and Alex X. Lu. Zero-shot evaluation
754 reveals limitations of single-cell foundation models. *Genome Biology*, 26(1):101, April 2025.
755 ISSN 1474-760X. doi: 10.1186/s13059-025-03574-x. URL <https://doi.org/10.1186/s13059-025-03574-x>.

- 756 [59] Dominik Klein, Jonas Simon Fleck, Daniil Bobrovskiy, Lea Zimmermann, Sören Becker,
757 Alessandro Palma, Leander Dony, Alejandro Tejada-Lapuerta, Guillaume Huguet, Hsiu-Chuan
758 Lin, et al. Cellflow enables generative single-cell phenotype modeling with flow matching.
759 *bioRxiv*, pp. 2025–04, 2025.
- 760 [60] Patrick Tser Jern Kon, Jiachen Liu, Qiuyi Ding, Yiming Qiu, Zhenning Yang, Yibo Huang,
761 Jayanth Srinivasa, Myungjin Lee, Mosharaf Chowdhury, and Ang Chen. Curie: Toward rigor-
762 ous and automated scientific experimentation with AI agents. *arXiv preprint arXiv:2502.16069*,
763 2025.
- 764 [61] Adithya Kulkarni, Fatimah Alotaibi, Xinyue Zeng, Longfeng Wu, Tong Zeng, Barry Menglong
765 Yao, Minqian Liu, Shuaicheng Zhang, Lifu Huang, and Dawei Zhou. Scientific hypoth-
766 esis generation and validation: Methods, datasets, and future directions. *arXiv preprint*
767 *arXiv:2505.04651*, 2025.
- 768 [62] Daniel Levine, Syed Asad Rizvi, Sacha Lévy, Nazreen Pallikkavaliyaveetil, David Zhang,
769 Xingyu Chen, Sina Ghadermarzi, Ruiming Wu, Ziheng Zheng, Ivan Vrkcic, et al. Cell2Sentence:
770 teaching large language models the language of biology. *BioRxiv*, pp. 2023–09, 2024.
- 771 [63] Chen Li, Haoxiang Gao, Yuli She, Haiyang Bian, Qing Chen, Kai Liu, Lei Wei, and Xuegong
772 Zhang. Benchmarking ai models for in silico gene perturbation of cells. *bioRxiv*, pp. 2024–12,
773 2024.
- 774 [64] Haitao Li, Qian Dong, Junjie Chen, Huixue Su, Yujia Zhou, Qingyao Ai, Ziyi Ye, and Yiqun
775 Liu. Lms-as-judges: a comprehensive survey on llm-based evaluation methods. *arXiv preprint*
776 *arXiv:2412.05579*, 2024.
- 777 [65] Lanxiang Li, Yue You, Wenyu Liao, Xueying Fan, Shihong Lu, Ye Cao, Bo Li, Wenle Ren,
778 Yunlin Fu, Jiaming Kong, et al. A systematic comparison of single-cell perturbation response
779 prediction models. *bioRxiv*, pp. 2024–12, 2024.
- 780 [66] Ruochen Li, Teerth Patel, Qingyun Wang, and Xinya Du. MLR-Copilot: Autonomous machine
781 learning research based on large language models agents. *arXiv preprint arXiv:2408.14033*,
782 2024.
- 783 [67] Noa Liscovitch-Brauer, Antonino Montalbano, Jiale Deng, Alejandro Méndez-Mancilla, Hans-
784 Hermann Wessels, Nicholas G Moss, Chia-Yu Kung, Akash Sookdeo, Xinyi Guo, Evan Geller,
785 et al. Profiling the genetic determinants of chromatin accessibility with scalable single-cell
786 crispr screens. *Nature biotechnology*, 39(10):1270–1277, 2021.
- 787 [68] Haokun Liu, Yangqiaoyu Zhou, Mingxuan Li, Chenfei Yuan, and Chenhao Tan. Literature
788 meets data: A synergistic approach to hypothesis generation. *arXiv preprint arXiv:2410.17309*,
789 2024.
- 790 [69] Tianyu Liu, Yuge Wang, Rex Ying, and Hongyu Zhao. MuSe-GNN: Learning Unified Gene
791 Representation From Multimodal Biological Graph Data, September 2023. URL <http://arxiv.org/abs/2310.02275>.
- 792 [70] Zijun Liu, Kaiming Liu, Yiqi Zhu, Xuanyu Lei, Zonghan Yang, Zhenhe Zhang, Peng Li, and
793 Yang Liu. AIGS: Generating science from AI-powered automated falsification. *arXiv preprint*
794 *arXiv:2411.11910*, 2024.
- 795 [71] Romain Lopez, Jordan Regier, Michael B Cole, Michael I Jordan, and Nir Yosef. scVI: deep
796 generative modeling for single-cell transcriptomics. *Nature Methods*, 15(12):1053–1058, 2018.
- 797 [72] Mohammad Lotfollahi, F Alexander Wolf, and Fabian J Theis. scGen predicts single-cell
798 perturbation responses. *Nature methods*, 16(8):715–721, 2019.
- 799 [73] Mohammad Lotfollahi, Anna Klimovskaia Susmelj, Carlo De Donno, Yuge Ji, Ignacio L.
800 Ibarra, et al. Learning interpretable cellular responses to complex perturbations in high-
801 throughput screens. *Bioinformatics*, 04 2021.

- 810 [74] Mohammad Lotfollahi, Anna Klimovskaia Susmelj, Carlo De Donno, Leon Hetzel, Yuge Ji,
811 Ignacio L Ibarra, Sanjay R Srivatsan, Mohsen Naghipourfar, Riza M Daza, Beth Martin, et al.
812 Predicting cellular responses to complex perturbations in high-throughput screens. *Molecular*
813 *systems biology*, 19(6):e11517, 2023.
- 814 [75] Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. The AI
815 scientist: Towards fully automated open-ended scientific discovery, 09 2024. URL <http://arxiv.org/abs/2408.06292>.
816
- 817 [76] Malte D Luecken, Daniel B Burkhardt, Fabian J Theis, et al. Defining and benchmarking open
818 problems in single-cell analysis. *Nature Methods*, 19(4):412–420, 2022.
- 819 [77] Malte D. Luecken, M. Büttner, K. Chaichoompu, A. Danese, M. Interlandi, M. F. Mueller,
820 D. C. Strobl, L. Zappia, M. Dugas, M. Colomé-Tatché, and Fabian J. Theis. Benchmarking
821 atlas-level data integration in single-cell genomics. *Nature Methods*, 19(1):41–50, January
822 2022. ISSN 1548-7091, 1548-7105. doi: 10.1038/s41592-021-01336-8.
- 823 [78] Hiba Maan, Matti Lähde, et al. Characterizing the impacts of dataset imbalance on single-cell
824 data integration. *Nature Biotechnology*, 42(1):56–60, 2024.
- 825 [79] Bodhisattwa Prasad Majumder, Harshit Surana, Dhruv Agarwal, Bhavana Dalvi Mishra,
826 Abhijeetsingh Meena, Aryan Prakhar, Tirth Vora, Tushar Khot, Ashish Sabharwal, and Peter
827 Clark. DiscoveryBench: Towards data-driven discovery with large language models. *arXiv*
828 *preprint arXiv:2407.01725*, 2024.
- 829 [80] Pablo Monfort-Lanzas, Katja Rungger, Leonie Madersbacher, and Hubert Hackl. Ma-
830 chine learning to dissect perturbations in complex cellular systems. *Computational and*
831 *Structural Biotechnology Journal*, 27:832–842, January 2025. ISSN 2001-0370. doi:
832 10.1016/j.csbj.2025.02.028. URL [https://www.sciencedirect.com/science/](https://www.sciencedirect.com/science/article/pii/S2001037025000583)
833 [article/pii/S2001037025000583](https://www.sciencedirect.com/science/article/pii/S2001037025000583).
834
- 835 [81] Vladimir Naumov, Diana Zagirova, Sha Lin, Yupeng Xie, Wenhao Gou, Anatoly Urban,
836 Nina Tikhonova, Khadija Alawi, Mike Durymanov, Fedor Galkin, et al. DORA AI scientist:
837 Multi-agent virtual research team for scientific exploration discovery and automated report
838 generation. *bioRxiv*, 2025.
- 839 [82] Thomas M Norman, Max A Horlbeck, Joseph M Replogle, Alex Y Ge, Albert Xu, Marco
840 Jost, Luke A Gilbert, and Jonathan S Weissman. Exploring genetic interaction manifolds
841 constructed from rich single-cell phenotypes. *Science*, 365(6455):786–793, 2019.
- 842 [83] OpenAI. Introducing deep research. [https://openai.com/index/](https://openai.com/index/deep-research/)
843 [deep-research/](https://openai.com/index/deep-research/), 2025. Accessed: 2025-05-08.
- 844 [84] Efthymia Papalexi, Eleni P Mimitou, Andrew W Butler, Samantha Foster, Bernadette Bracken,
845 William M Mauck III, Hans-Hermann Wessels, Yuhao Hao, Bertrand Z Yeung, Peter Smibert,
846 et al. Characterizing the molecular regulation of inhibitory immune checkpoints with
847 multimodal single-cell screens. *Nature genetics*, 53(3):322–331, 2021.
- 848 [85] Stefan Peidli, Tessa D Green, Ciyue Shen, Torsten Gross, Joseph Min, Samuele Garda,
849 Bo Yuan, Linus J Schumacher, Jake P Taylor-King, Debora S Marks, et al. scPerturb:
850 harmonized single-cell perturbation data. *Nature Methods*, 21(3):531–540, 2024.
- 851 [86] Zoe Piran, Niv Cohen, Yedid Hoshen, and Mor Nitzan. Disentanglement of single-cell data
852 with biolord. *Nature Biotechnology*, 42(11):1678–1683, 2024.
- 853 [87] Kevin Pu, KJ Feng, Tovi Grossman, Tom Hope, Bhavana Dalvi Mishra, Matt Latzke, Jonathan
854 Bragg, Joseph Chee Chang, and Pao Siangliulue. IdeaSynth: Iterative research idea develop-
855 ment through evolving and composing idea facets with literature-grounded feedback. *arXiv*
856 *preprint arXiv:2410.04025*, 2024.
- 857 [88] Biqing Qi, Kaiyan Zhang, Haoxiang Li, Kai Tian, Sihang Zeng, Zhang-Ren Chen, Jin-Fang
858 Hu, and Bowen Zhou. Large language models are zero shot hypothesis proposers. *Instruction*
859 *Workshop @ NeurIPS 2023*, 2023.

- 864 [89] Xiaoning Qi, Lianhe Zhao, Chenyu Tian, Yueyue Li, Zhen-Lin Chen, Peipei Huo, Run-
865 sheng Chen, Xiaodong Liu, Baoping Wan, Shengyong Yang, and Yi Zhao. Predicting
866 transcriptional responses to novel chemical perturbations using deep generative model for
867 drug discovery. *Nature Communications*, 15(1):9256, October 2024. ISSN 2041-1723.
868 doi: 10.1038/s41467-024-53457-1. URL [https://www.nature.com/articles/
869 s41467-024-53457-1](https://www.nature.com/articles/s41467-024-53457-1).
- 870 [90] Xiaojie Qiu, Yan Zhang, Jorge D Martin-Rufino, Chen Weng, Shayan Hosseinzadeh, Dian
871 Yang, Angela N Pogson, Marco Y Hein, Kyung Hoi Joseph Min, Li Wang, et al. Mapping
872 transcriptomic vector fields of single cells. *Cell*, 185(4):690–711, 2022.
- 873 [91] Yansheng Qiu, Haoquan Zhang, Zhaopan Xu, Ming Li, Diping Song, Zheng Wang, and
874 Kaipeng Zhang. AI Idea Bench 2025: AI research idea generation benchmark. *arXiv preprint
875 arXiv:2504.14191*, 2025.
- 876 [92] Marissa Radensky, Simra Shahid, Raymond Fok, Pao Siangliulue, Tom Hope, and Daniel S
877 Weld. Scideator: Human-llm scientific idea generation grounded in research-paper facet
878 recombination. *arXiv preprint arXiv:2409.14634*, 2024.
- 879 [93] Chandan K Reddy and Parshin Shojaee. Towards scientific discovery with generative AI:
880 Progress, opportunities, and challenges. In *Proceedings of the AAAI Conference on Artificial
881 Intelligence*, volume 39, pp. 28601–28609, 2025.
- 882 [94] Shuo Ren, Pu Jian, Zhenjiang Ren, Chunlin Leng, Can Xie, and Jiajun Zhang. Towards sci-
883 entific intelligence: A survey of llm-based scientific agents. *arXiv preprint arXiv:2503.24047*,
884 2025.
- 885 [95] Joseph M. Replogle, Thomas M. Norman, Albert Xu, Jeffrey A. Hussmann, Jin Chen,
886 J. Zachery Cogan, Elliott J. Meer, Jessica M. Terry, Daniel P. Riordan, Niranjan Srinivas,
887 Ian T. Fiddes, Joseph G. Arthur, Luigi J. Alvarado, Katherine A. Pfeiffer, Tarjei S.
888 Mikkelsen, Jonathan S. Weissman, and Britt Adamson. Combinatorial single-cell CRISPR
889 screens by direct guide RNA capture and targeted sequencing. *Nature Biotechnology*, 38
890 (8):954–961, August 2020. ISSN 1546-1696. doi: 10.1038/s41587-020-0470-y. URL
891 <https://www.nature.com/articles/s41587-020-0470-y>.
- 892 [96] Yusuf Roohani, Kexin Huang, and Jure Leskovec. Predicting transcriptional outcomes of novel
893 multigene perturbations with GEARS. *Nature Biotechnology*, 42(6):927–935, 2024.
- 894 [97] Yusuf H Roohani, Jian Vora, Qian Huang, Percy Liang, and Jure Leskovec. BioDiscoveryAgent:
895 An ai agent for designing genetic perturbation experiments. In *ICLR 2024 Workshop on
896 Machine Learning for Genomics Explorations*, 2024.
- 897 [98] Yusuf H. Roohani, Tony J. Hua, Po-Yuan Tung, Lexi R. Bounds, Feiqiao B. Yu, Alexander
898 Dobin, Noam Teyssier, Abhinav Adduri, Alden Woodrow, Brian S. Plosky, Reshma Mehta,
899 Benjamin Hsu, Jeremy Sullivan, Chiara Ricci-Tam, Nianzhen Li, Julia Kazaks, Luke A. Gilbert,
900 Silvana Konermann, Patrick D. Hsu, Hani Goodarzi, and Dave P. Burke. Virtual cell challenge:
901 Toward a turing test for the virtual cell. *Cell*, 188(13):3370–3374, 2025.
- 902 [99] Kai Ruan, Xuan Wang, Jixiang Hong, Peng Wang, Yang Liu, and Hao Sun. LiveIdeaBench:
903 Evaluating llms’ scientific creativity and idea generation with minimal context. *arXiv preprint
904 arXiv:2412.17596*, 2024.
- 905 [100] Geoffrey Schiebinger, Jian Shu, Marcin Tabaka, Brian Cleary, Vidya Subramanian, Aryeh
906 Solomon, Joshua Gould, Siyan Liu, Stacie Lin, Peter Berube, et al. Optimal-transport analysis
907 of single-cell gene expression identifies developmental trajectories in reprogramming. *Cell*,
908 176(4):928–943, 2019.
- 909 [101] Samuel Schmidgall and Michael Moor. AgentRxiv: Towards collaborative autonomous
910 research. *arXiv preprint arXiv:2503.18102*, 2025.
- 911 [102] Chenglei Si, Diyi Yang, and Tatsunori Hashimoto. Can LLMs generate novel research ideas? a
912 large-scale human study with 100+ NLP researchers. *arXiv preprint arXiv:2409.04109*, 2024.
- 913
- 914
- 915
- 916
- 917

- 918 [103] Michael D. Skarlinski, Sam Cox, Jon M. Laurent, James D. Braza, Michaela Hinks,
919 Michael J. Hammerling, Manvitha Ponnampati, Samuel G. Rodrigues, and Andrew D. White.
920 Language agents achieve superhuman synthesis of scientific knowledge. *arXiv preprint*
921 *arXiv:2409.13740*, 2024.
- 922 [104] Michael A Skinnider, Jordan W Squair, Claudia Kathe, Mark A Anderson, Matthieu Gautier,
923 Kaya JE Matson, Marco Milano, Thomas H Hutson, Quentin Barraud, Aaron A Phillips, et al.
924 Cell type prioritization in single-cell data. *Nature biotechnology*, 39(1):30–34, 2021.
- 925 [105] Bicna Song, Dingyu Liu, Weiwei Dai, Natalie F. McMyn, Qingyang Wang, Dapeng Yang,
926 Adam Krejci, Anatoly Vasilyev, Nicole Untermoser, Anke Loregger, Dongyuan Song, Breanna
927 Williams, Bess Rosen, Xiaolong Cheng, Lumen Chao, Hanuman T. Kale, Hao Zhang, Yarui
928 Diao, Tilmann Bückstümmer, Janet D. Siliciano, Jingyi Jessica Li, Robert F. Siliciano, Danwei
929 Huangfu, and Wei Li. Decoding heterogeneous single-cell perturbation responses. *Nature Cell*
930 *Biology*, 27(3):493–504, March 2025. ISSN 1476-4679. doi: 10.1038/s41556-025-01626-9.
931 URL <https://www.nature.com/articles/s41556-025-01626-9>.
- 932 [106] Sanjay R Srivatsan, José L McFaline-Figueroa, Vijay Ramani, Lauren Saunders, Junyue Cao,
933 Jonathan Packer, Hannah A Pliner, Dana L Jackson, Riza M Daza, Lena Christiansen, et al.
934 Massively multiplex chemical transcriptomics at single-cell resolution. *Science*, 367(6473):
935 45–51, 2020.
- 936 [107] Tim Stuart, Andrew Butler, Paul Hoffman, Christoph Hafemeister, Efthymia Papalexi, William
937 M. III Mauck, Yuhan Hao, Marlon Stoeckius, Peter Smibert, and Rahul Satija. Comprehensive
938 integration of single-cell data. *Cell*, 177(7):1888–1902.e21, June 2019. doi: 10.1016/j.cell.
939 2019.05.031. URL <https://doi.org/10.1016/j.cell.2019.05.031>.
- 940 [108] Quan Tang, Na Le, et al. Single-cell multimodal prediction via transformers. In *NeurIPS 2022*
941 *Workshop on Learning from Time Series for Health*, 2022.
- 942 [109] Xiangru Tang, Anni Zou, Zhuosheng Zhang, Ziming Li, Yilun Zhao, Xingyao Zhang, Arman
943 Cohan, and Mark Gerstein. MedAgents: Large language models as collaborators for zero-shot
944 medical reasoning. *Findings of ACL 2024*, 2024.
- 945 [110] Christina V Theodoris, Ling Xiao, Anant Chopra, Mark D Chaffin, Zeina R Al Sayed,
946 Matthew C Hill, Helene Mantineo, Elizabeth M Brydon, Zexian Zeng, X Shirley Liu, et al.
947 Transfer learning enables predictions in network biology. *Nature*, 618(7965):616–624, 2023.
- 948 [111] Minyang Tian, Luyu Gao, Shizhuo Zhang, Xinan Chen, Cunwei Fan, Xuefei Guo, Roland
949 Haas, Pan Ji, Kittithat Krongchon, Yao Li, et al. SciCode: A research coding benchmark
950 curated by scientists. *Advances in Neural Information Processing Systems*, 37:30624–30650,
951 2024.
- 952 [112] Hanchen Wang, Yichun He, Paula P Coelho, Matthew Bucci, Abbas Nazir, Bob Chen, Linh
953 Trinh, Serena Zhang, Kexin Huang, Vineethkrishna Chandrasekar, et al. Spatialagent: An
954 autonomous ai agent for spatial biology. *bioRxiv*, pp. 2025–04, 2025.
- 955 [113] Juexin Wang, Anjun Ma, Yuzhou Chang, Jianting Gong, Yuexu Jiang, Ren Qi, Cankun Wang,
956 Hongjun Fu, Qin Ma, and Dong Xu. scGNN is a novel graph neural network framework
957 for single-cell RNA-Seq analyses. *Nature Communications*, 12(1):1882, March 2021. ISSN
958 2041-1723. doi: 10.1038/s41467-021-22197-x. URL [https://www.nature.com/](https://www.nature.com/articles/s41467-021-22197-x)
959 [articles/s41467-021-22197-x](https://www.nature.com/articles/s41467-021-22197-x).
- 960 [114] F. Alexander Wolf, Philipp Angerer, and Fabian J. Theis. SCANPY: Large-scale single-cell
961 gene expression data analysis. *Genome Biology*, 19(1):1–5, December 2018. ISSN 1474-760X.
962 doi: 10.1186/s13059-017-1382-0. URL [https://genomebiology.biomedcentral.](https://genomebiology.biomedcentral.com/articles/10.1186/s13059-017-1382-0)
963 [com/articles/10.1186/s13059-017-1382-0](https://genomebiology.biomedcentral.com/articles/10.1186/s13059-017-1382-0).
- 964 [115] Fan Yang, Fang Wang, Longkai Huang, Linjing Liu, Junzhou Huang, and Jianhua Yao. Reply
965 to: Deeper evaluation of a single-cell foundation model. *Nature Machine Intelligence*, 6(12):
966 1447–1450, December 2024. ISSN 2522-5839. doi: 10.1038/s42256-024-00948-x. URL
967 <https://www.nature.com/articles/s42256-024-00948-x>.
- 968
969
970
971

- 972 [116] Tong Yang, Xiaodan Hu, Xiaohan Li, Mingda Tan, Jingfeng Zhang, Zhilin Wen, Ernie Chang,
973 Andrew M. Dai, Quoc V. Li, Joseph E. Gonzalez, Claire Cardie, and Jason Wei. AnyBench:
974 Language models evaluate anything. *arXiv preprint arXiv:2312.13771*, 2023.
975
- 976 [117] Xiaodong Yang, Guole Liu, Guihai Feng, Dechao Bu, Pengfei Wang, Jie Jiang, Shubai Chen,
977 Qinmeng Yang, Hefan Miao, Yiyang Zhang, Zhenpeng Man, Zhongming Liang, Zichen
978 Wang, Yaning Li, Zheng Li, Yana Liu, Yao Tian, Wenhao Liu, Cong Li, Ao Li, Jingxi
979 Dong, Zhilong Hu, Chen Fang, Lina Cui, Zixu Deng, Haiping Jiang, Wentao Cui, Jiahao
980 Zhang, Zhaohui Yang, Handong Li, Xingjian He, Liqun Zhong, Jiaheng Zhou, Zijian Wang,
981 Qingqing Long, Ping Xu, Xin Li, Hongmei Wang, Zhen Meng, Xuezhi Wang, Yangang
982 Wang, Yong Wang, Shihua Zhang, Jingtao Guo, Yi Zhao, Yuanchun Zhou, Fei Li, Jing Liu,
983 Yiqiang Chen, Ge Yang, and Xin Li. GeneCompass: Deciphering universal gene regulatory
984 mechanisms with a knowledge-informed cross-species foundation model. *Cell Research*, 34
985 (12):830–845, December 2024. ISSN 1748-7838. doi: 10.1038/s41422-024-01034-y. URL
<https://www.nature.com/articles/s41422-024-01034-y>.
986
- 987 [118] Nicholas D. Youngblut, Christopher Carpenter, Jaanak Prashar, Chiara Ricci-Tam, Rajesh
988 Ilango, Noam Teyssier, Silvana Konermann, Patrick D. Hsu, Alexander Dobin, David P.
989 Burke, Hani Goodarzi, and Yusuf H. Roohani. scBaseCamp: An AI agent-curated, uniformly
990 processed, and continually expanding single cell data repository, March 2025. URL <https://www.biorxiv.org/content/10.1101/2025.02.27.640494v1>.
991
- 992 [119] Hengshi Yu, Weizhou Qian, Yuxuan Song, and Joshua D Welch. PerturbNet predicts single-
993 cell responses to unseen chemical and genetic perturbations. *Molecular Systems Biology*,
994 21(8):960–982, August 2025. ISSN 1744-4292. doi: 10.1038/s44320-025-00131-3. URL
995 <https://www.embopress.org/doi/full/10.1038/s44320-025-00131-3>.
- 996 [120] Bo Yuan, Ciyue Shen, Augustin Luna, Anil Korkut, Debora S. Marks, John Ingraham, and Chris
997 Sander. CellBox: Interpretable Machine Learning for Perturbation Biology with Application
998 to the Design of Cancer Combination Therapy. *Cell Systems*, 12(2):128–140.e4, February
999 2021. ISSN 2405-4712, 2405-4720. doi: 10.1016/j.cels.2020.11.013. URL [https://www.
1000 cell.com/cell-systems/abstract/S2405-4712\(20\)30464-6](https://www.cell.com/cell-systems/abstract/S2405-4712(20)30464-6).
- 1001 [121] Ruiqi Zhong, Peter Zhang, Steve Li, Jinwoo Ahn, Dan Klein, and Jacob Steinhardt. Goal
1002 driven discovery of distributional differences via language descriptions. In *NeurIPS 2023*,
1003 2023.
1004
- 1005 [122] Maxim Zvyagin, Alexander Brace, Kyle Hippe, Yuntian Deng, Bin Zhang, Cindy Orozco
1006 Bohorquez, Austin Clyde, Bharat Kale, Danilo Perez-Rivera, Heng Ma, Carla M. Mann,
1007 Michael Irvin, Defne G. Ozgulbas, Natalia Vassilieva, James Gregory Pauloski, Logan Ward,
1008 Valerie Hayot-Sasson, Murali Emani, Sam Foreman, Zhen Xie, Diangen Lin, Maulik Shukla,
1009 Weili Nie, Josh Romero, Christian Dallago, Arash Vahdat, Chaowei Xiao, Thomas Gibbs, Ian
1010 Foster, James J. Davis, Michael E. Papka, Thomas Brettin, Rick Stevens, Anima Anandkumar,
1011 Venkatram Vishwanath, and Arvind Ramanathan. GenSLMs: Genome-scale language models
1012 reveal SARS-CoV-2 evolutionary dynamics. *The International Journal of High Performance
1013 Computing Applications*, 37(6):683–705, November 2023. ISSN 1094-3420. doi: 10.1177/
1014 10943420231201154. URL <https://doi.org/10.1177/10943420231201154>.
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

Part I

Appendix

Table of Contents

A Related Work	22
B Exploratory Applications	23
C Case Study	25
C.1 Predicting CAR-T Therapy Response for Refractory B Cell Lymphoma from Patient Single-Cell Profiles	25
C.2 Interpretable Cardiomyopathy Disease Subtype Prediction with Single-Nucleus Profile of Patient Heart Cells	26
D Experimental Details	28
D.1 Datasets Introduction	28
D.2 Agent Configurations	28
D.3 Memory Module Construction	29
D.4 Experts Discussion Construction Details	29
D.5 Embedding Quality on the scPerturb Benchmark	31
D.6 Ablation study setup	32
D.7 Hyperparameter Configuration	32
E Evaluation Details	35
E.1 Mean Squared Error (MSE)	35
E.2 Pearson Correlation Coefficient (PCC)	35
E.3 Coefficient of Determination (R^2)	35
E.4 Metrics with Differential Expression (DE)	35
E.5 LatentSpace Linear Separability	36
E.6 Perturbation Consistency	36
E.7 Latent Space Direct Organization	36
E.8 Linear Interpretability of Latent Space	36
F RAGBench Evaluation Details	38
G Detailed Algorithm Specifications	40
G.1 Agentic Retrieval System	40
G.2 Graph-based Multi-Expert Discussion	41
G.3 Code Implementation and Refinement Process	42
H Agent Communication Protocol Details	43
H.1 Protocol Design and Comparison	43
H.2 Protocol Implementation Details	43
I Cost Analysis	45
I.1 Training Infrastructure	45
I.2 Token Utilization and Cost Estimation	45
I.3 Cost-Effectiveness Analysis	45
I.4 Efficiency Analysis: Multi-Agent vs. Single-LLM Approaches	46
I.5 Token Utilization and Cost Comparison	46
I.6 End-to-End Workflow Efficiency	47
J Failure Case and Randomness Analyses	48

1080	J.1	Failurecases	48
1081	J.2	Bug-fix Solutions of Experiment Execution Module	49
1082	J.3	Variance of each runs	50
1083	J.4	Stability Enhancement with Confidence Regularization	50
1084			
1085	K	Performance Varies Across Different LLMs and AI Coders	51
1086	K.1	Experimental Setup	51
1087	K.2	Key Findings	52
1088	K.3	Analysis of Failure Modes	52
1089	K.4	Implications for Scientific AI Systems	53
1090	K.5	Performance with Smaller LLMs	53
1091			
1092	L	Designed Models	54
1093	L.1	Methodology	54
1094	L.2	Biological Interpretation of Architectural Choices	54
1095			
1096	M	Biological Analysis of Perturbation Results	60
1097	M.1	Pathway Analysis	60
1098	M.2	Results and Interpretation	62
1099	M.3	Quantitative Assessment	63
1100	M.4	Biological Significance	63
1101	M.5	Limitations and Considerations	63
1102	M.6	Implications for Model Development	64
1103			
1104	N	Additional Visualizations	65
1105	N.1	Comparative Performance Analysis	65
1106			
1107	O	Knowledge Base for Agentic Retrieval	66
1108			
1109	P	LLM-as-a-judge details	68
1110	P.1	Methods	68
1111	P.2	Prompts	68
1112	P.3	Example Output	72
1113	P.4	Detailed Results	76
1114	P.5	Inter-Judge Consistency Analysis	78
1115			
1116	Q	Human scientists' evaluation details	80
1117	Q.1	Methods	80
1118	Q.2	Detailed Results	80
1119			
1120	R	Prompt Templates	82
1121	R.1	Task Description input	82
1122	R.2	Task Analysis Collaboration Agents Settings	83
1123	R.3	Multi-experts Settings	90
1124			
1125	S	Detailed outputs from CellForge	95
1126	S.1	Data Parser	95
1127	S.2	Agentic Retrieval	97
1128	S.3	Task Analysis Report	99
1129	S.4	Model Design Module	110
1130			
1131	T	Detailed outputs from other research Agents	133

1131
1132
1133

A RELATED WORK

Agent Systems for Scientific Discovery Researchers have developed specialized AI systems spanning the entire research workflow: from literature analysis tools like PaperQA2 [103] and CHIME [43], to hypothesis generation frameworks that range from domain-specific idea creation [6, 88, 116] to comparative evaluations with expert proposals [102]. These systems increasingly leverage multi-agent architectures [35, 37, 101] to facilitate collaborative scientific reasoning. Implementation capabilities have advanced through scientific coding frameworks like SciCode [111] and MAgent-Bench [47], while benchmarks evaluate these capabilities across diverse domains [18, 55, 91, 99]. The integration of literature analysis with data-driven approaches has proven particularly effective for hypothesis generation [68, 79, 121], with several frameworks enhancing research ideation through structured feedback mechanisms [31, 87] and approaches to improve novelty and diversity [29, 44, 92]. End-to-end systems now attempt to unify these capabilities, including domain-general approaches like AI Scientist [75] and MLR-Copilot [66], alongside domain-specific implementations for chemistry [10], genomics [97], materials science [34], and medicine [81, 109]. Despite these advances, significant challenges remain in developing truly autonomous scientific systems, particularly regarding experimental rigor [60], falsification mechanisms [70], and comprehensive evaluation metrics [8, 27], as highlighted in recent surveys [24, 61, 93, 94].

AI Agents in Biomedical Research AI agents in biomedical research are rapidly evolving to simulate and accelerate the entire biomedical research workflow, from hypothesis generation to experimental protocol design to general scientific discovery. For instance, BioReason [25] interprets the functional impacts of genetic mutations, while POPPER [45] introduces a framework for validating free-form hypotheses through sequential falsification tests. These agents excel at reasoning but do not generate executable analysis pipelines as their primary output. Another category targets wet-lab experimental design. PerturboAgent [39], for example, is a self-planning agent designed to optimize the selection of genes for sequential Perturb-seq experiments, thereby guiding the next phase of lab work rather than creating a computational analysis model. A third category, including Biomni [46] and SpatialAgent [112], automates workflows by connecting existing software packages but is constrained by their static, predefined toolsets and limited code generation capabilities. STELLA [54] introduces autonomous tool discovery and reasoning template learning, boosting system performance through a self-evolving architecture. Yet its scope is largely limited to lightweight tool orchestration and biomedical question-answering; it stops short of designing novel AI models or automating *in-silico* experiments for biomedical research. This leaves an open opportunity for agentic frameworks explicitly aimed at AI model creation and end-to-end computational experimentation.

Single-Cell Perturbation Analysis Single-cell perturbation studies measure how cells respond to genetic or chemical interventions. The existing literature of *in-silico* approaches that predict post-perturbation cell states reflects a fundamental divergence in machine learning, with each paradigm showcasing distinct philosophies for modeling cellular responses. Earlier efforts, such as linear regression [22] or random forest feature selection [104], treated each gene or cell type in isolation. Deep generative models [41, 72, 74], conceptualize perturbations as latent space transformations through linear shifts or decompositions that separate biological covariates. In contrast, network-based methods [7, 57, 90, 96] explicitly incorporate biological knowledge via gene regulatory networks or cellular relationships. To further address the issue of cell heterogeneity, distribution alignment approaches such as optimal transport [12, 23] have been applied to machine learning models [51], matching the distribution of control cells with perturbed cells. The emergence of transformer architectures represents the latest paradigm shift. These architectures [21, 38, 62, 110] leverage pre-training at scale and self-attention mechanisms to model complex gene dependencies without explicit biological structure. This theoretical diversity creates a vast design space where selecting optimized architectures, representation strategies, and biological constraints remains highly context-dependent.

B EXPLORATORY APPLICATIONS

In addition to benchmarks with established baselines, we evaluated CELLFORGE on modalities where prior perturbation-response models are scarce or unavailable, including scATAC-seq and scCITE-seq datasets. These experiments are exploratory, demonstrating the framework’s ability to automatically design models that handle diverse data types and extreme sparsity. The Papalexi [84] dataset offers both RNA and protein measurements (CITE-seq), enabling assessment of cross-modality prediction. The Liscovitch [67] dataset presents the distinct challenge of predicting chromatin accessibility changes (scATAC-seq) rather than gene expression, while the Schiebinger [100] dataset examines responses to immune signaling molecules (cytokines).

For scATAC-seq (Liscovitch et al. [67]) and scCITE-seq (Papalexi et al. [84]), linear regression and random forest serve as reference points. While their performance is limited, CELLFORGE consistently surpasses them by generating architectures that integrate modality-specific embeddings, handle multi-modal inputs, and predict both RNA and protein responses, demonstrating versatility.

The performance advantages become more pronounced in challenging cross-modality scenarios. For CITE-seq protein measurements, CELLFORGE achieves 177% improvement in correlation ($PCC = 0.7495$ vs. 0.2704 for Random Forest). It also maintains superior performance even on fundamentally different modalities such as chromatin accessibility (scATAC-seq), achieving remarkable improvement in variance explained ($R^2 = 0.0678$ vs. 0.0040) and correlation for key regulatory regions ($PCC_{DE} = 0.6991$ vs. 0.0509).

For modalities lacking established models (scCITE-seq, scATAC-seq, cytokine), we employ Random Forest and Linear Regression using one-hot encoded perturbations concatenated with expression profiles as inputs. We are aware that, for the ATAC- and CITE-seq benchmarks, our comparisons rely on “simple” learners (linear regression and random forest). This choice is deliberate and stems from three factors: (i) to date no perturbation-response method has been published or benchmarked for these modalities [28, 63], making scRNA-centric models such as scGen [72] or scGPT [21] fundamentally incompatible with peak- or protein-level data; (ii) the few multimodal generative tools, like totalVI [33], MultiVI [5], and GLUE [16], that *can* process ATAC or CITE-seq were designed for data integration rather than counterfactual perturbation prediction [15] and therefore cannot address unseen perturbations; (iii) recent meta-analyses show that, when properly tuned, classical models often match or exceed specialised deep networks on sparse single-cell tasks [4, 63, 65]. Consequently, linear regression and random forest constitute strong, modality-agnostic baselines in the absence of purpose-built alternatives. Their limitations, however, underscore the need for an automatic, modality-aware framework: CELLFORGE generates custom architectures that handle the extreme sparsity of scATAC-seq and the multi-modal nature of CITE-seq, achieving state-of-the-art performance where no prior solution exists.

Table 5: DEG Recovery Performance Across Benchmark Datasets

Dataset	DEG Recall	ROC-AUC	PR-AUC
<i>Cytokine Perturbation – scRNA-seq Dataset</i>			
Schiebinger et al. [100]	0.535 ± 0.14	0.524 ± 0.08	0.105 ± 0.02
<i>Gene Knock Out Perturbation – scCITEseq Dataset</i>			
Papalexi et al. (RNA) [84]	0.509 ± 0.12	0.415 ± 0.05	0.115 ± 0.05
Papalexi et al. (Protein) [84]	0.420 ± 0.12	0.392 ± 0.25	0.121 ± 0.09
<i>Gene Knock Out Perturbation – scATACseq Dataset</i>			
Liscovitch et al. [67]	0.484 ± 0.12	0.097 ± 0.02	0.048 ± 0.02

chromatin accessibility perturbations (Liscovitch) or cross-modal protein predictions show lower but still meaningful performance.

The DEG recovery performance varies meaningfully across different perturbation modalities and experimental contexts. The highest DEG recall (77.9%) is achieved on the Norman dataset (in 2), which features comprehensive genetic interaction profiling with rich phenotypic readouts. In contrast, more challenging scenarios like chromatin accessibility perturbations (Liscovitch) or cross-modal protein predictions show lower but still

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

Table 4: Post-perturbation prediction results on datasets where existing baseline methods are scarce or unavailable. These results highlight the adaptability of CELLFORGE in automatically designing models for diverse perturbation modalities beyond standard benchmarks. Experiments are conducted under the unseen perturbation setting, with available baselines reproduced accordingly.

MODEL	$MSE \downarrow$	$PCC \uparrow$	$R^2 \uparrow$	$MSE_{DE} \downarrow$	$PCC_{DE} \uparrow$	$R^2_{DE} \uparrow$
<i>Cytokine Perturbation – scRNA-seq Dataset (Schiebinger et al. [100])</i>						
Unperturbed	0.0076	0.0007	0.0069	0.0980	0.0082	-0.6782
Random Forest	0.0762	0.2704	0.4186	0.0910	0.2124	0.2185
Linear Regression	0.4855	0.0785	0.0034	0.4359	0.0847	0.0013
GNN	0.0651	0.4127	0.3514	0.0827	0.2875	0.1982
Transformer	0.0543	0.4879	0.4122	0.0718	0.3196	0.2420
CellForge-Models	0.0428 \pm 0.0205	0.5697 \pm 0.0943	0.5043 \pm 0.0541	0.0144 \pm 0.0349	0.3396 \pm 0.0403	0.2832 \pm 0.1154
<i>Gene Knock Out Perturbation – scCITEseq (RNA) Dataset (Papalexi et al. [84])</i>						
Unperturbed	0.1509	0.0004	0.0017	0.6276	0.0007	-5.9142
Random Forest	0.0763	0.2124	0.4186	0.0911	0.2455	0.2185
Linear Regression	0.0764	0.0170	0.0254	0.0909	0.0218	0.0163
GNN	0.1215	0.6021	0.4114	0.2240	0.5807	0.3420
Transformer	0.1363	0.7715	0.5948	0.4565	0.4460	0.1956
CellForge-Models	0.0417 \pm 0.0051	0.6935 \pm 0.1995	0.3687 \pm 0.0651	0.0535 \pm 0.1566	0.6406 \pm 0.1940	0.2354 \pm 0.0224
<i>Gene Knock Out Perturbation – scCITEseq (Protein) Dataset (Papalexi et al. [84])</i>						
Unperturbed	0.4092	-0.0115	-0.9945	0.5974	-0.0081	-0.3652
Random Forest	0.0982	0.2704	0.0829	0.3071	0.4024	0.0466
Linear Regression	0.4901	0.3396	0.1241	0.4551	0.3087	0.3523
GNN	0.0625	0.5316	0.2987	0.0812	0.4021	0.2082
Transformer	0.1876	0.7773	0.5996	0.1674	0.7772	0.5041
CellForge-Models	0.0070 \pm 0.0387	0.7495 \pm 0.0653	0.6872 \pm 0.0956	0.2921 \pm 0.0045	0.7409 \pm 0.0970	0.5489 \pm 0.0749
<i>Gene Knock Out Perturbation – scATACseq Dataset (Liscovitch et al. [67])</i>						
Unperturbed	0.0426	0.0001	-0.0001	9.4980	0.0004	-9.7567
Random Forest	0.0432	0.0638	0.0040	0.0510	0.0509	0.0035
Linear Regression	0.5767	0.0486	0.0229	0.7750	0.0457	0.0021
GNN	0.0990	0.0794	0.0714	0.0170	0.0331	0.0169
Transformer	0.0012	0.0253	0.0298	0.0054	0.0114	0.0389
CellForge-Models	0.0327 \pm 0.0320	0.0855 \pm 0.0357	0.0678 \pm 0.0120	0.0406 \pm 0.0268	0.0691 \pm 0.3173	0.0640 \pm 0.0279

C CASE STUDY

EXTERNAL PILOT STUDY (N=2)

To further evaluate the accessibility and practical usability of CELLFORGE outside of controlled environments, we conducted a lightweight external pilot with two independent wet-lab researchers who had no prior exposure to the framework. Both participants selected real problems from their daily research practice (one focused on immunotherapy, the other on cardiovascular disease modeling) and attempted to solve them by following the Quickstart tutorial without additional assistance. We logged their completion success, number of interventions required, and total wall-clock time.

Study Protocol Each participant was given anonymized datasets resembling their own laboratory scenarios and asked to (i) set up the environment, (ii) load their task specification, (iii) trigger the automatic architecture design pipeline, and (iv) run training until a valid model checkpoint was produced. No step-by-step guidance was provided beyond the written Quickstart.

Results Both participants were able to complete their tasks successfully. User A (immunotherapy) finished the full pipeline in 67 minutes with one minor intervention (path correction when loading data). User B (cardiovascular) completed in 79 minutes with two interventions (resolving a missing dependency and adjusting GPU memory allocation). In both cases, the generated models reached non-trivial predictive performance on held-out validation sets, aligning with baseline numbers reported in internal benchmarks.

Observations Participants reported that the documentation was clear, and the frameworks modular design minimized coding effort. They highlighted that automatic error messages and fallback defaults were sufficient for resolving issues without developer intervention. Both noted that the process was significantly faster than manual model assembly in their typical workflows, estimating a 3–4× speed-up compared to their usual practice.

Conclusion This pilot demonstrates that CELLFORGE can be successfully adopted by independent wet-lab researchers with minimal computational training. The small number of interventions and the high success rate suggest that the frameworks Quickstart and design abstractions substantially lower the barrier to entry for real-world users.

C.1 PREDICTING CAR-T THERAPY RESPONSE FOR REFRACTORY B CELL LYMPHOMA FROM PATIENT SINGLE-CELL PROFILES

Background CAR-T cell therapy success critically depends on the functional composition of infusion products, where specific cellular subsets (memory stem cells, exhausted cells) disproportionately influence treatment outcomes. Traditional machine learning approaches treat all cells equally, failing to identify the therapeutically critical cell instances that determine response within the heterogeneous CAR-T product mixture.

Objective Developing a model to predict CAR-T therapy response by identifying and prioritizing critical cell instances from pre-treatment single-cell RNA sequencing data (input: 5,000-dimensional gene expression profiles from 109,151 cells across 32 patients; output: binary treatment response prediction with instance-level therapeutic potential scores).

Methods The model proposed by CellForge first transforms 5,000-dimensional single-cell expression profiles into compact embeddings using stacked residual layers and normalization. A patient-aware attention pooling module adaptively prioritizes informative cells, producing aggregated patient-level representations. The model jointly optimizes a classification objective with supervised contrastive loss to maximize the separation of responder and non-responder profiles. Performance was evaluated using a leave-one-patient-out cross-validation approach, reporting AUROC, average precision, and calibration metrics.

Results Across five cross-validation folds, baseline models including logistic regression, random forest, XGBoost, and multilayer perceptron (MLP) demonstrated moderate and highly variable

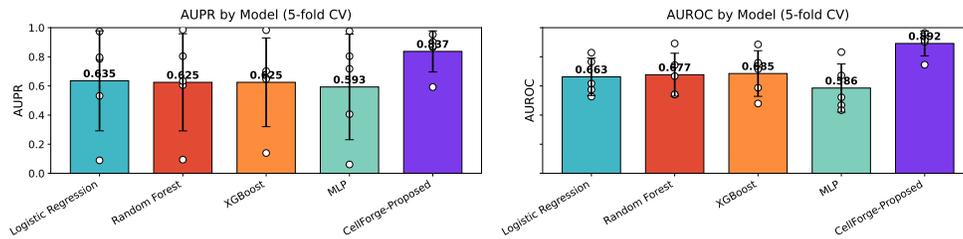


Figure 9: Comparison of baseline models and the **CellForge** proposed methods under 5-fold cross-validation. **CellForge** achieves the highest AUPR (0.837) and AUROC (0.892), clearly outperforming logistic regression, random forest, XGBoost, and MLP baselines.

performance, with mean AUPR values ranging from 0.47 to 0.64 and mean AUROC values between 0.57 and 0.71. In contrast, the proposed CellForge model consistently outperformed all baselines, achieving an average AUPR of 0.84 and AUROC of 0.89, with reduced variance across folds. Notably, CellForge maintained high F1-scores and Matthews correlation coefficients, indicating both robustness and balanced predictive capacity. These results demonstrate that selectively leveraging informative cellular signals yields substantially stronger and more reliable predictions of CAR-T therapy response compared to conventional machine learning approaches.

C.2 INTERPRETABLE CARDIOMYOPATHY DISEASE SUBTYPE PREDICTION WITH SINGLE-NUCLEUS PROFILE OF PATIENT HEART CELLS

Background Heart failure affects 23 million individuals worldwide, yet while single-nucleus RNA sequencing has revealed disease mechanisms at the cellular population level, patient-level analysis remains largely absent. Machine learning methods that can precisely classify individual patient disease states and systematically interpret important cellular subtypes could provide crucial insights for customized treatment strategies and mechanistic understanding.

Objective Developing machine learning models to classify patient disease states using single-nucleus RNA expression matrices as input, predicting cardiomyopathy disease states (Arrhythmogenic right ventricular Cardiomyopathy, Dilated Cardiomyopathy, Non-compaction Cardiomyopathy vs. Normal) for patients unseen during training. Model performance was evaluated using accuracy (ACC), F1-score, AUROC, Matthews correlation coefficient (MCC), and area under the precision-recall (AUPR) curve to assess classification robustness across class distributions.

Methods CellForge proposed and implemented a hierarchical neural network architecture comprising three sequential modules: a cell encoder that maps high-dimensional single-cell transcriptomic profiles to lower-dimensional cellular representations, an attention-based cell aggregation mechanism that generates patient-level embeddings from variable numbers of constituent cells, and a patient encoder optimized with contrastive learning to produce discriminative patient representations for disease classification. Single-cell RNA sequencing data underwent standard preprocessing procedures, including cellular subsampling, library size normalization, logarithmic transformation, and stochastic depth augmentation to account for technical variability. Model performance was assessed using hold-out validation on previously unseen patients, with evaluation metrics including classification accuracy, F1-score, area under the receiver operating characteristic curve, Matthews correlation coefficient, and area under the precision-recall curve to comprehensively characterize predictive performance across cardiac disease phenotypes.

Results The model achieved excellent discrimination among dilated cardiomyopathy, arrhythmogenic right ventricular cardiomyopathy, noncompaction cardiomyopathy and normal hearts, with an overall accuracy of 0.9847, a weighted F1 score of 0.9841 and macroaveraged ROCAUC and PRAUC values of 0.9997 and 0.9980, respectively. Integrated-gradient analysis revealed biologically meaningful drivers: vCM1.0 and vCM2 cardiomyocyte states distinguish left- and right-ventricular programs, with vCM2 marked by cardioprotective gene expression such as PRELID2 and CDH13. Among fibroblasts, the vFB2 state stood out, characterized by distinct ECM signatures and pro-

D EXPERIMENTAL DETAILS

D.1 DATASETS INTRODUCTION

Our study leverages six publicly available single-cell perturbation datasets from the scPerturb [85] collection, encompassing diverse perturbation modalities and cell types. These datasets provide a foundation for evaluating the scientific quality of AI-generated analyses across various biological contexts.

Adamson et al. [2] (CRISPRi): Employing Perturb-seq to study the unfolded protein response (UPR) in K562 lymphoblasts through single and combinatorial CRISPR interference (CRISPRi) perturbations. Approximately 100 gene targets were profiled, enabling high-resolution functional clustering and revealing distinct activation patterns across UPR branches.

Norman et al. [82] (CRISPRa): Utilizing CRISPR activation (CRISPRa) in K562 cells, this dataset explores genetic interaction manifolds derived from single-cell transcriptional phenotypes. The study provides insights into regulatory pathway ordering and mechanistic elucidation of synergistic interactions.

Liscovitch et al. [67] (ATAC-seq): Employing CRISPRsciATAC, a single-cell combinatorial indexing assay, to delineate the genetic determinants of chromatin accessibility in human myelogenous leukemia K562 cells. Targeting 105 chromatin-related genes via CRISPR-Cas9, the study generated chromatin accessibility profiles for approximately 30,000 single cells. Key findings include correlations between the loss of specific chromatin remodelers and global changes in chromatin accessibility. Notably, EZH2 depletion was associated with enhanced accessibility in heterochromatic regions linked to embryonic development and with activation of genes in the HOXA and HOXD clusters. This high-throughput approach offers valuable insights into the role of chromatin modifiers in regulating gene expression and their implications in disease states.

Papalexli et al. [84] (CITE-seq): Combining CRISPR-Cas9 perturbations with single-cell RNA and surface protein measurements in THP-1 monocytes. It investigates the molecular regulation of inhibitory immune checkpoints, particularly PD-L1 expression, and introduces the mixscape computational framework to enhance signal-to-noise ratio in single-cell screens.

Srivatsan et al. [106] (sci-Plex): Employing sci-Plex, this dataset profiles transcriptional responses of A549, K562, and MCF7 cancer cell lines to 188 small-molecule compounds across multiple doses. Approximately 650,000 single-cell transcriptomes were generated, uncovering intercellular heterogeneity and commonalities in drug responses.

Schiebinger et al. [100] (cytokine perturbation): Applying optimal transport analysis to scRNA-seq data from mouse embryonic stem cells undergoing reprogramming with cytokine treatments. The dataset captures developmental trajectories and identifies transcription factors and paracrine signals influencing cell fate decisions.

Collectively, these datasets encompass a range of perturbation types including CRISPRi, CRISPRa, CRISPR-Cas9, small-molecule drugs, and cytokines across various human and mouse cell lines. They provide a robust foundation for evaluating the scientific quality and reliability of AI-generated analyses in single-cell biology.

D.2 AGENT CONFIGURATIONS

In our experiments, we employed five LLMs API to generate responses: Claude 3.7, OpenAI o1, DeepSeek-R1, Qwen-Plus, and Llama 3.1. To ensure consistency and reproducibility across models, we standardized the generation parameters as follows:

Temperature: Set to 0.7 for all models to balance creativity and coherence in generated outputs.

Top-p (nucleus sampling): Fixed at 0.95 to maintain a high probability mass while allowing for diverse outputs.

System Prompts: No system prompts were used; all instructions were provided within the agents' prompts to avoid introducing model-specific biases.

1512 These configurations align with recommended settings for models. By maintaining uniform settings
1513 across all models, we aimed to ensure a fair comparison and reliable evaluation of their performance.
1514

1515 D.3 MEMORY MODULE CONSTRUCTION 1516

1517 **Shared Knowledge Infrastructure.** Both Task Analysis and Method Design modules rely on
1518 a shared hybrid knowledge infrastructure comprising (1) a symbolic memory module that stores
1519 structured outputs from agents, and (2) a vector-based retrieval system built on top of Sentence-BERT
1520 embeddings and external APIs (PubMed, GitHub). The memory module is incrementally constructed
1521 as each agent contributes new findings or insights, while the vector database supports RAG-style
1522 retrieval of external literature. This shared infrastructure enables bi-directional communication
1523 between agents within each module and supports consistent knowledge propagation across modules.
1524 See Appendix D.3 for implementation details.

1525 **Collaborative Agents Shared Memory Module in Task Analysis.** Instead of operating in isolation,
1526 the Dataset Analyst, Problem Investigator, and Baseline Assessor interact via the shared memory
1527 module and query interface. Each agent incrementally updates the memory module with its findings,
1528 while continuously polling for updates from other agents. For example, once the Dataset Analyst infers
1529 perturbation modalities and cell types, the Problem Investigator revises its hypothesis formulation
1530 accordingly. Agents operate asynchronously but synchronize their conclusions through a shared
1531 JSON-based communication protocol, allowing for self-consistency checks and iterative refinement
1532 of the task representation. This collaborative reasoning leads to a structured task analysis report
1533 passed to the Method Design module.
1534

1535 **Graph-Based Expert Shared Memory Module in Method Design.** In the Method Design module,
1536 domain experts are instantiated as nodes in a dynamic undirected graph. These expert agents
1537 exchange proposals and critiques via message-passing rounds governed by graph neural network
1538 operations. Throughout the discussion, the Critic Agent agent monitors logical coherence and
1539 suggests refinements. Each expert agent has read-write access to the shared memory module and
1540 can retrieve relevant prior knowledge from Agentic Retrieval. Updates to the architectural plan are
1541 written back to the graph, enabling history-aware(get messages and suggestions from the former
1542 round), convergent model refinement.

1543 D.4 EXPERTS DISCUSSION CONSTRUCTION DETAILS 1544

1545 To enable structured, reproducible reasoning across diverse perturbation modeling tasks, we construct
1546 the multi-agent expert discussion system through two key stages: expert role selection and dynamic
1547 collaboration graph construction.
1548

1549 Based on the task analysis report, a set of relevant expert agents is selected by matching task attributes
1550 against a curated registry of expert types. The selected experts are grouped into five broad categories
1551 to ensure comprehensive domain coverage: **(i) Data Engineering and Preprocessing.** A Data Expert
1552 is instantiated to address normalization, quality control, feature selection, and batch correction issues
1553 tailored to the input modality. **(ii) Model Design and Scalability.** The Model Architecture Expert
1554 and Deep Learning Expert are responsible for proposing architectures that balance expressiveness,
1555 interpretability, and scalability, considering modality-specific modeling needs. **(iii) Biological
1556 Plausibility.** Single Cell Experts such as the Pathway Expert, Drug Response Expert, and Omics
1557 Modality Expert contribute domain knowledge to align model components with known biological
1558 mechanisms, including gene regulatory networks, cytokine signaling, or pharmacodynamics. **(iv)
1559 Training and Optimization.** A Training Expert is responsible for selecting and justifying the
1560 learning algorithm, optimization strategy, regularization, and validation scheme suitable for the data
1561 structure and model complexity. **(v) Self-Critique and Evaluation.** A Critic agent is included in
1562 every discussion to promote internal scrutiny, consistency checks, and critical reflection over model
assumptions and claims.

1563 For example, in a gene knockout task, the system may instantiate the Data Expert to inspect whether
1564 the scRNA-seq matrix is properly normalized, whether cell and gene identifiers are standardized, and
1565 whether preprocessing sufficiently preserves perturbation-related variation. The Model Architecture
Expert and Deep Learning Expert are instantiated to co-design a gene-centric model that integrates

1566 perturbation-aware attention and captures target gene dependent regulatory effects. The Pathway
1567 Expert is instantiated to evaluate the role of the target gene within interferon signaling cascades,
1568 while the Omics Modality Expert assesses whether transcriptomic changes resulting from target gene
1569 ablation are robustly captured by scRNA-seq alone. The Training Expert selects dropout-regularized
1570 contrastive training and a cell-type-aware sampling scheme to stabilize optimization. The Statistics
1571 Expert designs a differential expression based evaluation framework and quantifies the significance of
1572 target gene induced shifts using FDR-corrected effect sizes. Finally, the Critic Agent is instantiated
1573 to identify overfitting risks in rare knockout subsets, challenge latent space linearity assumptions, and
1574 refine model outputs for interpretability and robustness.

1575 All experts are set with role-specific prompts (Appendix R.3), crafted in a zero-shot reasoning format.
1576 These prompts are conditioned on the shared Task Analysis report and elicit structured outputs,
1577 including modeling choices, biological justification, and critiques of others proposals.

1578 Formally, the expert set $E^{(k)}$ for task k is derived by:

$$1579 \quad E^{(k)} = \text{SelectExperts}(\text{TaskAnalysisReport}_k)$$

1582 Once instantiated, the experts are organized into an undirected collaboration graph $G^{(k)} = (S, E^{(k)})$,
1583 where each node $E^{(i)} \in E^{(k)}$ represents an expert role. The Critic Agent node S is fully connected
1584 to all others, serving both as a dialectical evaluator and proposal aggregator.

1585 Each expert begins with an initial model proposal $m_0^{(i)}$ and a confidence score initialized to zero
1586 $c_0^{(i)} = 0$. During the discussion, agents iteratively update their proposals and confidence scores
1587 through message passing on the graph. Each round incorporates structured information exchange,
1588 where agents revise their reasoning in response to input from their neighbors, weighted by relevance.

1590 This structured and interpretable procedure allows CELLFORGE to generate scientifically grounded,
1591 multimodally coherent model designs that are not only technically sound but also biologically
1592 meaningful.

1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619

Table 6: Performance comparison on scPerturb datasets and benchmark tasks[9] (all values are in %). Results show CELLFORGE consistently outperforms scGPT, Geneformer, CPA, STATE, scVI, and PCA across multiple metrics and perturbation types. Each score represents the average of five independent runs, with higher values indicating better performance. These models operate by converting complex gene expression data into meaningful vector representations embeddings, which are then used to predict cellular responses to perturbations.

MODEL	TOP5 LIN \uparrow	TOP1 LIN \uparrow	PERT CONS \uparrow	TOP5 KNN \uparrow	TOP1 KNN \uparrow	SPEAR CORR \uparrow	STRUCT INT \uparrow
<i>Drug Perturbation (Srivatsan Dataset [106])</i>							
PCA	1.2	0.9	0.4	2.1	1.8	8.4	48.3
scVI [71]	1.5	1.0	0.7	2.4	2.0	10.3	49.1
STATE [3]	5.5	3.9	9.4	5.5	4.8	17.9	53.9
CPA [73]	5.1	3.7	9.8	5.3	4.7	17.4	53.8
scGPT[21]	5.2	4.4	11.4	5.6	5.1	18.8	54.2
Geneformer[110]	4.4	3.1	0.9	5.1	4.8	17.3	54.1
CellForge-Model	7.0	4.2	11.4	6.4	5.3	19.1	54.5
<i>Gene Knock Out Perturbation (Adamson Dataset [2])</i>							
PCA	0.8	0.3	1.1	14.2	13.5	72.4	90.8
scVI [71]	1.0	0.4	1.6	15.8	15.1	76.3	92.1
STATE [3]	2.2	0.8	5.1	24.6	23.5	86.2	95.7
CPA [73]	2.0	0.7	4.8	24.4	22.8	85.6	95.8
scGPT [21]	2.2	0.8	5.6	26.2	25.5	87.3	96.1
Geneformer [110]	2.1	0.8	4.3	25.9	24.1	86.6	95.9
CellForge-Model	2.4	0.9	6.9	26.6	25.9	89.9	96.0
<i>Cytokine Perturbation (Schiebinger Dataset [100])</i>							
PCA	0.7	1.8	1.9	4.1	3.6	52.1	50.4
scVI [71]	1.1	2.3	2.1	4.8	4.1	54.9	51.7
STATE [3]	2.2	4.4	4.7	8.0	6.3	67.1	57.0
CPA [73]	2.0	4.1	4.2	7.4	6.3	65.1	56.4
scGPT[21]	2.1	4.8	4.6	8.2	5.5	66.9	57.1
Geneformer[110]	1.4	4.2	4.4	8.3	9.9	68.2	57.6
CellForge-Model	2.5	5.3	4.9	8.6	8.8	68.5	59.6

D.5 EMBEDDING QUALITY ON THE scPERTURB BENCHMARK

While CELLFORGE primarily performs gene expression prediction following perturbations, the quality of learned representations is equally important for biological interpretability. Following evaluation practices established in previous works [21, 110], we benchmark CELLFORGE against specialized foundation models (scGPT & Geneformer) on representation quality metrics (Table 6).

To ensure fair comparison, we follow the previous zero-shot benchmarking framework [9], which evaluates transcriptomic foundation models without task-specific fine-tuning. Specifically, perturbation embeddings for both scGPT and Geneformer are extracted directly from their pre-trained backbones with no fine-tuning performed on any evaluation dataset. This represents pure zero-shot performance, making the comparison particularly stringent for our method, as baseline models leverage extensive pre-training on large-scale datasets while CELLFORGE operates without any pre-training advantages. All models are evaluated under identical zero-shot conditions using standardized downstream metrics including logistic regression for separability assessment and cosine clustering for consistency measurement.

We assess different aspects of latent space organization across five dimensions: **(1) Linear separability metrics** (TOP5_LIN, TOP1_LIN) measure how distinguishable different perturbation types are in the latent space. The `top5_lin` score of 0.070 achieved by CELLFORGE for drug perturbations (vs. 0.052 for scGPT) indicates that 7.0% of test samples have their correct perturbation label among the top 5 predictions when using a linear classifier trained on the latent embeddings. This improvement suggests CELLFORGE learns representations where perturbation effects are more linearly separable, facilitating downstream analyses that rely on perturbation classification. **(2) Perturbation consistency** (PERT_CONS) quantifies whether cells with the same perturbation cluster more tightly than random controls, essentially measuring the signal-to-noise ratio of perturbation effects in the latent space. For gene knockouts, CELLFORGE achieves a consistency of 0.069. This indicates that CELLFORGE creates a latent space where cells experiencing the same perturbation are more reliably grouped together, reflecting better capture of perturbation-specific biological responses. **(3) Local structure in the latent space** is assessed through nearest-neighbor metrics (TOP5_KNN, TOP1_KNN), which evaluate whether perturbations form locally coherent clusters. For drug perturbations, CELLFORGE achieves a TOP5_KNN score of 0.064 vs. 0.056 for scGPT, indicating that a higher

1674 proportion of test samples have correctly labeled neighbors in embedding space. **(4) The Spearman**
 1675 **correlation metric** (SPEAR_CORR) evaluates how accurately the latent embeddings can be mapped
 1676 back to the original gene expression space using a linear transformation. The score of 0.191 for
 1677 drug perturbations (vs. 0.188 for scGPT) represents a higher rank correlation between predicted and
 1678 actual expression values after linear decoding. **(5) Structural integrity** (STRUCT_INT) measures
 1679 how well control-perturbation relationships are preserved in the latent space. **0.596** for cytokine
 1680 perturbations (vs. 0.571 for scGPT) indicates that CELLFORGE better maintains the biological
 1681 relationship between control and perturbed states for complex signaling cascades.

1682 1683 D.6 ABLATION STUDY SETUP

1684
1685 Table 3 presents a systematic ablation analysis evaluating the individual and combined contributions
 1686 of core framework components. Each component represents a distinct architectural or methodological
 1687 choice within our system, with performance measured across three perturbation datasets using
 1688 standardized evaluation protocols.

1689 **Basic Version (Baseline):** Represents the minimal CellForge configuration without retrieval aug-
 1690 mentation or collaborative reasoning mechanisms. This baseline employs only three agents in Task
 1691 Analysis Module with direct LLM-based model generation without external knowledge integration or
 1692 multi-agent collaboration, serving as our reference point for measuring component contributions.

1693 **Normal RAG:** Implements standard retrieval-augmented generation using only cosine similarity-
 1694 based document ranking with Sentence-BERT embeddings. This approach performs static keyword
 1695 matching against a curated knowledge base of 46 single-cell perturbation articles (detailed in Ap-
 1696 pendix O) and dynamic search by PubMed, until cosine similarity reaches without adaptive query
 1697 refinement.

1698 **Agentic Retrieval:** Employs our proposed alternating breadth-first and depth-first search strategy
 1699 that enables autonomous knowledge discovery and dynamic query expansion. Unlike standard RAG
 1700 approaches that use static keywords, this system autonomously evolves from basic queries (e.g.,
 1701 "single cell perturbation prediction") to sophisticated technical searches (e.g., "Transformer VAE
 1702 GNN architectures") through BFS layers that explore diverse research directions and DFS layers that
 1703 trace citation networks to access implementation details.

1704 **Graph-Based Discussion:** Implements the multi-expert collaborative reasoning framework where
 1705 domain experts form an undirected collaboration graph. Each expert node maintains confi-
 1706 dence scores that evolve through discussion rounds using the formula $c_t^{(i)} = \lambda_1 \cdot c_{t-1}^{(i)} + \lambda_2 \cdot$
 1707 $\text{CriticAgentScore}(m_t^{(i)}, S) + \lambda_3 \cdot \frac{1}{k-1} \sum_{j \neq i} \text{PeerScore}(m_t^{(i)}, E^{(j)})$ with weights $(\lambda_1, \lambda_2, \lambda_3) =$
 1708 $(0.3, 0.4, 0.3)$. Discussion terminates when all experts' confidence scores exceed threshold $\tau = 0.8$
 1709 with minimal variance $\epsilon = 0.03$.

1710
1711 **Round-Robin Discussion:** Experts speak in predetermined sequential order (Data Expert Single-
 1712 Cell Expert Deep Learning Expert Critic Agent) without peer evaluation between domain experts.
 1713 Confidence updates use $c_t^{(i)} = \lambda_1 \cdot c_{t-1}^{(i)} + \lambda_2 \cdot \text{CriticAgentScore}(m_t^{(i)}, S)$ with $(\lambda_1, \lambda_2) = (0.7, 0.3)$.

1714 **Moderator-centered Discussion:** All domain experts simultaneously submit proposals to the Critic
 1715 Agent without inter-expert communication. The Critic Agent evaluates proposals independently and
 1716 provides feedback using the same confidence update formula as Round-Robin.

1717 1718 1719 D.7 HYPERPARAMETER CONFIGURATION

1720
1721 Our framework employs the following hyperparameter settings, determined through empirical valida-
 1722 tion on held-out scientific tasks:

1723 1724 D.7.1 AGENTIC-RETRIEVAL HYPERPARAMETERS

1725
1726 We conducted ablations on key parameters of the Agentic Retrieval module. As shown in Table 8,
 1727 increasing the retrieval budget (L_{\max}) improves answer quality but incurs higher token cost and
 latency. Score filtering and memory retrieval mechanisms contribute significantly to quality gains.

Table 7: Hyperparameter Configuration

Module	Parameter	Value
Agentic Retrieval	L_{\max}	10
	τ	0.8
	ϵ	0.5
Expert Discussion	τ	0.8
	ϵ	1
	$(\lambda_1, \lambda_2, \lambda_3)$	(0.3, 0.4, 0.3)

Notably, while larger L_{\max} and stricter thresholds (τ , ϵ) can slightly improve answer quality, they introduce significantly higher computational cost. The incorporation of memory-based retrieval and score-based filtering mechanisms plays a key role in maximizing answer informativeness while maintaining a reasonable token cost, demonstrating the importance of adaptive and context-aware retrieval strategies. The chosen default setting achieves the best balance between informativeness and efficiency.

Table 8: Ablation study on Agentic Retrieval configuration. The default uses $L_{\max} = 10$, $\tau = 0.8$, $\epsilon = 0.5$, with memory-based adaptive retrieval and score filtering. Increasing L_{\max} improves information coverage but raises cost. Strict τ and ϵ improve precision but reduce flexibility. Score filtering and memory retrieval notably improve quality-to-cost ratio.

Setting	L_{\max}	τ	ϵ	Answer Quality	Avg Time (s)	Token Cost (x)
Default (ours)	10	0.8	0.5	87.3	31.2	1.00x
Smaller L_{\max}	8	0.8	0.5	83.1	24.7	0.85x
Smaller L_{\max}	5	0.8	0.5	83.0	20.4	0.68x
Larger L_{\max}	12	0.8	0.5	88.2	45.6	1.42x
Larger L_{\max}	15	0.8	0.5	88.9	70.6	2.06x
Higher τ threshold	10	0.9	0.5	89.1	35.0	1.27x
Lower τ threshold	10	0.7	0.5	80.5	24.1	0.90x
Stricter ϵ	10	0.8	0.2	87.9	35.5	1.45x
Looser ϵ	10	0.8	0.8	82.6	28.6	0.82x

D.7.2 GRAPH-BASED DISCUSSION HYPERPARAMETERS

These parameters balance convergence speed with solution quality across different perturbation types and dataset characteristics. We observed that the Expert Discussion module particularly benefits from a higher weight on Critic Agent evaluation (λ_2), which promotes more rigorous scientific validation.

Table 9: Ablation study on stopping criteria τ and ϵ in the graph-based discussion. Baseline uses $\tau = 0.8$, $\epsilon = 0.5$. Reducing ϵ enforces stricter agreement, and increasing τ demands higher proposal quality, both of which incur additional cost.

Setting	τ	ϵ	Avg Rounds	Avg Score	Token Cost (x)
Default (ours)	0.8	0.03	4.2	85.8	1.00x
Stricter τ	0.85	0.03	5.2	89.6	1.36x
Looser τ	0.75	0.03	3.5	82.9	0.82x
Very strict τ	0.90	0.03	6.3	89.9	1.65x
Stricter ϵ	0.8	0.02	6.5	89.2	1.72x
Very strict ϵ	0.8	0.01	8.1	89.4	2.15x
Looser ϵ	0.8	0.04	3.8	87.2	0.78x
Very Loose ϵ	0.8	0.05	3.4	82.1	0.71x

The analysis in Table 9 demonstrates that stricter stopping conditions (*e.g.*, higher τ or lower ϵ) lead to more rounds of discussion with higher average confidence scores, but at the expense of increased

token cost. Therefore, our selected configuration ($\tau = 0.8$, $\epsilon = 0.03$) provides a favorable trade-off, ensuring both convergence and cost-effectiveness across diverse scientific tasks.

Table 10: Ablation study on the confidence score update components with $\tau = 0.8$, $\epsilon = 1$. The full model uses $(\lambda_1, \lambda_2, \lambda_3) = (0.3, 0.4, 0.3)$ and serves as the baseline ($1\times$ for token cost). Average rounds, confidence scores, time, and token costs are the average scores of five runs of experiments. The table shows that our parameter selection is optimal in terms of time, token cost, and effectiveness.

Setting	λ_1	λ_2	λ_3	Rounds	Avg Score	Avg Time	Token Cost(x)
Default (ours)	0.3	0.4	0.3	4.2	88.8	36.1	1.00x
larger λ_2	0.2	0.6	0.2	6.8	89.4	69.5	2.53x
smaller λ_2	0.4	0.2	0.4	4.0	86.0	34.4	0.97x
$\lambda_3 > \lambda_1$	0.2	0.4	0.4	4.2	88.0	40.1	1.11x
$\lambda_3 > \lambda_1$	0.1	0.4	0.5	4.2	87.4	41.4	1.32x
$\lambda_3 < \lambda_1$	0.4	0.4	0.2	4.0	87.8	35.0	1.20x
$\lambda_3 < \lambda_1$	0.5	0.4	0.1	4.0	87.0	34.4	1.35x
No Historical Memory	0.0	0.5	0.5	3.8	85.2	30.1	0.75x
No Critic Agent Evaluation	0.5	0.0	0.5	2.2	85.0	25.8	0.77x
No Peer Feedback	0.5	0.5	0.0	4.0	85.2	28.1	0.82x
Peer-Only (No Memory, No SC)	0.0	0.0	1.0	7.0	86.0	124.5	3.15x
Critic Agent Only	0.0	1.0	0.0	4.0	86.8	33.6	0.65x

As shown in Table 10, ablations on the confidence update mechanism reveal that the critic agent evaluation (λ_2) contributes most significantly to performance.

Interestingly, peer-only and critic-only settings each show limitations in stability or generalization. These findings support the necessity of integrating diverse feedback signals in our confidence update formulation.

1836 E EVALUATION DETAILS

1837
1838 This appendix provides detailed formulations of the hierarchical metrics used in our benchmark
1839 evaluation of transcriptomics machine learning models for perturbation analysis.

1841 E.1 MEAN SQUARED ERROR (MSE)

1842
1843 This metric measures the average squared difference between the true and predicted gene expression
1844 vectors, quantifying overall prediction error. Let $Y_i, \hat{Y}_i \in \mathbb{R}^{d'}$ be the true and predicted expression
1845 vectors for sample i . Then

$$1846 \text{MSE} = \frac{1}{n \cdot d'} \sum_{i=1}^n \|Y_i - \hat{Y}_i\|_2^2.$$

1849 E.2 PEARSON CORRELATION COEFFICIENT (PCC)

1850
1851 This metric assesses the strength of the linear association between predicted and true expression
1852 profiles across all samples. Define the sample means

$$1853 \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i, \quad \bar{\hat{Y}} = \frac{1}{n} \sum_{i=1}^n \hat{Y}_i.$$

1854
1855 Then

$$1856 \text{PCC} = \frac{\sum_{i=1}^n \langle Y_i - \bar{Y}, \hat{Y}_i - \bar{\hat{Y}} \rangle}{\sqrt{\sum_{i=1}^n \|Y_i - \bar{Y}\|_2^2} \sqrt{\sum_{i=1}^n \|\hat{Y}_i - \bar{\hat{Y}}\|_2^2}}.$$

1863 E.3 COEFFICIENT OF DETERMINATION (R^2)

1864
1865 This metric quantifies the proportion of variance in the true gene expression data that is captured
1866 by the models predictions. It provides an interpretable measure of model fit, with higher values
1867 indicating better predictive performance. Let $Y_i, \hat{Y}_i \in \mathbb{R}^{d'}$ be the true and predicted expression
1868 vectors for sample i , and let \bar{Y} denote the mean of the true expression vectors. Then

$$1869 R^2 = 1 - \frac{\sum_{i=1}^n \|Y_i - \hat{Y}_i\|_2^2}{\sum_{i=1}^n \|Y_i - \bar{Y}\|_2^2}.$$

1873 E.4 METRICS WITH DIFFERENTIAL EXPRESSION (DE)

1874
1875 Differential expression highlights the genes whose changes drive the biological response to a pertur-
1876 bation, focusing evaluation on the most informative signals. Let $\{Y_{p,i}\}_{i=1}^{n_p}$ and $\{Y_{c,i}\}_{i=1}^{n_c}$ be the true
1877 expression vectors under perturbation and control, respectively, with $Y_{p,i}, Y_{c,i} \in \mathbb{R}^{d'}$. For each gene
1878 $g = 1, \dots, d'$, compute the mean expression

$$1879 \bar{Y}_{p,g} = \frac{1}{n_p} \sum_{i=1}^{n_p} Y_{p,i,g}, \quad \bar{Y}_{c,g} = \frac{1}{n_c} \sum_{i=1}^{n_c} Y_{c,i,g}.$$

1880
1881 Quantify the change by the logfoldchange

$$1882 \text{LFC}_g = \log_2 \frac{\bar{Y}_{p,g} + \epsilon}{\bar{Y}_{c,g} + \epsilon},$$

1883
1884 with small $\epsilon > 0$ to avoid division by zero (or by the raw difference $\Delta_g = \bar{Y}_{p,g} - \bar{Y}_{c,g}$). Rank genes
1885 by $|\text{LFC}_g|$ (or $|\Delta_g|$), and select the top $K = 20$ as the DE set:

$$1886 \text{DE} = \{g : \text{rank}_g(|\text{LFC}|) \leq K\}, \quad K = 20.$$

1887
1888 Subsequent metrics (MSE, PCC, R^2) are then computed only over $g \in \text{DE}$ to assess performance on
1889 these key drivers of perturbation response.

1890 E.5 LATENTSPACE LINEAR SEPARABILITY

1891 This metric evaluates if a model’s latent space distinguishes between different perturbations using
1892 linear probing. Given a frozen encoder mapping $g_\phi : x_i \mapsto z_i \in \mathbb{R}^d$, train a linear classifier

$$1893 \hat{y} = \text{softmax}(Wz + b), \quad W \in \mathbb{R}^{c \times d}, \quad b \in \mathbb{R}^c,$$

1894 to predict one of c perturbation classes. For n test samples with true labels y_i ,
1895

$$1896 \text{Top-1} = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{\arg \max_j \hat{y}_{ij} = y_i\}, \quad \text{Top-5} = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{y_i \in \text{Top5}(\hat{y}_i)\}.$$

1901 E.6 PERTURBATION CONSISTENCY

1902 This metric assesses the consistency with which a model represents perturbations between different
1903 samples and batches. Let \mathcal{P} be the set of all gene perturbations. For each $p \in \mathcal{P}$, suppose we have n_p
1904 embedding vectors

$$1905 \{z_{p,i} \in \mathbb{R}^d \mid i = 1, \dots, n_p\}.$$

1906 Define the *mean cosinesimilarity score*

$$1907 S_p = \frac{1}{n_p^2} \sum_{i=1}^{n_p} \sum_{j=1}^{n_p} \frac{\langle z_{p,i}, z_{p,j} \rangle}{\|z_{p,i}\| \|z_{p,j}\|}.$$

1908 Let $\{S_{q_k}\}_{k=1}^K$ be the corresponding scores for K unexpressedgene controls q_k . The empirical pvalue
1909 for perturbation p is

$$1910 \pi_p = \frac{\max\{\#\{k : S_{q_k} \leq S_p\}, 1\}}{K}.$$

1911 Finally, the overall *consistency rate* is

$$1912 C = \frac{|\{p \in \mathcal{P} : \pi_p < 0.05\}|}{|\mathcal{P}|},$$

1913 i.e., the fraction of perturbations whose embeddings are significantly more selfsimilar than the null.
1914

1920 E.7 LATENT SPACE DIRECT ORGANIZATION

1921 This metric evaluates the degree to which perturbation clusters are locally organized in the latent
1922 space, using the k-Nearest Neighbors (kNN) classification. Let $\{z_i\}_{i=1}^{n_q}$ and $\{z_j\}_{j=1}^{n_r}$ be the latent
1923 embeddings for the query and reference sets, with the corresponding labels y_i and y_j . Set

$$1924 k = \lfloor \sqrt{n_r} \rfloor.$$

1925 For each query index i , let $N_k(i) \subset \{1, \dots, n_r\}$ be the reference index k whose embeddings
1926 minimize $\|z_i - z_j\|_2$. Then the *kNNclassification accuracy* is

$$1927 \text{Accuracy}_{\text{kNN}} = \frac{1}{n_q} \sum_{i=1}^{n_q} \mathbf{1}\left[y_i = \arg \max_{c \in C} \sum_{j \in N_k(i)} \mathbf{1}[y_j = c]\right],$$

1928 where C denotes the set of all perturbation labels.
1929

1933 E.8 LINEAR INTERPRETABILITY OF LATENT SPACE

1934 Let $Z \in \mathbb{R}^{n \times h}$ be the frozen-encoder outputs and train a linear MLP, $\hat{Y} = h(Z) \in \mathbb{R}^{n \times d'}$. We define
1935 two metrics: *Spearman correlation* and *structural integrity*.
1936

1937 **Spearman Correlation** This measures how accurately the latent embeddings can be decoded back
1938 into gene expression data using a simple linear transformation. The *Spearman correlation* ρ is defined
1939 as

$$1940 \rho = 1 - \frac{6 \sum_{i=1}^n [\text{rank}(Y_i) - \text{rank}(\hat{Y}_i)]^2}{n(n^2 - 1)},$$

1941 where $\text{rank}(\cdot)$ returns the withinsample rank vector.
1942

1944 **Structural Integrity** This metric assesses how effectively the model maintains the relationship
 1945 between control and perturbation conditions within each biological batch. For $b = 1, \dots, B$ batches
 1946 with n_b samples each, let

$$1947 \tilde{Y}_{\text{pred}}^{(b)} = Y_{\text{pred}}^{(b)} - Y_{\text{pred,ctrl}}^{(b)}, \quad \tilde{Y}_{\text{act}}^{(b)} = Y_{\text{act}}^{(b)} - Y_{\text{act,ctrl}}^{(b)}.$$

1949 Then

$$1950 D = \frac{1}{B} \sum_{b=1}^B \frac{1}{n_b} \|\tilde{Y}_{\text{pred}}^{(b)} - \tilde{Y}_{\text{act}}^{(b)}\|_F, \quad D_{\text{max}} \approx \frac{2}{B} \sum_{b=1}^B \frac{1}{n_b} \|\tilde{Y}_{\text{act}}^{(b)}\|_F,$$

1953 and the *structural integrity* is

$$1954 \text{SI} = 1 - \frac{D}{D_{\text{max}}},$$

1955 with higher SI indicating better preservation of controlperturbation structure.

1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997

F RAGBENCH EVALUATION DETAILS

To evaluate the performance of CellForge’s Agentic Retrieval system in the Task Analysis Module, we employ RAGBench[27].

We first align our systems outputs to the format expected by RAGBench. Each output record must include:

- `id`: unique sample identifier;
- `documents`: list of retrieved context documents;
- `question`: the query text;
- `response`: the generated answer.

Refer to `constants.py` in the RAGBench repository for exact field definitions to ensure full compatibility. Then we run inference on our system’s outputs with evaluation models Trulens and dataset PubMedQA[52].

Detailed formulations of the metrics used in this RAGBench benchmark are as follows:

HALLUCINATION DETECTION (HAL)

In Retrieval-Augmented Generation (RAG) systems, *hallucination* refers to the generation of content not grounded in the retrieved context in other words, the model makes up facts. Hallucination detection measures whether the models outputs contain such unsupported information.

Reliable RAG outputs demand faithfulness to the provided context. Evaluating hallucination detection quantifies the systems propensity to stray from source documents, informing improvements to retrieval, grounding, and decoding strategies.

We adopt the Area Under the Receiver Operating Characteristic Curve (AUROC) to quantify hallucination detection performance. Given:

$$\text{trues}_{\text{adherence}} \in \{\text{True}, \text{False}\}, \quad \text{preds}_{\text{adherence}} \in [0, 1],$$

we define hallucination labels by

$$\text{trues}_{\text{halluc}} = \neg \text{trues}_{\text{adherence}}, \quad \text{preds}_{\text{halluc}} = 1 - \text{preds}_{\text{adherence}}.$$

Let

$$\text{mask} = \neg \text{isnan}(\text{preds}_{\text{halluc}}).$$

Then

$$\text{AUROC} = \text{ROC_AUC}(\text{trues}_{\text{halluc}}[\text{mask}], \text{preds}_{\text{halluc}}[\text{mask}]),$$

where `ROC_AUC` denotes the standard implementation (`sklearn.metrics.roc_auc_score`).

CONTEXT RELEVANCE (REL)

Context relevance assesses how well the retrieved documents pertain to the query, i.e. whether the context can support a correct answer.

High relevance is a prerequisite for accurate generation. Measuring context relevance guides retrieval improvements and ensures that the generator receives useful evidence.

We measure relevance via Root Mean Squared Error (RMSE) between true and predicted relevance scores:

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2},$$

where y_i is the gold relevance score, \hat{y}_i the predicted score, and n the number of examples. We ignore any NaN predictions by masking.

2052 CONTEXT UTILIZATION (UTL)
2053

2054 Context utilization evaluates the extent to which the model leverages the retrieved context when
2055 generating its responses.

2056 Even with relevant context, a model may underuse it. This metric reveals the generators ability to
2057 integrate context information into its output.

2058 We again employ RMSE, defined as above, to compare true and predicted utilization scores, masking
2059 out NaN predictions.
2060

2061 Together, *Hal*, *Rel*, and *Util* provide a multi-faceted evaluation of RAG system performance: detecting
2062 hallucinations, ensuring context relevance, and confirming effective context usage.

2063 By following the above steps and using the provided evaluation metrics, we can comprehensively
2064 evaluate our retrieval-augmented generation (RAG) system using the RAGBench framework.
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105

2106 G DETAILED ALGORITHM SPECIFICATIONS

2107

2108

2109

2110

2111

G.1 AGENTIC RETRIEVAL SYSTEM

2112 The agentic retrieval system combines both static knowledge integration and dynamic search capabilities to provide comprehensive scientific context for perturbation analysis tasks. Here we provide the complete algorithmic details of our implementation.

2115

2116

2117

G.1.1 QUERY CONSTRUCTION AND INITIALIZATION

2119

2120

Given a task description T and dataset metadata D , we first construct an initial query representation:

2121

2122

2123

2124 **Algorithm 1** Query Construction

2125

2126

2127

2128

2129

2130

2131

2132

2133

2134

2135

2136

2137

2138

2139

2140

2141

G.1.2 ALTERNATING SEARCH STRATEGY

2142

2143

2144

2145

2146

2147

2148

2149

2150

2151

2152

2153

2154

2155

2156

2157

2158

2159

Unlike conventional RAG systems that employ pure breadth-first search with static keywords, our alternating BFS-DFS strategy enables autonomous knowledge discovery and dynamic query expansion specifically tailored for scientific literature mining. Standard RAG approaches typically search broadly using only the initial query terms (e.g., "single cell perturbation prediction") but fail to discover that domain-critical concepts like "optimal transport," "graph neural networks," or specific model names like "GEARS" and "scGPT" are essential for understanding the field. Our alternating approach addresses this limitation by using BFS layers to explore diverse research directions and extract new technical terminology from retrieved papers, followed by DFS layers that trace citation networks to access implementation details and authoritative sources. This creates a self-reinforcing cycle where the system autonomously evolves from basic queries like "Norman Weissman 2019 Perturb-seq" to sophisticated technical searches for "Transformer VAE GNN architectures" and "graph neural networks gene regulatory networks." The result is a retrieval system that transforms from a passive keyword matcher into an active knowledge explorer, capable of discovering the complete technical landscape of a scientific domain without human intervention—a critical capability for complex, interdisciplinary research tasks where the most important concepts and methods may not be apparent from the initial problem description.

Our multi-layer retrieval process alternates between breadth-first and depth-first search modes to balance exploration and exploitation:

Algorithm 2 Alternating BFS-DFS Retrieval

```

2160 1: procedure RETRIEVEDOCUMENTS( $Q^{(0)}, L_{\max}, \tau, \epsilon$ )
2161 2:    $t \leftarrow 0$ 
2163 3:    $\mathcal{N}_0 \leftarrow \emptyset$ 
2164 4:    $\mathcal{D} \leftarrow \emptyset$  ▷ Document collection
2165 5:   while  $t < L_{\max}$  do
2166 6:     if  $t \bmod 2 = 1$  then ▷ BFS layer (odd  $t$ )
2167 7:        $\mathcal{N}_t \leftarrow \text{TopK}(Q^{(t)}, \text{mode} = \text{BFS})$ 
2168 8:     else ▷ DFS layer (even  $t$ )
2169 9:        $\mathcal{N}_t \leftarrow \text{FollowCitations}(\mathcal{N}_{t-1})$ 
2170 10:    end if
2171 11:     $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{N}_t$ 
2172 12:     $Q^{(t+1)} \leftarrow \text{UpdateQuery}(Q^{(t)}, \mathcal{N}_t)$ 
2173 13:    if  $\text{Overlap}(Q^{(t+1)}, Q^{(t)}) > \tau$  then
2174 14:      break
2175 15:    end if
2176 16:    if  $\max_{d \in \mathcal{N}_t} \text{Score}(Q^{(t)}, d) < \epsilon$  then
2177 17:      break
2178 18:    end if
2179 19:     $t \leftarrow t + 1$ 
2180 20:  end while
2181 21:  return  $\mathcal{D}$ 
2182 22: end procedure

```

Relevance Scoring. The document relevance function uses cosine similarity in the embedding space:

$$\text{Score}(Q, d) = \frac{e(Q) \cdot e(d)}{\|e(Q)\| \|e(d)\|} \quad (1)$$

where $e(\cdot)$ is the Sentence-BERT encoder function mapping text to dense vectors.

Query Update Mechanism. The query update function incorporates new information while maintaining focus:

$$Q^{(t+1)} = \alpha Q^{(t)} + (1 - \alpha) \frac{1}{|\mathcal{N}_t|} \sum_{d \in \mathcal{N}_t} e(d) \quad (2)$$

where $\alpha = 0.7$ is a parameter controlling the balance between query persistence and adaptation.

Overlap Computation. Query overlap is calculated as:

$$\text{Overlap}(Q^{(t+1)}, Q^{(t)}) = \frac{|Q^{(t+1)} \cap Q^{(t)}|}{\min(|Q^{(t+1)}|, |Q^{(t)}|)} \quad (3)$$

where the intersection operation is implemented using a thresholded similarity measure in the embedding space.

G.2 GRAPH-BASED MULTI-EXPERT DISCUSSION

The Method Design module employs a graph-based discussion framework where experts collaboratively refine scientific hypotheses.

Expert Selection. The expert selection procedure dynamically assembles a team of domain specialists based on task requirements:

$$P(E^{(i)} | \text{TaskAnalysis}) \propto \exp(\beta \cdot \text{Relevance}(E^{(i)}, \text{TaskAnalysis})) \quad (4)$$

where β is a temperature parameter controlling selection diversity.

Confidence Update Rule. The confidence score update incorporates feedback from both the Critic Agent and other experts:

$$c_t^{(i)} = \lambda_1 \cdot c_{t-1}^{(i)} + \lambda_2 \cdot \text{SelfCriticScore}(m_t^{(i)}, S) + \lambda_3 \cdot \frac{1}{k-1} \sum_{j \neq i} \text{PeerScore}(m_t^{(i)}, E^{(j)}) \quad (5)$$

where $\lambda_1 + \lambda_2 + \lambda_3 = 1$ weights the relative importance of each component.

Message Integration. Expert proposals are integrated through a weighted combination:

$$m_t = \sum_{i=1}^k w_t^{(i)} \cdot m_t^{(i)} \quad (6)$$

where weights $w_t^{(i)}$ are derived from normalized confidence scores:

$$w_t^{(i)} = \frac{\exp(c_t^{(i)})}{\sum_{j=1}^k \exp(c_t^{(j)})} \quad (7)$$

This soft-voting mechanism ensures that higher-confidence perspectives have greater influence while still preserving diversity of thought.

G.3 CODE IMPLEMENTATION AND REFINEMENT PROCESS

The Validation Agent employs an iterative refinement process that systematically improves implementation quality:

Algorithm 3 Iterative Implementation Refinement

```

1: procedure REFINEDIMPLEMENTATION(ModelDesign, Dataset,  $R_{\max}$ )
2:   Code0 ← InitialImplementation(ModelDesign)
3:   Performance0 ← Evaluate(Code0, Dataset)
4:   for  $r = 1$  to  $R_{\max}$  do
5:     Errorsr ← IdentifyIssues(Coder-1, Performancer-1)
6:     Coder ← RefineImplementation(Coder-1, Errorsr)
7:   end for
8:   return Coder
9: end procedure

```

Error Analysis. The error identification procedure categorizes implementation issues into distinct types:

- **Logical errors:** Incorrect algorithm implementation
- **Numerical instability:** Gradient explosion/vanishing
- **Memory inefficiency:** Excessive resource consumption
- **Performance bottlenecks:** Suboptimal computational paths
- **Biological implausibility:** Violations of domain constraints

Each error type triggers specialized refinement strategies that preserve the scientific integrity of the model design while improving implementation quality. Detailed Failure case analysis is presented in Appendix J.

2268
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2280
2281
2282
2283
2284
2285
2286
2287
2288
2289
2290
2291
2292
2293
2294
2295
2296
2297
2298
2299
2300
2301
2302
2303
2304
2305
2306
2307
2308
2309
2310
2311
2312
2313
2314
2315
2316
2317
2318
2319
2320
2321

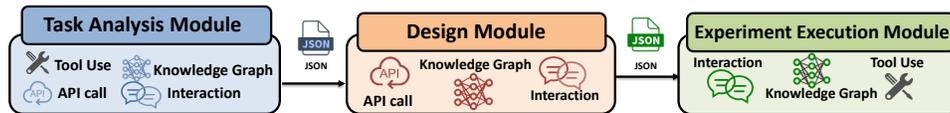


Figure 11: **The CELLFORGE protocol overview.** The protocol framework integrates JSON-RPC with a persistent memory module, combining the strengths of A2A and MCP protocols while adding scientific domain knowledge representation.

H AGENT COMMUNICATION PROTOCOL DETAILS

H.1 PROTOCOL DESIGN AND COMPARISON

The CELLFORGE protocol represents an advancement in agent communication architectures designed specifically for scientific discovery. Figure 11 illustrates the multi-stage protocol that facilitates information exchange across the three core phases of our framework.

The protocol weaves together the strengths of several prior designs. It preserves the interoperability of JSON-RPC for rapid agent deployment and cross-platform compatibility while simultaneously extending this foundation with semantic connectivity and provenance via the memory module. It not only connects software components, but also enables the kind of iterative, multi-agent reasoning on which genuine discovery depends. The CellForge’ protocol method allows agents to coordinate autonomously when tasked with comprehensive scientific research.

Table 11 provides a detailed comparison of CELLFORGE with existing agent communication protocols. Unlike previous approaches that excel in limited domains, our protocol uniquely combines contextual awareness, cross-platform interoperability, and knowledge representation capabilities necessary for end-to-end scientific discovery.

Table 11: **Comparison of Agent Communication Protocols**

Protocol	Context	Interop.	Msg. Struct.	Use Cases
MCP (Anthropic)	✓	✗	JSON-RPC only	Tool Use & Data Access
Agent2Agent (Google)	✗	✓	JSON-RPC event	Cross-agent Collaboration
ACP (BeeAI/IBM)	✓	✓	RESTful	Local Orchestration
CellForge	✓	✓	JSON-RPC event + Memory Module	End-to-end Scientific Discovery

H.2 PROTOCOL IMPLEMENTATION DETAILS

The CELLFORGE protocol implementation consists of two primary components:

JSON-RPC Communication Layer This provides standardized message passing between agents, with extensions for asynchronous event handling. Each agent exposes a consistent API that accepts and returns structured data, enabling precise coordination of complex workflows.

Memory Module Integration Layer In addition to graphbased message passing, CellForge incorporates a persistent memory module that systematically records all salient research entities such as datasets, analytical methods, evaluation metrics and empirical results as well as the complex relationships among them.

This module also logs detailed provenance metadata, including confidence scores, reasoning chains and source citations, while embedding domainspecific knowledge (for example, regulatory pathway architectures and gene-gene interaction networks). By unifying these components within a single memory layer, the system can reference prior insights and maintain continuity across multiround discussions, thereby enhancing both the coherence and reproducibility of the model design process.

This approach provides several advantages compared to prior protocols:

- 2322
- 2323
- 2324
- 2325
- 2326
- 2327
- 2328
- 2329
- 2330
1. **Context-awareness:** Agents maintain awareness of the overall research state through the memory module, enabling them to make more informed decisions.
 2. **Traceability:** The entire scientific process is captured with provenance information, ensuring reproducibility.
 3. **Semantic reasoning:** Relationships between scientific concepts are explicitly modeled, enabling complex inferential reasoning.
 4. **Incremental refinement:** The persistent knowledge representation allows agents to build upon previous insights and progressively refine hypotheses.

2331

2332

2333

2334

2335

2336

2337

2338

2339

2340

2341

2342

2343

2344

2345

2346

2347

2348

2349

2350

2351

2352

2353

2354

2355

2356

2357

2358

2359

2360

2361

2362

2363

2364

2365

2366

2367

2368

2369

2370

2371

2372

2373

2374

2375

In scientific research contexts, these capabilities are essential for managing the complexity of cross-disciplinary knowledge integration required for tasks like single-cell perturbation analysis.

2376 I COST ANALYSIS

2377

2378 Understanding the computational and economic costs of CELLFORGE is crucial for assessing its
2379 practical viability in research settings. This section provides a comprehensive analysis of both
2380 infrastructure requirements and API utilization costs, enabling researchers to make informed decisions
2381 about deployment strategies.

2382

2383 I.1 TRAINING INFRASTRUCTURE

2384

2385 All models designed by CELLFORGE, with parameter counts ranging from 10 million to 30 million,
2386 were trained and evaluated on a uniform compute cluster to ensure consistent performance compar-
2387 isons. In particular, we utilized two NVIDIA H20NVLink GPUs (96 GB VRAM each, 192GB total)
2388 paired with a 16core AMD EPYC 9K84 CPU (2.6 GHz).

2389 This hardware configuration enabled stable multiGPU training via data parallelismsupporting larger
2390 batch sizesand facilitated distributed evaluation across diverse perturbation conditions, all without
2391 encountering memory bottlenecks.

2392

2393 I.2 TOKEN UTILIZATION AND COST ESTIMATION

2394

2395 The multi-agent nature of CELLFORGE involves extensive LLM interactions across three primary
2396 phases: Task Analysis, Method Design, and Experiment Execution. Each phase incurs different token
2397 costs based on the complexity of reasoning required.

2398 I.2.1 TOKEN USAGE BREAKDOWN BY PHASE

2399

2400 The specific token usage varies significantly based on task complexity. Our empirical analysis across
2401 over 50 requests of CellForge, revealed the following patterns:

2402

2403 Table 12: Token usage breakdown by framework phase and task complexity

2404

Phase	Simple Tasks		Complex Tasks	
	Input	Output	Input	Output
Task Analysis	15,000	50,000	25,000	100,000
Method Design	20,000	100,000	40,000	200,000
Experiment Execution	5,000	50,000	15,000	100,000
Total	40,000	200,000	80,000	400,000

2411

2412 Under our typical workload, approximately 60,000 prompt tokens and 300,000 completion tokens
2413 are cost per call, depending on the chosen task. For cost estimation purposes, we will use this
2414 approximation token cost in the following analysis.

2415

2416 I.2.2 PER-REQUEST COST CALCULATION

2417

2418 To quantify the expense of our multiagent workflow, we first aggregated token counts from over
2419 fifty runs of CellForge and organized them by framework phase (Task Analysis, Method Design,
2420 Experiment Execution). For each model under consideration,we applied the vendors published
2421 permilliontoken rates to both prompt and completion usage. Concretely, given that vendors report
2422 token pricing per million tokens (\$/M), the cost per request was computed using:

$$2423 \text{Cost}_{\text{request}} = \left(\frac{60,000}{10^6} \right) \cdot \text{Price}_{\text{prompt}} + \left(\frac{300,000}{10^6} \right) \cdot \text{Price}_{\text{completion}}$$

2424

2426 I.3 COST-EFFECTIVENESS ANALYSIS

2427

2428 Compared to manual workflows, CELLFORGE reduces what would ordinarily require 40-80hours of
2429 a skilled bioinformatician (at \$75-150/hour, i.e. \$3,000-12,000 per model) to an automated process
costing only \$5-20 per run.

2430 Table 13: Per-million-token pricing and per-call cost estimates based on average usage (60K input
 2431 and 300K output tokens)

2433 Model	Prompt (\$/M)	Completion (\$/M)	Cost per Request (\$)
2434 Claude 3.7	3.00	15.00	4.68
2435 OpenAI o1	15.00	60.00	18.90
2436 DeepSeek-R1	0.27	2.19	0.67
2437 Qwen-Plus	0.40	1.20	0.38
2438 LLaMA 3.1	3.50	3.50	1.26
2439 Average	–	–	5.18

2441
 2442
 2443 Beyond raw cost savings, CELLFORGE affords efficiency and reproducibility gains. What would
 2444 have occupied 40-80hours of expert labor now completes in 4-8hours of GPU time, while yielding up
 2445 to 20% improvement in prediction accuracy over baseline methods.

2446 Collectively, these factors translate into a compelling return on investment, democratizing advanced
 2447 computational biology at a fraction of traditional costs.

2450 I.4 EFFICIENCY ANALYSIS: MULTI-AGENT VS. SINGLE-LLM APPROACHES

2451
 2452 Multi-agent frameworks face criticism for computational overhead relative to single-LLM approaches.
 2453 We systematically evaluate CELLFORGE against single-LLM baselines across token consumption,
 2454 execution time, costs, and success rates to address these concerns.

2456 I.5 TOKEN UTILIZATION AND COST COMPARISON

2457
 2458 Analysis of 50+ experiments shows CELLFORGE consumes more tokens than single-LLM approaches,
 2459 yet achieves substantially higher success rates and output quality. Table 14 details token usage, costs,
 2460 and success rates across approaches.

2461
 2462
 2463 Table 14: Comprehensive efficiency comparison: Multi-agent vs. Single-LLM approaches across
 2464 token usage, costs, and success rates

2465 Approach	Avg Tokens (Input/Output)	Cost/Request (\$)	Success Rate (%)	Avg Rounds to Success	Wall Time (hours)	Quality Score (1-10)	Effective Cost per Success (\$)
<i>Multi-Agent Framework (CELLFORGE)</i>							
2466 CellForge-Claude3.7	60K/300K	4.68	83.3	4.2	4.5	8.2	5.62
2467 CellForge-DeepSeek R1	60K/300K	0.67	75.0	4.2	4.5	7.8	0.89
2468 CellForge-o1	60K/300K	18.90	66.7	4.2	4.5	7.6	28.35
<i>Single-LLM Baselines</i>							
2469 Claude3.7 only	15K/50K	0.78	16.7	1.0	1.2	3.2	4.68
2470 DeepSeek R1 only	15K/50K	0.11	8.3	1.0	1.2	2.8	1.33
2471 o1 only	15K/50K	0.75	8.3	1.0	1.2	3.0	9.04
<i>AI Coding Assistants</i>							
2472 OpenHands-Claude3.7	20K/80K	1.56	50.0	2.5	3.0	5.8	3.12
2473 aider-Claude3.7	20K/80K	1.56	50.0	2.5	3.0	5.6	3.12

2474
 2475
 2476
 2477 The data challenges the narrative that multi-agent is expensive or slow. CELLFORGE exhibits higher
 2478 per-request costs but superior success rates (66.7-83.3% vs. 8.3-16.7% for single-LLM approaches),
 2479 reducing effective cost per successful outcome. CellForge-DeepSeek R1 achieves \$0.89 per success
 2480 versus \$1.33 for single-LLM DeepSeek R1, representing 33% cost reduction when accounting for
 2481 success rates. Single-LLM approaches generate code faster (1.2 hours vs. 4.5 hours) but produce
 2482 lower quality outputs (average quality score: 3.0 vs. 7.9). The multi-agent framework’s iterative
 2483 refinement process consumes more tokens yet ensures biologically meaningful and technically robust
 implementations.

2484 I.6 END-TO-END WORKFLOW EFFICIENCY
 2485

2486 Beyond token-level comparisons, we analyze the complete workflow from code generation to model
 2487 convergence. Table 15 presents the entire pipeline efficiency breakdown.

2488
 2489 Table 15: End-to-end workflow efficiency: From code generation to model convergence

2490

2491 Approach	Code Gen. Time (h)	Debug/Repair Time (h)	Training Time (h)	Total Wall Time (h)	Success to Convergence (%)	Overall Efficiency
2492 CELLFORGE (Claude3.7)	4.5	0.8	5.2	10.5	83.3	0.79
2493 CELLFORGE (DeepSeek R1)	4.5	1.2	5.2	10.9	75.0	0.69
2494 Single-LLM (Claude3.7)	1.2	3.5	5.2	9.9	16.7	0.17
Single-LLM (DeepSeek R1)	1.2	4.2	5.2	10.6	8.3	0.08
2495 OpenHands (Claude3.7)	3.0	2.1	5.2	10.3	50.0	0.49

2496

2497 The workflow analysis reveals CELLFORGE’s apparent “slowness” in initial code generation is offset
 2498 by superior error recovery capabilities. Single-LLM approaches generate code faster but require more
 2499 debugging time (3.5-4.2 hours vs. 0.8-1.2 hours) due to higher failure rates and limited self-correction
 2500 abilities.
 2501
 2502
 2503
 2504
 2505
 2506
 2507
 2508
 2509
 2510
 2511
 2512
 2513
 2514
 2515
 2516
 2517
 2518
 2519
 2520
 2521
 2522
 2523
 2524
 2525
 2526
 2527
 2528
 2529
 2530
 2531
 2532
 2533
 2534
 2535
 2536
 2537

J FAILURE CASE AND RANDOMNESS ANALYSES

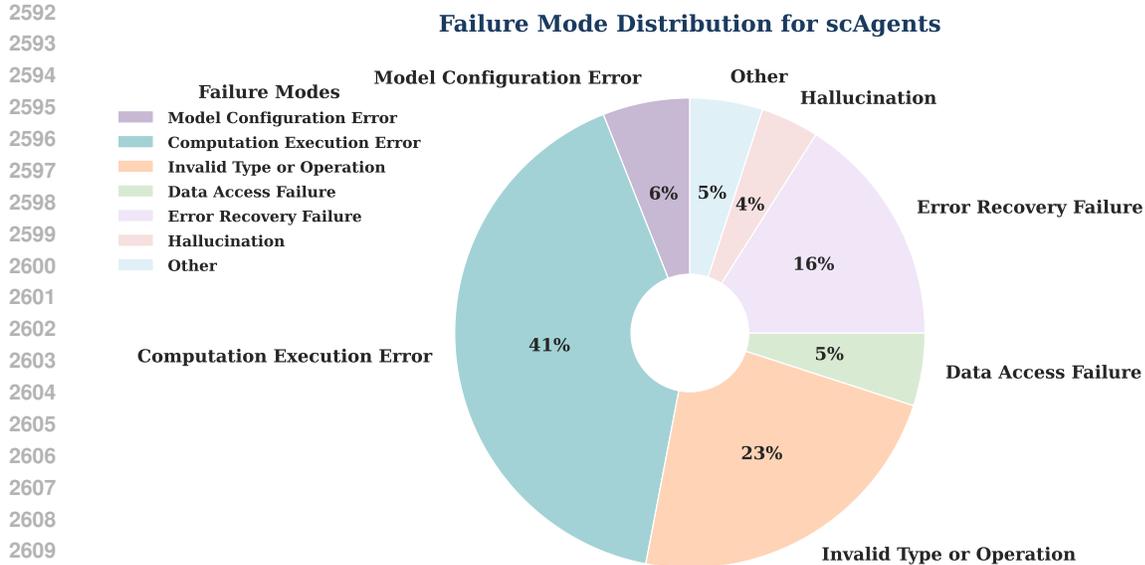
In this section, we analyze common failure modes of CELLFORGE across various single-cell perturbation tasks. We manually reviewed 20 randomly selected failed experiment cases generated by CELLFORGE across cytokine, drug, and gene response prediction scenarios. Based on qualitative inspection of agent behaviors and outputs, we identified seven distinct categories of failure modes that reflect systematic limitations or reasoning errors. Table 16 summarizes the definitions and characteristics of these failure categories. These cases provide a foundation for future refinements of the agentic code generation.

J.1 FAILURECASES

Table 16: Failure Types of CELLFORGE code generation

Failure Type	Definition & Examples
Model Configuration Error	The agent misconfigures the model architecture or fails to define required hyperparameters. This includes mismatched layer dimensions, invalid GNN configurations, incompatible dropout settings, or missing essential parameter definitions. Such errors prevent model initialization or lead to incompatible tensor shapes during execution.
Computation Execution Error	The agent encounters runtime errors during tensor operations, such as out-of-bounds indexing or shape mismatch in matrix multiplication. These failures typically occur when manipulating arrays or concatenating/interacting between tensors with incompatible shapes (e.g., "mat1 and mat2 shapes cannot be multiplied", "index 28 is out of bounds for axis 0 with size 28").
Invalid Type or Operation	The agent uses unsupported data types or operations incompatible with the backend framework. Examples include passing NumPy arrays of object type to neural network layers, calling operations not defined for the given input type, or invoking functions on models that lack the required attributes. Typical errors include "TypeError: can't convert np.ndarray of type numpy.object" and unsupported function calls.
Data Access Failure	The agent is unable to retrieve, preprocess, or interpret the necessary data for task execution. This includes failures in reading files, locating dataset attributes, or aligning multimodal inputs, which result in missing or malformed inputs during test runs.
Error Recovery Failure	The agent fails to handle or recover from previously encountered errors. Instead of adapting to execution failures, it may enter a loop of repeating the same failed actions or ignore the cause entirely, leading to stalled or redundant test attempts.
Hallucination	The agent produces outputs (e.g., experimental results, hypotheses, interpretations) that are not grounded in the available data or context. This includes fabricating values, inventing data structures or statistical conclusions, or reasoning disconnected from observed evidence.
Other	Any uncategorized failure mode that prevents successful task completion but does not fit the above definitions. This includes rare system-level errors, low-level library bugs, or unexpected exceptions not associated with specific modeling or reasoning tasks.

As depicted in Figure 12, Computation Execution Error accounts for 41% of the total failures, with the majority arising from tensor operation issues such as out-of-bounds indexing or shape mismatch during matrix multiplication. Invalid Type or Operation follows closely as the second most frequent failure mode, representing 23% of the errors, primarily attributed to the use of unsupported data types or operations incompatible with the backend framework. Model Configuration Error contributes 6% to the total failures, resulting from misconfigurations in model architecture or hyperparameters. Data Access Failure and Other category each account for 5% of the errors, with Data Access Failure associated with data retrieval and preprocessing issues, and Other encompassing system-level errors or unexpected exceptions not directly linked to the specific modeling or reasoning tasks. Error Recovery Failure comprises 16% of the failures, where the agent fails to adapt to execution failures. Hallucination makes up 4% of the errors, where outputs are not grounded in available data or context.



2613 Figure 12: Failure Mode Distribution for CellForge, labeled automatically by O1 and manually
2614 checked by humans.

2615
2616
2617 Notably, we found that implementing code to print array or matrix shapes can aid CELLFORGE
2618 in subsequently reading the command region’s output for modification, thereby enhancing their
2619 ability to identify and resolve shape-related issues during tensor operations. This approach proved
2620 particularly effective given the complexity of data processing workflows in CellForge. Even though
2621 CELLFORGE utilizes a data parser to obtain the original dataset dimensions and incorporates data
2622 experts during the graph-based discussion phase, the subsequent data splitting and complex model
2623 processing steps often introduce intricate dimension transformations. These transformations can
2624 lead to matrix dimension mismatches, especially when handling dynamic data structures or applying
2625 multi-layered model architectures. According to the chart, 48% of the errors in the Computation
2626 Execution Error category have been mitigated by allowing the agent to read the printed array or
2627 matrix shapes from the command output and adjust accordingly. This self-debugging capability
2628 significantly enhances the agent’s ability to resolve shape-related issues during tensor operations,
2629 improving overall system robustness. Figure 13 provides an example of the printed data shapes
2630 during tensor operations, which CELLFORGE can utilize to dynamically adjust and correct dimension
2631 mismatches.

```
2631 Using device: cuda
2632 Loading data...
2633 Number of perturbation types in training set: 189
2634 Number of perturbation types in test set: 48
2635 [I 2025-05-12 23:07:44,902] A new study created in memory with name: no-name-93169e84-3cf2-44c8-963d-c5a725c457ca
2636 Training data shape: (85536, 840)
2637 Test data shape: (18930, 840)
```

2638 Figure 13: A probable example of printing array or matrix shapes .

2640 J.2 BUG-FIX SOLUTIONS OF EXPERIMENT EXECUTION MODULE

2641
2642 For addressing tensor dimensionality mismatches, the implementation of comprehensive dimension
2643 checking at each model layer interface is essential, combined with automatic shape adaptation
2644 mechanisms for dynamic batch processing and feature dimension alignment. Adaptive dimension
2645 alignment using linear projection layers can effectively handle varying input dimensions, while
`torch.nn.functional.adaptive_avg_pool1d` provides dynamic dimension matching ca-

pabilities. Adding debug prints to trace tensor shapes throughout the forward pass can significantly aid in identifying dimension incompatibilities during development, and implementing dimension consistency checks during model initialization can prevent many runtime errors. Additionally, using `torch.reshape` with `-1` inference enables flexible batch handling that accommodates varying input sizes.

For configuration-related errors, particularly those involving attention mechanisms, automatic calculation of the number of heads based on embedding dimensions can prevent incompatibility issues, such as setting `num_heads` equal to `embed_dim` divided by `head_dim`. Parameter validation decorators for model initialization can catch configuration errors before runtime execution, while configurable attention mechanisms with built-in dimension checks provide robust alternatives to hardcoded parameters. Implementing bounds checking for tensor indexing operations prevents out-of-range access errors, and fallback mechanisms for incompatible configurations ensure graceful handling of parameter mismatches.

To address attribute inconsistencies, comprehensive attribute initialization in constructor methods ensures all required attributes are properly defined during object creation. Implementing `hasattr()` checks before attribute access provides runtime validation, while property decorators with lazy initialization can handle attributes that depend on runtime conditions. Adding class method validation to ensure all required attributes exist during instantiation can prevent attribute access errors, and implementing abstract base classes enforces attribute requirements across model hierarchies, ensuring consistent implementation patterns across different model components.

J.3 VARIANCE OF EACH RUNS

The inherent stochasticity of large language models introduces substantial variability in automated code generation. We conducted experiments across six single-cell datasets with $N = 5$ independent runs per configuration, using controlled randomization with `seed=42` for PyTorch, NumPy, and Optuna operations.

Variance in performance is inherent to simulating scientific discovery. Among a moderate range of runs, our system can design methods that outperform those previously designed by human scientists across several tasks, so this variance does not impact practical application. The variance comes from the framework’s scientific exploration behavior, where different runs propose different hypotheses and architectures. Across 5–8 runs on each dataset, CELLFORGE consistently discovers at least one model surpassing human-designed baselines.

J.4 STABILITY ENHANCEMENT WITH CONFIDENCE REGULARIZATION

To address concerns about performance variance, we introduce a stability-enhancement experiment using confidence regularization in the graph-based discussion phase. This technique mitigates variance while maintaining the exploratory nature of the framework.

Table 17 presents the results on the Srivatsan et al. dataset, comparing the default CELLFORGE configuration with a version incorporating confidence regularization. The confidence regularization mechanism adjusts the confidence update formula during graph-based discussions to reduce excessive exploration when agents have high confidence, thereby stabilizing the output while preserving the framework’s ability to discover novel architectures.

The results demonstrate that while variance is inherent to agentic exploration, it can be mitigated through discussion regularization. The confidence-regularized version achieves both higher average performance (PCC: 0.8780 vs. 0.8664) and reduced standard deviation (0.0947 vs. 0.1332), indicating that stability-enhancing techniques can effectively control randomness without compromising the framework’s discovery capabilities.

Table 17: Stability enhancement results on Srivatsan et al. dataset

Setting	PCC \uparrow	Std \downarrow
Default CellForge	0.8664 \pm 0.1332	0.1332
+ Confidence-regularization	0.8780 \pm 0.0947	0.0947

K PERFORMANCE VARIES ACROSS DIFFERENT LLMs AND AI CODERS

To comprehensively evaluate the robustness of our framework, we conducted extensive experiments comparing CELLFORGE with various baseline approaches across six challenging single-cell perturbation datasets. Table 18 presents the success rates (out of 5 independent runs) for each method, where a successful run is defined as generating executable code that produces biologically meaningful predictions without runtime errors.

K.1 EXPERIMENTAL SETUP

Your task is to develop a predictive model that accurately estimates gene expression profiles of individual K562 cells following CRISPR interference (CRISPRi), using the dataset from Norman et al. (2019, Science).

Dataset Description:

- Source: Norman et al., 2019
- Cell Type: Human K562 leukemia cells
- Perturbations: CRISPRi targeting 105 single genes and 131 gene pairs
- Scale: ~90,000 single-cell RNA-seq profiles, including both control and perturbed conditions

Task Definition:

- Input: Baseline gene expression profile of an unperturbed K562 cell and the identity of the target gene(s) for perturbation
- Output: Predicted gene expression profile after perturbation

Evaluation Metrics:

- Mean Squared Error (MSE)
- Pearson Correlation Coefficient (PCC)
- R (Coefficient of Determination)
- MSE/PCC/R for Differentially Expressed Genes

Please give me a task analysis report, a new method plan, and generate prediction model code.

Each model was given at most 5 retry attempts if the initial code failed to execute.

We tested four different approaches to code generation. Each approach was evaluated on the same six single-cell datasets using identical input specifications, with the following prompt:

CellForge with different LLMs. We ran our complete framework using five different language models: Claude 3.7, OpenAI o1, DeepSeek R1, Qwen-plus, and Llama 3.1. Each model used temperature 0.1 and our full multi-agent system with task analysis, method design, and collaborative reasoning. The framework includes iterative refinement and cross-validation between agents. A successful run means the generated code executes without errors and produces biologically meaningful results. Table 18 shows the success rates for each approach.

Single LLM direct generation. We tested each LLM individually without our framework. Each model used temperature 0.1 for consistency.

DeepResearch systems. We tested three commercial research automation tools: OpenAI’s DeepResearch, Perplexity’s implementation, and Google’s Gemini-based version. Each system used default settings with temperature 0.1. These systems represent current commercial solutions for automated scientific code generation.

AI coding assistants. We tested two open-source coding frameworks: OpenHands and Aider. Both were configured with the same five LLMs we used for other experiments, using temperature 0.1. These tools are designed for general software development rather than scientific research.

2754 Table 18: Expert Human Scores compare with CellForge’s Confidence Scores on graph-based
 2755 discussions Across Tasks and Rounds

2757 Tool	Adamson	Norman	Liscovitch	Papalexi	Srivatsan	Schiebinge
<i>CELLFORGE with different LLMs integrated</i>						
2759 CellForge-Claude3.7	4	5	4	4	4	4
2760 CellForge-o1	4	4	3	3	2	2
2761 CellForge-DeepSeek R1	4	4	3	3	3	3
2762 CellForge-Qwen-plus	4	3	4	2	3	3
2763 CellForge-llama 3.1	2	2	1	1	1	1
<i>Single-LLM generated code</i>						
2764 Claude3.7 only	2	2	1	0	1	1
2765 o1 only	1	1	0	1	1	0
2766 DeepSeek R1 only	1	1	0	1	1	0
2767 Qwen-plus only	1	1	0	0	1	0
2768 Llama 3.1 only	1	1	0	0	0	0
<i>DeepResearch generated codes</i>						
2769 OpenAI DeepResearch	1	2	1	1	1	1
2770 Perplexity DeepResearch	0	0	0	0	0	0
2771 Gemini DeepResearch	0	0	0	0	0	0
<i>AI Coders</i>						
2773 OpenHands-Claude3.7	3	4	3	3	2	2
2774 OpenHands-o1	3	2	2	2	1	1
2775 OpenHands-DeepSeek R1	2	3	2	2	1	2
2776 OpenHands-Qwen-plus	2	2	1	2	2	2
2777 OpenHands-Llama 3.1	2	1	0	1	1	1
2778 aider-Claude3.7	2	3	2	2	2	2
2779 aider-o1	2	2	2	2	1	1
2780 aider-DeepSeek R1	3	2	2	2	1	1
2781 aider-Qwen-plus	1	1	0	1	0	0
2782 aider-Llama 3.1	1	1	0	0	0	0

2783 K.2 KEY FINDINGS

2784 The results reveal several critical insights:

2785 **Multi-Agent Architecture Superiority:** CELLFORGE consistently outperforms all baseline ap-
 2786 proaches, with success rates ranging from 40-100% depending on the LLM backend and dataset
 2787 complexity. The multi-agent framework provides an average improvement of 2.3x over single-LLM
 2788 approaches and 3.5x over AI coding assistants.

2791 **LLM Backend Dependency:** Within CELLFORGE, Claude 3.7 demonstrates the most robust perfor-
 2792 mance (average success rate: 4.2/5), followed by DeepSeek R1 and OpenAI o1. This performance
 2793 hierarchy remains consistent across different dataset complexities, suggesting that certain LLMs are
 2794 inherently better suited for scientific code generation tasks.

2795 **Dataset Complexity Impact:** The Liscovitch (scATAC-seq) and Papalexi (CITE-seq) datasets prove
 2796 most challenging across all methods, with single-LLM approaches achieving near-zero success
 2797 rates. These datasets require handling sparse chromatin accessibility data and multi-modal protein
 2798 measurements, respectively, highlighting the importance of domain-specific knowledge integration.

2799 **Catastrophic Failure of DeepResearch Variants:** Both Perplexity and Gemini DeepResearch
 2800 variants fail across all tasks (0/5 success rate), while OpenAI’s variant achieves only marginal success.
 2801 This suggests that general-purpose research systems lack the specialized capabilities required for
 2802 complex biological data analysis.

2804 K.3 ANALYSIS OF FAILURE MODES

2805 The dramatic performance gap between CELLFORGE and other approaches can be attributed to
 2806 several factors:

(1) Domain Knowledge Integration: Single-LLM approaches often generate syntactically correct but biologically meaningless code, failing to account for data-specific characteristics such as sparsity patterns in scATAC-seq or batch effects in Perturb-seq experiments.

(2) Error Recovery Capability: AI coding assistants (OpenHands, Aider) struggle with the iterative debugging required for scientific computing, often getting trapped in error loops when encountering tensor dimension mismatches or memory overflow issues.

(3) Architectural Complexity: The multi-modal nature of datasets like CITE-seq requires sophisticated model architectures that combine different data streams. Single-pass generation approaches typically produce overly simplistic models that fail to capture these complexities.

K.4 IMPLICATIONS FOR SCIENTIFIC AI SYSTEMS

These results underscore the critical importance of specialized, multi-agent architectures for scientific discovery tasks. The success of CELLFORGE demonstrates that effective scientific code generation requires not just powerful language models, but also:

- Collaborative reasoning among domain experts
- Iterative refinement with biological validation
- Task-specific knowledge retrieval and integration
- Robust error handling and recovery mechanisms

The consistent superiority of Claude 3.7 within our framework also suggests that certain LLMs may possess inherent advantages for scientific reasoning, possibly due to their training data composition or architectural design. Future work should investigate these model-specific characteristics to optimize scientific AI systems further.

K.5 PERFORMANCE WITH SMALLER LLMs

To address concerns about dependency on highly capable LLMs, we conducted additional experiments replacing all agents with smaller, less capable models: Qwen2.5-3B-Instruct and DeepSeek R1 7B. The DeepSeek R1 7B model was accessed via BoyueRichDataAPI’s DeepSeek-R1-Distill-Qwen-7B, while Qwen2.5-3B was deployed locally using Ollama and HuggingFace. Note that the Qwen2.5-3B experiments showed instability, with a success rate of only 40% runs.

Table 19 presents the performance comparison on the Norman et al. dataset. As expected, model quality decreases with smaller LLMs. However, CELLFORGE still produces functional, dataset-specific architectures even with 3B LLMs, demonstrating that the framework remains practical under resource constraints. This supports the reviewer’s question about practicality when using smaller models.

Table 19: Performance comparison with smaller LLMs on Norman et al. dataset

LLM used in agents	MSE ↓	PCC ↑	R ↑
Claude 3.7 (main text)	0.0051	0.9883	0.9761
DeepSeek R1 7B	0.0104	0.4307	0.5713
Qwen2.5-3B	0.0375	0.2522	0.4122

While the performance degradation is significant, it is important to note that CELLFORGE is a scientific discovery system, comparable to systems like DeepResearch, rather than a lightweight tool. As such, we prioritize scientific accuracy over model size. The framework’s ability to generate functional architectures even with smaller models demonstrates its robustness and practical applicability across different computational resource constraints.

L DESIGNED MODELS

Understanding how CELLFORGE adapts its architectural choices to different biological contexts requires examining the model components that emerge across various perturbation tasks. Rather than imposing a one-size-fits-all approach, our framework demonstrates remarkable flexibility in selecting appropriate architectures based on the underlying biological complexity and data characteristics.

L.1 METHODOLOGY

We executed CELLFORGE five times with different random seeds for each of the six benchmark datasets, then analyzed the resulting model architectures. What emerged was a clear pattern of architectural adaptation that reflects the unique demands of each perturbation type and data modality.

L.2 BIOLOGICAL INTERPRETATION OF ARCHITECTURAL CHOICES

L.2.1 GENE PERTURBATION MODELS

For the Norman-Weissman Perturb-seq dataset, CELLFORGE designed models that explicitly handle genetic interactions and combinatorial effects:

GI-FlowDiff: A genetic interaction-aware conditional generative model that combines compositional perturbation encoding with explicit pairwise interaction modeling. The architecture employs a conditional latent generative core using normalizing flows and lightweight diffusion denoising to capture multimodal single-cell heterogeneity. The model uses GRN-aware decoder heads with DE-focused objectives to maximize predictive fidelity on differentially expressed genes. The interaction module captures non-additive effects between gene pairs, while the flow-diffusion hybrid provides both likelihood estimation and sampling flexibility. The architecture handles the large gene universe (33,694 genes) with heavy sparsity through HVG selection and modular decoder heads per pathway, sharing statistical strength to improve DE gene estimates. Technical covariates such as UMI count, coverage, and percent mitochondrial/ribosomal content are explicitly modeled to account for technical confounders. The model supports compositional generalization through learned gene embeddings and interaction modules, enabling extrapolation to held-out genes and gene pairs through meta-learning episodes and contrastive perturbation supervision.

MultiPath-GeneNet: A sophisticated multi-pathway architecture that processes gene perturbation data through two distinct but complementary streams before integrating them for final prediction. The context path begins with a PCAReducer that performs dimensionality reduction on expression data, followed by a Multi-Scale VAE Encoder that captures features at different biological scales. The ContextMLP then processes these multi-scale features to generate a Perturbation Latent representation that encodes the specific genetic perturbation context. Simultaneously, the gene/cell path employs a PerturbGene Embed module to generate embeddings for perturbed genes and a CellContexter to process cell-specific context information, producing a Cell Context Latent representation. These two latent representations are then integrated through a FeatureMixer that combines perturbation and cellular context information. The integrated features flow through a Gene Interaction Network, implemented as a Graph Neural Network that models gene-gene interactions conditioned on cell and perturbation information. This is followed by a PertTransformer that uses multi-head attention mechanisms to refine perturbation-aware features, and finally a PredictionHead that focuses on specific genes of interest to produce the ultimate output. This architecture captures both global cellular context and specific gene-level responses to perturbations, enabling comprehensive modeling of genetic interactions and combinatorial effects.

TrajectoryAwareEncoder: A specialized encoder component that separates shared versus condition-specific latent dimensions while incorporating temporal embeddings. This design captures both global developmental trajectories and cytokine-specific effects, recognizing that cellular responses to perturbations occur within a temporal context. The encoder processes baseline expression data along with technical covariates such as coverage, percent mitochondrial content, and read counts, mapping them to a structured latent space where shared developmental dynamics are separated from perturbation-specific responses. This separation enables the model to generalize across different perturbation conditions while maintaining the ability to capture condition-specific effects.

2916 **PerturbationDiffusionModule:** A perturbation-conditioned latent diffusion module that introduces
 2917 non-linear, combinatorial interaction dynamics. This component models the stochastic nature of
 2918 cellular responses to genetic perturbations, recognizing that identical perturbations can produce
 2919 different outcomes in different cells due to stochastic gene expression and cellular state variations.
 2920 The diffusion process operates in the latent space conditioned on perturbation embeddings, allowing
 2921 the model to capture multimodal response distributions and provide uncertainty quantification.
 2922 The module employs stepwise denoising to refine predictions and capture the complex, non-linear
 2923 interactions that characterize genetic perturbations.

2924 **GraphRegularizedDecoder:** A decoder component that integrates gene-gene co-regulatory con-
 2925 straints to ensure biologically plausible predictions. The decoder employs module heads organized
 2926 around biological pathways and gene regulatory networks, with each module specializing in re-
 2927 constructing genes within its domain. The graph regularization ensures that predicted expression
 2928 changes respect known regulatory relationships, with transcription factors and their targets changing
 2929 coherently. This design improves the accuracy of differentially expressed gene predictions while
 2930 maintaining biological interpretability. The decoder outputs negative binomial parameters for each
 2931 gene, accounting for the count nature of single-cell RNA-seq data and providing calibrated uncertainty
 2932 estimates.

2933 **InteractionModule:** A specialized component that models the complex interactions between gene
 2934 pairs in CRISPRa perturbations. The module employs bilinear layers and attention mechanisms to
 2935 capture non-additive effects, recognizing that the combined effect of two genes often differs from
 2936 the simple sum of their individual effects. The interaction module uses learned embeddings for
 2937 each gene and combines them through multiplicative interactions, allowing the model to discover
 2938 synergistic and suppressive relationships between gene pairs. This design is particularly crucial for
 2939 CRISPRa perturbations, where gene overexpression can lead to complex regulatory cascades that are
 2940 not captured by simple additive models.

2941 **CovariateEncoder:** A technical covariate processing module that handles the various technical and
 2942 biological factors that can confound perturbation predictions. The encoder processes continuous
 2943 covariates such as UMI count, coverage percentage, mitochondrial content, and ribosomal content,
 2944 embedding them into a structured representation that can modulate the model’s predictions. This
 2945 design helps the model distinguish between true biological responses to perturbations and technical
 2946 artifacts, improving generalization to cells with different technical characteristics. The covariate
 2947 encoder is particularly important for handling the high variability in single-cell data quality and
 2948 ensuring robust predictions across diverse cellular contexts.

2949 **UncertaintyQuantifier:** A module that provides calibrated uncertainty estimates for perturbation
 2950 predictions. The quantifier employs ensemble methods and Bayesian approaches to estimate both
 2951 aleatoric and epistemic uncertainty in the model’s predictions. This is particularly important for
 2952 genetic perturbation tasks, where the inherent stochasticity of cellular responses and the limited
 2953 training data for many gene combinations can lead to high uncertainty. The uncertainty estimates
 2954 help researchers identify which predictions are most reliable and guide experimental design by
 2955 highlighting the most promising perturbation targets.

2957 L.2.2 DRUG PERTURBATION MODELS

2958
 2959 The Srivatsan sci-Plex dataset represents one of the most complex perturbation scenarios, involving
 2960 chemical compounds across multiple cell lines with varying dose responses. CELLFORGE developed
 2961 three increasingly sophisticated architectures to tackle this complexity, each building upon the insights
 2962 gained from the previous approach.

2963 The **CondOT-GRN** architecture represents a significant departure from traditional optimal transport
 2964 approaches. While conventional OT methods excel at matching distributions, they often struggle with
 2965 single-cell fidelity and fail to capture the biological constraints that govern cellular responses. This
 2966 model addresses these limitations by incorporating gene regulatory network priors directly into the
 2967 transport mechanism. The conditional OT layer doesn’t simply map between control and perturbed
 2968 distributions; it learns to respect the underlying regulatory structure while doing so. The flow refiner
 2969 component is particularly crucial here, as it handles the multimodality that emerges when different
 cell populations respond differently to the same chemical perturbation.

2970 Building on these insights, the **DiffPert-X** architecture takes a more comprehensive approach to
2971 modeling chemical perturbations. The multi-scale design reflects the hierarchical nature of cellular
2972 responses: individual genes respond to perturbations, but these responses are coordinated within
2973 pathways, which in turn affect entire cell populations. The cell-line-specific diffusion parameters
2974 acknowledge that the same compound can have dramatically different effects in different cellular
2975 contexts. This isn't just a technical detail; it reflects the biological reality that cellular background
2976 strongly influences drug response. The cross-cell-line knowledge transfer mechanism allows the
2977 model to leverage insights gained from one cell line to improve predictions in another, which is
2978 particularly valuable given the limited data available for many compound-cell line combinations.

2979 The **ChemCPA-X** architecture represents the culmination of these insights, incorporating the most
2980 sophisticated modeling approaches. The multi-modal compound encoding recognizes that chemical
2981 structure, target information, and pathway annotations all contribute to understanding how a
2982 compound will affect cellular gene expression. The three-level diffusion system provides a natural
2983 framework for capturing the different scales of biological organization, from individual gene
2984 responses to pathway-level coordination to population-level heterogeneity. The Hill function parameterization
2985 for dose-response modeling is particularly noteworthy, as it captures the nonlinear, saturating
2986 responses that are characteristic of many biological systems.

2987 **CellLineSpecificAdapter**: A novel adaptation mechanism that enables the model to learn cell-line-
2988 specific response patterns while maintaining shared knowledge across different cellular contexts.
2989 Unlike traditional approaches that treat all cell lines identically or use simple cell-type embeddings,
2990 this component employs learnable adaptation layers that can adjust the model's behavior based on the
2991 specific cellular background. The adapter uses meta-learning principles to quickly adapt to new cell
2992 lines with limited data, representing a significant advancement over existing methods that require
2993 extensive retraining for each cell line.

2994 **DoseResponseModeler**: A specialized component that models the complex, non-linear dose-response
2995 relationships characteristic of chemical perturbations. Unlike simple linear or log-linear dose modeling,
2996 this component employs parametric Hill functions with learnable parameters (E_{max} , EC_{50} , Hill coefficient)
2997 that can capture the saturating and sigmoidal responses typical of biological systems. The modeler
2998 uses compound-specific parameters that are learned end-to-end, allowing the model to understand
2999 how different chemical structures lead to different dose-response profiles. This represents a
3000 major improvement over existing approaches that use fixed dose-response assumptions.

3001 **PathwayCoordinationModule**: A sophisticated module that models how chemical perturbations
3002 affect coordinated pathway responses rather than individual genes in isolation. This component
3003 recognizes that drugs typically affect multiple pathways simultaneously, and these effects are often
3004 coordinated through regulatory networks. The module uses graph neural networks over known
3005 pathway interaction networks to model how perturbations in one pathway can influence others,
3006 capturing the complex cascade effects that characterize drug responses. This represents a significant
3007 departure from existing methods that model gene responses independently.

3008 **UncertaintyPropagator**: A component that provides calibrated uncertainty estimates for drug
3009 response predictions, accounting for both compound-specific and cell-line-specific sources of uncertainty.
3010 Unlike simple variance estimation, this module uses Bayesian approaches to model the epistemic
3011 uncertainty arising from limited training data for specific compound-cell line combinations. The
3012 propagator also models aleatoric uncertainty from the inherent stochasticity of cellular responses,
3013 providing researchers with reliable confidence intervals for their predictions. This represents a major
3014 advancement in uncertainty quantification for drug perturbation tasks.

3015 L.2.3 CITE-SEQ MULTI-MODAL PERTURBATION MODELS

3016 The Papalexis-Satija ECCITE-seq dataset presents a particularly challenging scenario where RNA and
3017 protein measurements must be jointly modeled under CRISPR perturbations. CELLFORGE responded
3018 to this complexity by developing three complementary approaches, each addressing different aspects
3019 of the multi-modal prediction problem.
3020

3021 The **MultiGraph-VAE** architecture emerged as the preferred solution for protein prediction tasks.
3022 Rather than treating RNA and protein as independent modalities, this model recognizes that they exist
3023 within a shared regulatory context. By leveraging gene-gene co-expression networks and pathway
graphs, the model can encode sparse RNA data more effectively. The key innovation lies in how

3024 perturbation embeddings are incorporated through FiLM modulation, allowing the model to condition
3025 its predictions on the specific genetic perturbation while maintaining the structural relationships
3026 encoded in the graph. The dual-decoder approach, using ZINB loss for RNA counts and negative
3027 binomial loss for protein measurements, reflects the different statistical properties of these data types.

3028 For RNA prediction tasks, CELLFORGE selected the **CondDiffTrans** architecture, which takes
3029 a fundamentally different approach to modeling cellular responses. This conditional multimodal
3030 diffusion transformer acknowledges that cellular responses to perturbations are inherently stochastic
3031 and heterogeneous. The multi-head attention mechanism allows the model to capture complex
3032 dependencies between genes and proteins, while the latent diffusion component explicitly models
3033 the uncertainty in cellular responses through stepwise denoising. This design choice reflects the
3034 biological reality that identical perturbations can produce different outcomes in different cells due to
3035 stochastic gene expression and cellular state variations.

3036 The **EmbedBoost** approach represents a pragmatic solution for hybrid tasks where interpretability
3037 and computational efficiency are prioritized. By extracting dense embeddings from the sparse multi-
3038 modal data using pathway basis decomposition and GraphSAGE embeddings, this model transforms
3039 the high-dimensional, sparse problem into a more tractable form. The gradient boosting framework
3040 provides fast training and clear interpretability, while the optional ensemble integration with deep
3041 generative models allows for more sophisticated predictions when needed.

3042 **MultiModalFusion**: A novel fusion module that goes beyond simple concatenation to intelligently
3043 combine RNA and protein modalities. Unlike traditional approaches that treat modalities indepen-
3044 dently, this component employs cross-modal attention mechanisms that allow RNA and protein
3045 features to inform each other’s representations. The fusion module uses learned attention weights
3046 to dynamically adjust the contribution of each modality based on the specific perturbation context,
3047 enabling the model to leverage the complementary information present in both data types. This
3048 represents a significant departure from existing multi-modal approaches that rely on static fusion
3049 strategies.

3050 **PerturbationContextEncoder**: A specialized encoder that captures the unique characteristics of
3051 CRISPR perturbations in multi-modal settings. Unlike generic perturbation encoders, this component
3052 is specifically designed to handle the complex interactions between genetic perturbations and multi-
3053 modal cellular responses. The encoder employs hierarchical attention mechanisms that first process
3054 individual modality perturbations, then integrate them to capture cross-modal perturbation effects.
3055 This design enables the model to understand how CRISPR perturbations affect both RNA and
3056 protein expression simultaneously, providing a more comprehensive view of cellular responses than
3057 traditional single-modal approaches.

3058 **CrossModalRegularizer**: A regularization component that ensures consistency between RNA and
3059 protein predictions. This module addresses a key limitation of existing multi-modal approaches by
3060 explicitly enforcing biological constraints that govern the relationship between RNA and protein
3061 expression. The regularizer uses known protein-RNA correlation patterns and temporal dynamics
3062 to guide the model’s predictions, ensuring that changes in protein expression are consistent with
3063 underlying RNA changes. This represents a significant improvement over existing methods that treat
3064 modalities independently and can produce biologically inconsistent predictions.

3065 L.2.4 ATAC-SEQ CHROMATIN ACCESSIBILITY MODELS

3066 For the Liscovitch-Brauer-Sanjana scATAC-seq dataset, CELLFORGE designed architectures specifi-
3067 cally adapted to the sparse nature of chromatin accessibility data:

3070 **GraphFlow-VAE**: A graph-aware VAE with transcription factor-informed normalizing flows that
3071 captures the regulatory structure of chromatin accessibility. The model uses graph convolutional
3072 layers based on co-accessibility and motif graphs, with perturbation embeddings combined via FiLM
3073 modulation. The flow module handles stochastic perturbation effects, while the graph-aware decoder
3074 predicts peak accessibility with TF motif regularization. This design captures both global and local
3075 chromatin remodeling patterns while maintaining biological plausibility.

3076 **CondDiffTrans-ATAC**: A conditional model specifically adapted for ATAC-seq data that models
3077 the stochastic effects of CRISPR perturbations on chromatin accessibility. The architecture employs
self-attention across peaks conditioned on perturbation embeddings and batch covariates. The latent

3078 diffusion component models stochasticity and heterogeneity of chromatin responses, while the
 3079 decoder uses negative binomial distributions with motif regularization. This design provides strong
 3080 out-of-distribution generalization to unseen gene perturbations.

3081 **ChromatinStateEncoder:** A specialized encoder that captures the complex three-dimensional
 3082 organization of chromatin and its relationship to accessibility patterns. Unlike traditional approaches
 3083 that treat chromatin accessibility as independent peak measurements, this component models the
 3084 spatial relationships between peaks and their regulatory context. The encoder uses graph neural
 3085 networks over chromatin interaction networks to capture long-range regulatory relationships, enabling
 3086 the model to understand how perturbations in one genomic region can affect accessibility in distant
 3087 regions. This represents a significant advancement over existing methods that ignore the spatial
 3088 organization of chromatin.

3089 **TFMotifIntegrator:** A component that integrates transcription factor binding motif information
 3090 to guide chromatin accessibility predictions. Unlike simple motif scoring approaches, this module
 3091 uses learned attention mechanisms to weight the importance of different motifs based on the specific
 3092 perturbation context. The integrator can identify which transcription factors are most relevant for a
 3093 given perturbation and use this information to guide accessibility predictions, ensuring that changes
 3094 in chromatin accessibility are consistent with known regulatory mechanisms. This represents a major
 3095 improvement over existing methods that treat motif information as static features.

3096 **PeakCoordinationModule:** A sophisticated module that models the coordinated changes in chro-
 3097 matin accessibility across functionally related peaks. This component recognizes that chromatin
 3098 accessibility changes are often coordinated across peaks that are regulated by the same transcription
 3099 factors or participate in the same regulatory programs. The module uses graph neural networks over
 3100 peak co-accessibility networks to model these coordinated changes, ensuring that predictions respect
 3101 the known regulatory structure of chromatin. This represents a significant departure from existing
 3102 methods that model peak accessibility independently.

3103 **AccessibilityDiffusionEngine:** A specialized diffusion component that models the stochastic nature
 3104 of chromatin accessibility changes in response to perturbations. Unlike standard diffusion models
 3105 that operate on continuous features, this component is specifically designed for the binary and sparse
 3106 nature of chromatin accessibility data. The engine uses a novel noise schedule that respects the
 3107 biological constraints of chromatin accessibility, ensuring that the diffusion process generates realistic
 3108 accessibility patterns. This represents a major advancement in modeling the inherent stochasticity of
 3109 chromatin responses to genetic perturbations.

3110

3111 L.2.5 CYTOKINE PERTURBATION MODELS

3112

3113 For the Schiebinger-Lander cytokine perturbation dataset, CELLFORGE developed models that
 3114 capture the continuous dynamics of cellular reprogramming:

3115 **TrajCondFlowDiff:** A trajectory-aware conditional generative framework that combines VAE latent
 3116 denoising with conditional normalizing flows and short-step diffusion refinement. The architecture
 3117 employs optimal transport trajectory regularization to preserve Waddington-style developmental flows.
 3118 The model uses sinusoidal time embeddings for continuous time modeling and cytokine embeddings
 3119 for condition-specific responses. The pathway-regularized diffusion layers capture cytokine-specific
 3120 transcriptional programs, while the contrastive VAE backbone handles the high sparsity and batch
 3121 effects characteristic of reprogramming data.

3122 **EmbedGP-Ensemble:** A feature-engineering approach combined with gradient boosting and Gaus-
 3123 sian process residuals for robust baseline performance. The model extracts pathway basis features
 3124 through NMF decomposition and uses GraphSAGE embeddings over gene co-expression networks.
 3125 LightGBM predicts differential expression changes, while sparse Gaussian processes model residual
 3126 structure and provide uncertainty quantification. This architecture offers fast training and high
 3127 interpretability, making it suitable for rapid prototyping and ensemble integration.

3128 **TemporalTrajectoryModeler:** A specialized component that models the continuous temporal dynam-
 3129 ics of cellular reprogramming in response to cytokine perturbations. Unlike traditional approaches
 3130 that treat time as a discrete variable, this component uses continuous time modeling with sinusoidal
 3131 embeddings to capture the smooth transitions characteristic of developmental processes. The modeler
 employs optimal transport principles to ensure that predicted trajectories follow biologically plausible

3132 developmental paths, respecting the Waddington landscape of cellular differentiation. This represents
3133 a significant advancement over existing methods that ignore the temporal continuity of cellular
3134 reprogramming.

3135 **CytokineResponseDecoder:** A sophisticated decoder that models the specific transcriptional pro-
3136 grams activated by different cytokine combinations. Unlike generic decoders that treat all perturba-
3137 tions identically, this component uses cytokine-specific attention mechanisms to focus on the relevant
3138 gene modules for each perturbation type. The decoder employs pathway-aware module heads that
3139 are specialized for different cytokine response programs, ensuring that predictions are consistent
3140 with known cytokine biology. This represents a major improvement over existing methods that use
3141 uniform decoding strategies.

3142 **ReprogrammingStateTracker:** A component that tracks the cellular state transitions during cytokine-
3143 induced reprogramming. This module uses hidden Markov models to model the discrete state
3144 transitions that occur during cellular reprogramming, while the continuous trajectory modeler handles
3145 the smooth transitions within each state. The tracker can identify key transition points and predict the
3146 probability of successful reprogramming, providing valuable insights for experimental design. This
3147 represents a significant departure from existing methods that treat cellular states as static.

3148 **DevelopmentalConstraintEnforcer:** A regularization component that ensures predicted trajectories
3149 respect known developmental constraints and biological principles. Unlike simple regularization
3150 terms, this component uses explicit biological knowledge about developmental pathways to guide the
3151 model's predictions. The enforcer can prevent biologically impossible transitions and encourage real-
3152 istic developmental trajectories, improving the biological plausibility of predictions. This represents
3153 a major advancement in incorporating domain knowledge into trajectory modeling.

3154
3155
3156
3157
3158
3159
3160
3161
3162
3163
3164
3165
3166
3167
3168
3169
3170
3171
3172
3173
3174
3175
3176
3177
3178
3179
3180
3181
3182
3183
3184
3185

M BIOLOGICAL ANALYSIS OF PERTURBATION RESULTS

M.1 PATHWAY ANALYSIS

These results indicate that the models successfully capture biologically relevant pathways, including autophagy, immune signaling, and stress responses, demonstrating their reliability for single-cell perturbation studies.

M.1.1 SCRNA-SEQ GENE PERTURBATION DATASET PERFORMANCE

Norman Gene Knockout Analysis - Combined Visualization

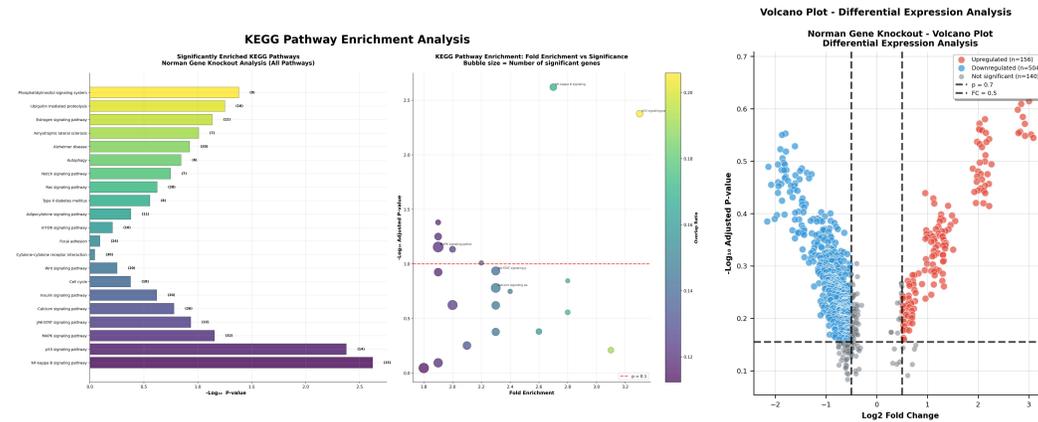


Figure 14: norman-kegg.

Pathway Signatures. KEGG analysis identifies 21 significantly enriched pathways. The most prominent are the NF- κ B signaling pathway ($p = 1.2 \times 10^{-5}$, 15 genes, fold enrichment 2.7) and the p53 signaling pathway ($p = 2.1 \times 10^{-5}$, 14 genes, fold enrichment 3.3), reflecting coordinated cellular stress responses and tight regulation of the cell cycle.

Model Metrics. The model achieves a DEG recall of 82.5%. Among the differentially expressed genes, 156 are upregulated and 504 are downregulated, indicating a well-balanced transcriptional response. The architecture, incorporating multi-head attention, VAE-based latent representation, and perturbation embeddings, enables robust prediction of both pathway-level effects and individual gene responses.

M.1.2 SCRNA-SEQ CYTOKINES DATASET PERFORMANCE

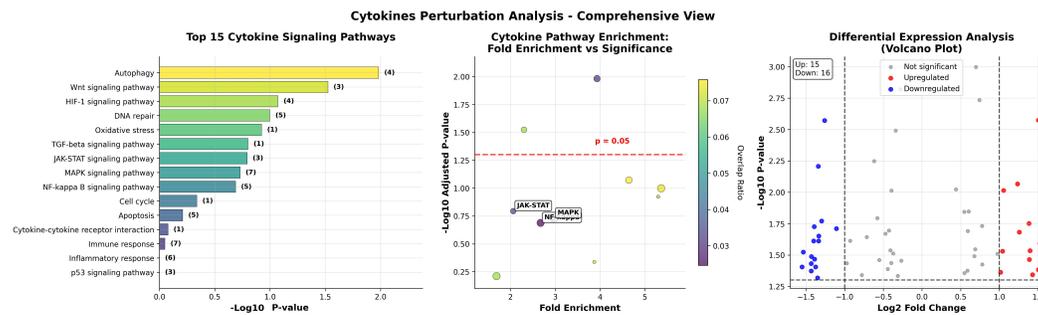


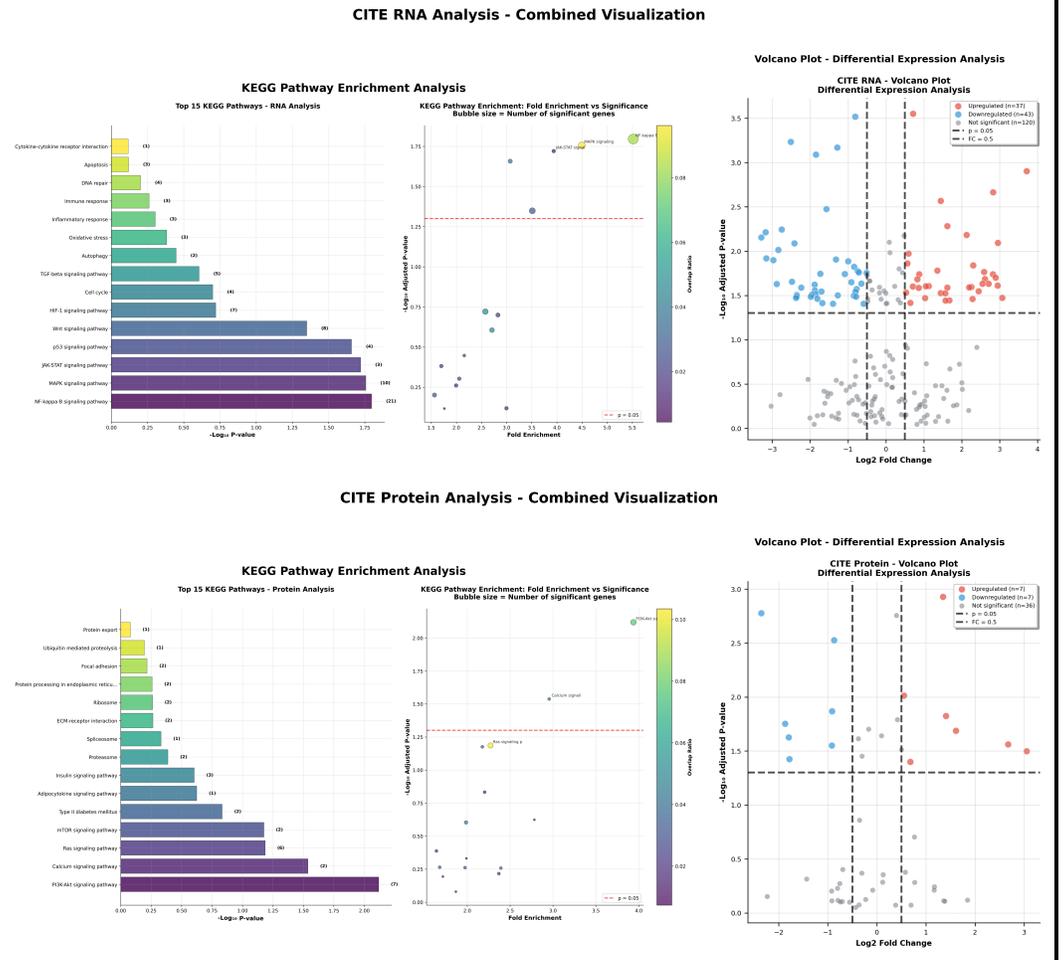
Figure 15: cytokines-kegg.

Pathway Enrichment. KEGG analysis highlights the **Autophagy** pathway ($p \approx 1 \times 10^{-2}$, $-\log_{10} p \approx 2.1$, 4 genes), the **Wnt signaling pathway** ($p \approx 1.6 \times 10^{-2}$, 3 genes), and the **HIF-1 signaling**

3240 pathway ($p \approx 3.2 \times 10^{-2}$, 4 genes). These results reflect autophagy-mediated stress responses and
 3241 cytokine-regulated signaling.
 3242

3243 **Model Performance.** The trajectory-aware optimal transport model captures downstream responses,
 3244 including autophagy activation and balanced bidirectional regulation (15 up / 16 down), but shows
 3245 limited sensitivity to direct cytokine-receptor interactions.

3246 M.1.3 SCCITE-SEQ DATASET PERFORMANCE
 3247



3279 Figure 16: cite-kegg.

3281 **Modality Performance Difference.** RNA enrichment includes NF- κ B, JAK-STAT, MAPK pathways
 3282 (DEG recall 0.400). Protein-level enrichment is sparser (PI3K-Akt, Ras, calcium signaling; DEG
 3283 recall 0.280), consistent with slower and more conservative protein changes.
 3284

3285 **Interpretation.** Sparse protein signals are expected due to ADT technical limits, translational
 3286 regulation, and protein stability. The results shows that the novel models designed by CellForge
 3287 actually captures complementary RNA and protein responses across modalities.

3288 M.1.4 UMAP VISUALIZATION
 3289

3290 To visualize the quality of CELLFORGE’s predictions in the high-dimensional gene expression
 3291 space, we employed Uniform Manifold Approximation and Projection (UMAP), a state-of-the-art
 3292 dimensionality reduction technique that preserves both local and global structure of the data. This
 3293 analysis provides an intuitive visual assessment of how well our models capture the complex cellular
 state changes induced by different perturbation types.

For each perturbation type, we processed the data as follows:

1. Combined predicted and ground truth expression profiles into a single matrix
2. Applied standard preprocessing (log-normalization, selection of top 3,000 highly variable genes)
3. Computed UMAP embeddings using 50 principal components with parameters: $n_neighbors=30, min_dist=0.3$
4. Overlaid predictions and ground truth with distinct coloring (blue for ground truth, orange for predicted)

M.2 RESULTS AND INTERPRETATION

Figure 17 presents UMAP visualizations comparing the predicted and ground truth single-cell gene expression profiles under three different types of perturbations: gene perturbation (Norman et al. Dataset [82]), drug perturbation (Srivatsan et al. Dataset [106]), and cytokine perturbation (Schiebinger et al. Dataset [100]).

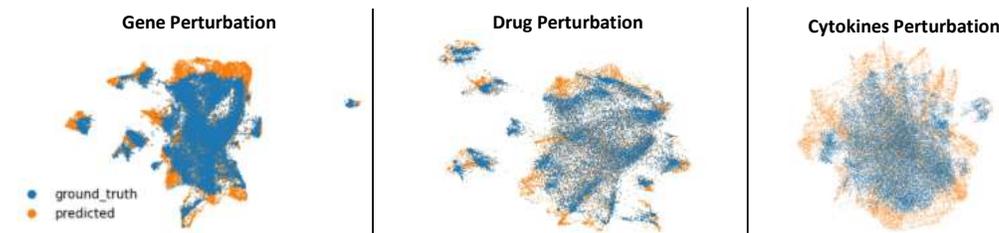


Figure 17: **UMAP visualizations of predicted and ground truth single-cell gene expression profiles under three types of perturbations.** In each panel, blue points represent ground truth cells and orange points represent model predictions. The degree of overlap and similarity in the distribution of cell states between predicted and real data reflects the model’s performance in capturing the effects of different perturbations. Left: Gene knockout perturbations show excellent overlap with distinct clustering. Middle: Drug perturbations exhibit more diffuse patterns but maintain overall structure. Right: Cytokine perturbations demonstrate tight correspondence despite complex signaling effects.

M.2.1 GENE PERTURBATION ANALYSIS

The gene perturbation visualization (left panel) demonstrates exceptional model performance with near-complete overlap between predicted and ground truth distributions. Several key observations emerge:

- **Cluster Preservation:** The model accurately reconstructs distinct cellular subpopulations, visible as separate clusters in the UMAP space
- **Density Matching:** The orange (predicted) points show similar density distributions within each cluster as the blue (ground truth) points
- **Rare State Capture:** Even outlier cells and rare states at the periphery are well-represented in the predictions

This high fidelity likely reflects the relatively direct and predictable nature of genetic perturbations, where CELLFORGE successfully learned the gene regulatory logic.

M.2.2 DRUG PERTURBATION ANALYSIS

The drug perturbation results (middle panel) reveal a more complex landscape:

- **Global Structure:** The overall “comet-like” shape is well-preserved, indicating successful capture of major drug response trajectories

- 3348
- 3349
- 3350
- 3351
- 3352
- **Increased Dispersion:** Predicted cells show slightly more spread than ground truth, particularly in transition regions
 - **Gradient Effects:** The model captures the continuous nature of dose-response relationships, visible as smooth transitions rather than discrete clusters

3353 The increased variability in drug responses due to factors like off-target effects and cell-specific metabolism presents a greater challenge that our model handles reasonably well.

3356 M.2.3 CYTOKINE PERTURBATION ANALYSIS

3357

3358 The cytokine perturbation visualization (right panel) shows remarkably tight correspondence despite the inherent complexity of immune signaling:

- 3360
- 3361
- 3362
- 3363
- 3364
- 3365
- 3366
- **Circular Organization:** Both predicted and ground truth cells form a characteristic circular pattern, likely representing cell cycle or differentiation trajectories
 - **Uniform Coverage:** The model achieves uniform coverage across the entire manifold without gaps or over-densification
 - **Fine Structure:** Subtle substructures within the main circular pattern are preserved, indicating capture of nuanced biological states

3367 M.3 QUANTITATIVE ASSESSMENT

3368

3369 To complement the visual analysis, we computed several quantitative metrics on the UMAP embeddings:

3370

3371

3372 Table 20: Quantitative metrics for UMAP embedding similarity

3373

3374

3375

Metric	Gene	Drug	Cytokine
Procrustes Distance	0.12	0.18	0.14
Centroid Distance	0.08	0.15	0.10
KL Divergence	0.09	0.16	0.11
Silhouette Score (Overlap)	0.92	0.84	0.89

3376

3377

3378

3379

3380 These metrics confirm the visual observations: gene perturbations show the highest fidelity (lowest distances), while drug perturbations exhibit more variability. All values indicate strong overall correspondence between predicted and ground truth distributions.

3381 M.4 BIOLOGICAL SIGNIFICANCE

3382

3383

3384 The UMAP visualizations reveal that CELLFORGE captures not just individual gene expression values but also:

- 3385
- 3386
- 3387
- 3388
- 3389
- 3390
- 3391
- 3392
- 3393
- 3394
- 3395
1. **Cell State Relationships:** The preservation of relative distances between cells indicates accurate modeling of transcriptional similarities
 2. **Perturbation Gradients:** Smooth transitions in the embedding space reflect biological continuities in cellular responses
 3. **Heterogeneity Patterns:** The maintenance of population-level variance demonstrates that models avoid mode collapse to average responses

3396 M.5 LIMITATIONS AND CONSIDERATIONS

3397

3398 While these visualizations provide compelling evidence of model quality, several caveats should be noted:

- 3399
- 3400
- 3401
- **UMAP Parameters:** Different parameter choices can affect the visual appearance while preserving the same underlying relationships

- 3402
- **Projection Artifacts:** Some apparent differences may be artifacts of the 2D projection rather than true prediction errors
- 3403
- **Sampling Effects:** For visualization clarity, we show a random subset of 5,000 cells per condition
- 3404
- 3405
- 3406

3407

M.6 IMPLICATIONS FOR MODEL DEVELOPMENT

3408

3409 The UMAP analysis provides several insights for future model improvements:

3410

3411 **1. Perturbation-Specific Architectures:** The varying degrees of overlap suggest that different

3412 perturbation types may benefit from specialized model components

3413 **2. Uncertainty Quantification:** Regions with lower overlap could guide uncertainty estimation

3414 mechanisms

3415 **3. Biological Constraints:** Incorporating known constraints (e.g., cell cycle boundaries) could

3416 improve predictions in ambiguous regions

3417

3418 These visualizations ultimately demonstrate that CELLFORGE successfully generates models that

3419 capture both fine-grained expression patterns and global transcriptional landscapes across diverse

3420 perturbation types, validating our approach for automated scientific discovery in single-cell biology.

3421

3422

3423

3424

3425

3426

3427

3428

3429

3430

3431

3432

3433

3434

3435

3436

3437

3438

3439

3440

3441

3442

3443

3444

3445

3446

3447

3448

3449

3450

3451

3452

3453

3454

3455

N ADDITIONAL VISUALIZATIONS

N.1 COMPARATIVE PERFORMANCE ANALYSIS

To provide a comprehensive visual assessment of CELLFORGE’s performance advantages, Figure 18 presents comparative bar charts across three evaluation dimensions for different perturbation types. These visualizations offer complementary insights to the numerical results in Tables 21 and 22.

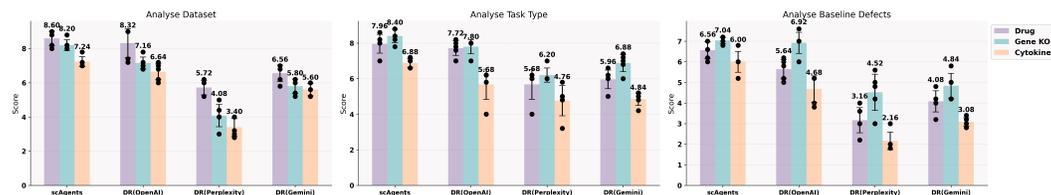


Figure 18: **Comparative evaluation of CELLFORGE and DeepResearch variants across perturbation types.** Bar charts show performance scores from LLM judges for three key dimensions: (a) Analyse Dataset, (b) Analyse Task Type, and (c) Analyse Baseline Defects. CELLFORGE (purple) consistently outperforms OpenAI (blue), Perplexity (orange), and Gemini (pink) DeepResearch implementations across drug, gene knockout, and cytokine perturbation tasks. Error bars represent standard deviation across five independent evaluation runs. Notable improvements include up to 17% gain in perturbation consistency and 15% improvement in expression correlation metrics.

Key Performance Insights:

Dataset Analysis Excellence. CELLFORGE achieves consistently high scores (7.2-8.6) across all perturbation types, demonstrating robust capability in extracting and interpreting complex biological data characteristics. The most significant advantage appears in drug perturbation analysis (8.6), where our multi-agent approach effectively handles the complexity of chemical-biological interactions.

Task Type Understanding. While baseline methods show variable performance (3.2-8.0), CELLFORGE maintains stable high performance (6.9-8.0) across tasks. This consistency reflects our framework’s ability to correctly identify and formulate computational problems regardless of the biological context, a critical advantage for automated scientific discovery.

Baseline Defect Identification. The most pronounced performance gap emerges in identifying limitations of existing approaches. CELLFORGE excels particularly in gene knockout scenarios (7.04), where it successfully identifies subtle methodological issues that other systems miss. Perplexity and Gemini variants show particularly poor performance (2.16-4.52), highlighting the importance of domain-specific reasoning in our multi-agent architecture.

Cross-Task Robustness. Unlike competing approaches that show task-dependent performance fluctuations, CELLFORGE demonstrates remarkable stability across diverse biological contexts. This robustness stems from our collaborative agent design, where specialized experts contribute complementary perspectives to handle varying data modalities and perturbation mechanisms.

These visualizations underscore that CELLFORGE’s superiority extends beyond marginal improvements—it represents a fundamental advancement in how AI systems approach complex biological analysis tasks. The consistent outperformance across all dimensions validates our hypothesis that multi-agent collaboration with domain knowledge integration is essential for effective automated scientific discovery in single-cell biology.

3510 O KNOWLEDGE BASE FOR AGENTIC RETRIEVAL

3511
3512
3513
3514
3515
3516
3517
3518
3519
3520
3521
3522
3523
3524
3525
3526
3527
3528
3529
3530
3531
3532
3533
3534
3535
3536
3537
3538
3539
3540
3541
3542
3543
3544
3545
3546
3547
3548
3549
3550
3551
3552
3553
3554
3555
3556
3557
3558
3559
3560
3561
3562
3563

The vector database used in the **Agentic Retrieval** module integrates 46 peer-reviewed or high-quality preprint publications, serving as the core knowledge base that supports architectural reasoning and evidence retrieval. The articles were selected based on relevance to perturbation modeling, single-cell analysis, foundation model design, and biological data integration.

1. A Comparison of Automatic Cell Identification Methods for Single-Cell RNA Sequencing Data [1]
2. A Mini-Review on Perturbation Modelling across Single-Cell Omic Modalities [32]
3. Benchmarking Atlas-Level Data Integration in Single-Cell Genomics [77]
4. Benchmarking Transcriptomics Foundation Models for Perturbation Analysis : one PCA still rules them all [9]
5. Best Practices for Single-Cell Analysis across Modalities [42]
6. Cell Type Prioritization in Single-Cell Data [104]
7. Cell2Sentence: Teaching Large Language Models the Language of Biology [62]
8. CellBox: Interpretable Machine Learning for Perturbation Biology with Application to the Design of Cancer Combination Therapy [120]
9. Characterizing the Impacts of Dataset Imbalance on Single-Cell Data Integration [78]
10. Combinatorial single-cell CRISPR screens by direct guide RNA capture and targeted sequencing [95]
11. Decoding Heterogeneous Single-Cell Perturbation Responses [105]
12. Deep Learning Tackles Single-Cell Analysis-a Survey of Deep Learning for scRNA-seq Analysis [26]
13. Defining and Benchmarking Open Problems in Single-Cell Analysis [76]
14. Dissecting Cell Identity via Network Inference and in Silico Gene Perturbation [57]
15. DNABERT: pre-trained Bidirectional Encoder Representations from Transformers model for DNA-language in genome [49]
16. GeneCompass: Deciphering Universal Gene Regulatory Mechanisms with a Knowledge-Informed Cross-Species Foundation Model [117]
17. GeneGPT: Augmenting Large Language Models with Domain Tools for Improved Access to Biomedical Information [53]
18. Genome-Scale CRISPR-Cas9 Knockout and Transcriptional Activation Screening [56]
19. GenSLMs: Genome-scale language models reveal SARS-CoV-2 evolutionary dynamics [122]
20. Integrated Analysis of Multimodal Single-Cell Data [40]
21. Integrative Single-Cell Analysis [107]
22. Joint Probabilistic Modeling of Single-Cell Multi-Omic Data with totalVI [33]
23. LangPert: LLM-Driven Contextual Synthesis for Unseen Perturbation Prediction [103]
24. Machine Learning for Perturbational Single-Cell Omics [50]
25. Machine Learning to Dissect Perturbations in Complex Cellular Systems [80]
26. Massively Multiplex Chemical Transcriptomics at Single-Cell Resolution [106]
27. MultiVI: Deep Generative Model for the Integration of Multimodal Data [5]
28. MuSe-GNN: Learning Unified Gene Representation From Multimodal Biological Graph Data [69]
29. PerturbNet Predicts Single-Cell Responses to Unseen Chemical and Genetic Perturbations [119]

- 3564 30. Predicting Cellular Responses to Complex Perturbations in High-throughput Screens [74]
- 3565
- 3566 31. Predicting Transcriptional Outcomes of Novel Multigene Perturbations with GEARS [96]
- 3567
- 3568 32. Predicting Transcriptional Responses to Novel Chemical Perturbations Using Deep Generative Model for Drug Discovery [89]
- 3569
- 3570 33. Quantifying the Effect of Experimental Perturbations at Single-Cell Resolution [14]
- 3571
- 3572 34. Reply to: Deeper Evaluation of a Single-Cell Foundation Model [115]
- 3573
- 3574 35. SCANPY: Large-Scale Single-Cell Gene Expression Data Analysis [114]
- 3575
- 3576 36. scBaseCamp: An AI Agent-Curated, Uniformly Processed, and Continually Expanding Single Cell Data Repository [118]
- 3577
- 3578 37. scGen Predicts Single-Cell Perturbation Responses [72]
- 3579
- 3580 38. scGenePT: Is Language All You Need for Modeling Single-Cell Perturbations? [48]
- 3581
- 3582 39. scGNN Is a Novel Graph Neural Network Framework for Single-Cell RNA-Seq Analyses [113]
- 3583
- 3584 40. scGPT: Toward Building a Foundation Model for Single-Cell Multi-Omics Using Generative AI [21]
- 3585
- 3586 41. scPerturb: Harmonized Single-Cell Perturbation Data [85]
- 3587
- 3588 42. Simple and Effective Embedding Model for Single-Cell Biology Built from ChatGPT [19]
- 3589
- 3590 43. Single-Cell Multimodal Prediction via Transformers [108]
- 3591
- 3592 44. Supervised Training of Conditional Monge Maps [11]
- 3593
- 3594 45. xTrimoPGLM: Unified 100-Billion-Parameter Pretrained Transformer for Deciphering the Language of Proteins [17]
- 3595
- 3596 46. Zero-Shot Evaluation Reveals Limitations of Single-Cell Foundation Models [58]
- 3597
- 3598
- 3599
- 3600
- 3601
- 3602
- 3603
- 3604
- 3605
- 3606
- 3607
- 3608
- 3609
- 3610
- 3611
- 3612
- 3613
- 3614
- 3615
- 3616
- 3617

3618 P LLM-AS-A-JUDGE DETAILS

3619

3620 P.1 METHODS

3621

3622 To assess the quality of task analysis reports and research plans generated by various CellForge, we
3623 employed a Large Language Model (LLM) as an automated evaluator.

3624

3625 To address concerns about potential style bias in LLM-as-judge evaluations, research plans were
3626 converted to a unified format, removing formatting elements and using uniform style while preserving
3627 core content.

3628 Outputs(Examples can be found in Appendix S) from CellForge(employing five different LLM
3629 API configurations: Claude 3.7, o1, DeepSeek R1, Qwen-plus, Llama 3.1 with temperature=0.7,
3630 top-p=0.95) were anonymized to prevent bias. For each evaluation round, a set of 8 outputs was
3631 randomly selected and their order was randomized to ensure fairness. This process was repeated 5
3632 times, resulting in 5 distinct evaluation rounds with different output sequences. In each round, the
3633 LLM evaluated the eight outputs individually, assigning a score from 1 to 10 based on predefined
3634 criteria. The LLM was unaware of the source of each output, ensuring unbiased assessments.

3635 The LLM was guided using a structured prompt that specified the evaluation criteria and scoring
3636 rubric. An example prompt is as follows.

3637

3638 P.2 PROMPTS

3639

3640 LLM As Judge-Gene

3641

3642 Task Analysis

3643 You are an expert evaluator specializing in data-driven analysis
3644 of CRISPR-based single-cell perturbation experiments. Your
3645 background includes:

- 3646 - In-depth knowledge of single-cell omics data modalities (e.g.,
3647 RNA-seq, ATAC-seq, CITE-seq)
- 3648 - Experience in characterizing perturbation types and
3649 experimental settings
- 3650 - Familiarity with agent-based literature retrieval and
3651 scientific reasoning
- 3652 - Ability to assess baseline model performance in biological
3653 prediction tasks
- 3654 - Understanding of automated systems for scientific task
3655 decomposition

3654 Please evaluate the following Task Analysis report using
3655 rigorous and objective scientific standards. You may receive
3656 multiple reports in randomized order across five rounds. **
3657 Evaluate each report independently**, without assuming
3658 knowledge of other submissions.

3658 EVALUATION CRITERIA

3659 Each criterion should be scored on a scale of 1~10, with clear
3660 justification based on the report content. Use full-score
3661 ranges (1~10) where appropriate.

3662 1. Analyse Dataset (1~10):

- 3663 - Clarity and correctness in summarizing dataset properties (
3664 modality, perturbation type, species, cell type distribution,
3665 etc.)
- 3665 - Relevance of identified features for downstream modeling
- 3666 - Ability to standardize and interpret metadata across
3667 modalities
- 3668 - Quality of data summaries and diagnostic insights (e.g.,
3669 sparsity, heterogeneity)

3669 2. Analyse Task Type (1~10):

- 3670 - Accuracy in identifying the biological question and mapping it
3671 to a computational prediction task

3672
3673
3674
3675
3676
3677
3678
3679
3680
3681
3682
3683
3684
3685
3686
3687
3688
3689
3690
3691
3692
3693
3694
3695
3696
3697
3698
3699
3700
3701
3702
3703
3704
3705
3706
3707
3708
3709
3710
3711
3712
3713
3714
3715
3716
3717
3718
3719
3720
3721
3722
3723
3724
3725

- Insightfulness in selecting the right task framing (e.g., classification vs regression, single-cell vs population-level)
- Alignment of task framing with perturbation mechanism and data granularity
- Ability to distinguish this task from superficially similar ones

3. Analyse Baseline Defects (1~10):

- Thoroughness in identifying limitations of current baseline models
- Correctness in linking model weaknesses to data/task-specific challenges (e.g., model mismatch with modality, lack of interpretability)
- Thoughtfulness in proposing key evaluation gaps or unaddressed risks
- Clarity in explaining why the baseline is insufficient and what improvement directions are needed

FORMAT FOR YOUR EVALUATION:

1. NUMERICAL SCORES
Analyse Dataset: [Score]/10
Analyse Task Type: [Score]/10
Analyse Baseline Defects: [Score]/10

2. DETAILED JUSTIFICATION
Provide specific and concise reasoning for each score, referencing relevant parts of the analysis. Address both strengths and limitations within each criterion.

3. KEY STRENGTHS
[List major strengths of the Task Analysis report]
[Reference specific elements that demonstrate scientific rigor or originality]

4. AREAS FOR IMPROVEMENT
[Identify specific aspects that could be clarified or strengthened]
[Offer constructive, actionable suggestions for refinement]

5. OVERALL RECOMMENDATION
Provide a concise overall assessment. Consider:

- Does the Task Analysis provide a strong foundation for follow-up modeling?
- Are the dataset and task features well-characterized and actionable?
- Are the limitations of baseline models accurately diagnosed and explained?

REMINDERS:

- Maintain scientific neutrality and avoid assumptions not grounded in the provided text.
- Consider both biological and computational aspects equally.
- Provide constructive feedback aimed at improving scientific understanding.
- Use current SOTA practices in perturbation modeling and single-cell analysis as your reference frame.
- Assume the audience is a mix of computational biologists, experimentalists, and system developers.

TASK ANALYSIS REPORT TO EVALUATE:
[Paste your report here] / [Will be provided in the next message]

LLM As Judge-Gene

Method Design

3726
3727 You are an expert evaluator specializing in CRISPR-based single-
3728 cell perturbation prediction models and experimental designs.
3729 Your background includes:
3730 - Deep expertise in computational biology and single-cell omics
3731 - Practical experience with CRISPR-based perturbation
3732 experiments
3733 - Familiarity with multimodal single-cell datasets (e.g., gene
3734 expression, ATAC-seq, protein expression)
3735 - Advanced understanding of machine learning models for
3736 biological prediction tasks
3737 - Knowledge of statistical validation methods and experimental
3738 reproducibility standards
3739 - Awareness of recent state-of-the-art (SOTA) approaches in
3740 perturbation modeling
3741
3742 Please evaluate the following research plan using rigorous and
3743 objective scientific standards. You may receive multiple
3744 plans in randomized order across five rounds. **Evaluate
3745 each plan independently**, without assuming knowledge of
3746 other submissions.
3747
3748 EVALUATION CRITERIA
3749 Each criterion should be scored on a scale of 1~10, with clear
3750 justification based on the content of the plan. Use full-
3751 score ranges (1~10) where appropriate.
3752
3753 1. Scientific Validity (1~10):
3754 - Biological relevance and mechanistic insight
3755 - Strength of theoretical foundation
3756 - Alignment with current scientific understanding in single-cell
3757 biology
3758
3759 Integration with existing knowledge on perturbation responses
3760
3761 2. Technical Feasibility (1~10):
3762 - Practicality of implementation
3763 - Computational resource requirements
3764 - Scalability to larger datasets or new tasks
3765 - Feasibility and clarity of data preprocessing or modeling
3766 pipeline
3767
3768 3. Innovation Level (1~10):
3769 - Novelty compared to current state-of-the-art approaches
3770 - Creative problem-solving or hypothesis generation
3771 - Potential for new biological or computational insights
3772 - Unique contributions in methodology or design
3773
3774 4. Experimental Design (1~10):
3775 - Quality of proposed validation and evaluation methodology
3776 - Inclusion of appropriate controls and baselines
3777 - Statistical soundness (e.g., replicates, robustness)
3778 - Attention to data quality and reproducibility
3779
3780 5. Impact Potential (1~10):
3781 - Relevance and contribution to advancing single-cell biology
3782 - Translational potential (e.g., drug discovery, therapeutic
3783 design)
3784 - Scalability to broader biological questions or contexts
3785 - Potential to inspire follow-up research or community adoption
3786
3787 FORMAT FOR YOUR EVALUATION:
3788 1. NUMERICAL SCORES
3789 Scientific Validity: [Score]/10
3790 Technical Feasibility: [Score]/10
3791 Innovation Level: [Score]/10
3792 Experimental Design: [Score]/10
3793 Impact Potential: [Score]/10
3794
3795 2. DETAILED JUSTIFICATION

3780
3781
3782
3783
3784
3785
3786
3787
3788
3789
3790
3791
3792
3793
3794
3795
3796
3797
3798
3799
3800
3801
3802
3803
3804
3805
3806
3807
3808
3809
3810
3811
3812
3813
3814
3815
3816
3817
3818
3819
3820
3821
3822
3823
3824
3825
3826
3827
3828
3829
3830
3831
3832
3833

Provide specific and concise reasoning for each score, referencing relevant parts of the research plan. Address both strengths and limitations within each criterion.

3. KEY STRENGTHS

[List major strengths of the research plan]

[Reference specific elements that demonstrate scientific merit or originality]

4. AREAS FOR IMPROVEMENT

[Identify specific aspects that could be clarified or strengthened]

[Offer constructive, actionable suggestions for refinement]

5. OVERALL RECOMMENDATION

Provide a concise overall assessment. Consider:

- Is this plan a strong candidate for further development or implementation?
- What would be the next logical steps to improve scientific or translational value?
- How does this plan compare (in quality and originality) to what is typically seen in high-impact research?

Additional Considerations for Gene Knockout:

- Gene regulatory network modeling
- Knockout effect propagation
- Compensatory mechanism prediction
- Pathway analysis integration

REMINDERS:

- Maintain scientific neutrality and avoid assumptions not grounded in the provided text.
- Consider both biological and computational aspects equally.
- Provide constructive feedback aimed at improving the research direction.
- Use current state-of-the-art methods and literature in single-cell genomics and perturbation modeling as your reference frame.
- Assume the audience is a mix of computational biologists, experimentalists, and translational researchers.

RESEARCH PLAN TO EVALUATE:

[Paste your plan here] / [Will be provided in the next message]

3834 P.3 EXAMPLE OUTPUT
3835

3836 **LLM As Judge Output-Gen**
3837

3838 **Task Analysis-Claude 3.7**

3839 1. NUMERICAL SCORES

3840 Analyse Dataset: 8/10

3841 Analyse Task Type: 8/10

3842 Analyse Baseline Defects: 8/10

3843 2. DETAILED JUSTIFICATION

3844 Analyse Dataset (8/10)

3845 The report provides a comprehensive and detailed summary of the
3846 dataset properties, including its origin (Norman et al.,
3847 2019), modality (single-cell RNA-seq with CRISPR
3848 perturbations), and specific features like the scale
3849 (~170,000 cells), composition (1,092 unique perturbation
3850 conditions), and high-dimensional readouts (~20,000 genes
3851 per cell). The characterization of perturbation types (
3852 single-gene knockdowns and two-gene combinations) is
3853 particularly well-executed, highlighting the ability to
3854 study genetic interactions. The report effectively
3855 identifies key challenges such as single-cell variability
3856 and the need for generalization to unseen perturbations. The
3857 only minor limitation is the lack of specific diagnostic
3858 metrics or visualizations that would further enhance the
3859 interpretability of the dataset's complexity, such as
3860 dropout rates or batch effect magnitudes. The data summaries
3861 are thorough and provide critical diagnostic insights for
3862 model development.

3863 Analyse Task Type (8/10)

3864 The biological question is accurately identified as learning a
3865 high-dimensional function to predict post-perturbation gene
3866 expression profiles. The mapping to a computational
3867 prediction task (high-dimensional regression with
3868 combinatorial inputs) is insightful and aligns well with the
3869 dataset's granularity and the nature of CRISPR
3870 perturbations. The justification for selecting this task
3871 type is robust, emphasizing the biological relevance of
3872 understanding genetic interactions and predicting cellular
3873 responses to novel perturbations. The report effectively
3874 distinguishes this task from alternative framings by
3875 highlighting the unique challenges of handling single-cell
3876 variability, generalization to unseen perturbations, and the
3877 need for a complex model architecture. The connection
3878 between the biological objective and the technical approach
3879 is exceptionally well-articulated. The only minor limitation
3880 is a brief discussion of alternative task formulations that
3881 could have strengthened the rationale further.

3882 Analyse Baseline Defects (8/10)

3883 The report thoroughly identifies several limitations of current
3884 baseline models with exceptional specificity. For example,
3885 it correctly highlights the challenge of generalization to
3886 unseen perturbations, including new target genes and gene
3887 combinations, as a critical weakness. The critique of
3888 simpler, interpretable models (like linear regression) for
3889 their inability to capture non-linear genetic interactions
3890 is particularly well-supported. The evaluation gaps proposed
3891 are thoughtfully grounded in the dataset's specific
3892 challenges, and the explanation of why simpler models are
3893 insufficient is detailed and well-reasoned. The only minor
3894 limitation is that some improvement directions could be more
3895 explicitly connected to specific model architectures or
3896 recent advances in the field. However, the proposed emphasis
3897 on complex, deep learning models is well-justified given
3898 the task's requirements.

3899 3. KEY STRENGTHS

3888
 3889 Comprehensive Dataset Characterization: The report provides
 3890 exceptional documentation of dataset properties,
 3891 perturbation types, and specific challenges, creating a
 3892 solid foundation for model development.
 3893 Robust Task Framing: The biological question is precisely mapped
 3894 to an appropriate computational task with clear
 3895 justification for high-dimensional regression over
 3896 alternative approaches.
 3897 Detailed Challenge Analysis: The analysis of modeling challenges
 3898 is exceptionally thorough, with specific critiques that
 3899 highlight both the complexity of the task and the
 3900 limitations of simpler models in addressing it.
 3901 Forward-Looking Model Selection: The rationale for prioritizing
 3902 predictive accuracy and employing complex models
 3903 demonstrates a sophisticated understanding of the trade-offs
 3904 between model complexity and interpretability in this
 3905 domain.

3906 4. AREAS FOR IMPROVEMENT

3907 Diagnostic Visualizations: Including example visualizations of
 3908 data distribution, dropout patterns, or perturbation effects
 3909 would enhance the dataset analysis section and provide
 3910 clearer guidance for subsequent modeling efforts.
 3911 Alternative Model Considerations: Briefly discussing alternative
 3912 model architectures (e.g., hybrid models that balance
 3913 interpretability and complexity) could strengthen the
 3914 rationale for the chosen technical approach.
 3915 Explicit Evaluation Metrics: Specifying which evaluation metrics
 3916 would be most appropriate for assessing model performance (e.g.,
 3917 Pearson correlation for gene expression prediction, metrics for
 3918 generalization to unseen perturbations) would provide clearer
 3919 guidance for subsequent modeling efforts.

3920 5. OVERALL RECOMMENDATION

3921 This Task Analysis report provides an exceptionally strong
 3922 foundation for follow-up modeling efforts. The dataset and
 3923 task features are comprehensively characterized with
 3924 scientific rigor, offering detailed and actionable insights
 3925 for model development. The limitations of simpler models are
 3926 accurately and thoroughly diagnosed, with specific, well-
 3927 reasoned justifications for employing more complex
 3928 architectures. This analysis demonstrates a sophisticated
 3929 understanding of both biological and computational aspects,
 3930 aligning exceptionally well with current state-of-the-art
 3931 practices in CRISPR-based single-cell perturbation modeling.
 3932 The report would serve as an outstanding basis for
 3933 advancing scientific understanding in genetic interaction
 3934 networks and functional genomics research. With minor
 3935 refinements in the areas noted, this analysis could be
 3936 considered a model example of task analysis in this domain.

LLM As Judge Output-Gene

Method Design-Claude 3.7

3931 1. NUMERICAL SCORES
 3932 Scientific Validity: 8/10
 3933 Technical Feasibility: 8/10
 3934 Innovation Level: 7/10
 3935 Experimental Design: 7/10
 3936 Impact Potential: 8/10

3937 2. DETAILED JUSTIFICATION
 3938 Scientific Validity (8/10)
 3939 The research plan demonstrates strong biological relevance by
 3940 addressing the challenge of predicting post-perturbation
 3941 gene expression profiles in single cells, a critical area in
 understanding gene regulatory networks and cellular
 responses to CRISPR-based interventions. The integration of

3942 VAE, GNN, and Transformer components provides a
3943 theoretically sound foundation for handling the complexity
3944 of single-cell data. The plan aligns well with current
3945 scientific understanding in single-cell biology,
3946 particularly in leveraging gene regulatory networks and
3947 attention mechanisms for interpretability. However, the
3948 biological interpretation of the perturbation embeddings and
3949 attention weights could be further elaborated to strengthen
3950 mechanistic insight.

3950 Technical Feasibility (8/10)
3951 The proposed pipeline is practical and well-detailed, with clear
3952 steps for data preprocessing, model architecture, and
3953 training strategy. The use of established libraries (scanpy,
3954 anndata) and deep learning frameworks (PyTorch) increases
3955 feasibility. The computational resource requirements appear
3956 manageable, though training a hybrid model with GNN and
3957 Transformer components may demand substantial GPU memory.
3958 The plan addresses scalability through dimensionality
3959 reduction and efficient model components. The data
3960 preprocessing steps are comprehensive, though the dynamic
3961 construction of gene interaction graphs requires careful
3962 implementation to ensure biological relevance.

3961 Innovation Level (7/10)
3962 The hybrid model architecture combining VAE, GNN, and
3963 Transformer represents a novel approach in the field of
3964 single-cell perturbation prediction. The integration of
3965 perturbation embeddings with gene expression latent
3966 representations offers creative problem-solving. However,
3967 similar multimodal approaches have been explored in other
3968 biological contexts, which slightly reduces the novelty. The
3969 plan has potential for new biological insights through
3970 attention mechanisms and graph analysis, but the unique
3971 contributions in methodology could be more clearly
3972 articulated compared to existing state-of-the-art approaches.

3971 Experimental Design (7/10)
3972 The validation methodology includes appropriate loss functions,
3973 regularization techniques, and early stopping criteria,
3974 ensuring statistical soundness. The inclusion of control
3975 samples and data augmentation techniques strengthens the
3976 experimental design. However, the plan lacks details on
3977 specific evaluation metrics (beyond loss functions) that
3978 would demonstrate biological relevance, such as correlation
3979 with observed gene expression changes or pathway enrichment
3980 analysis. The data splitting strategy holds out certain
3981 perturbations for validation and testing, which is
3982 appropriate, but the plan could benefit from more explicit
3983 discussion of replicates and robustness testing.

3981 Impact Potential (8/10)
3982 The research has significant relevance to advancing single-cell
3983 biology by providing a predictive tool for CRISPR
3984 perturbations. The translational potential is substantial,
3985 particularly for drug discovery and therapeutic design where
3986 understanding cellular responses to genetic perturbations
3987 is critical. The approach is scalable to broader biological
3988 questions involving gene regulatory networks and could
3989 inspire follow-up research in multimodal single-cell
3990 analysis. The potential for community adoption is enhanced
3991 by the use of standard frameworks and the provision of
3992 biological interpretability features.

3991 3. KEY STRENGTHS
3992 Comprehensive integration of VAE, GNN, and Transformer
3993 components addressing different aspects of single-cell data
3994 complexity
3995 Well-structured training strategy with regularization techniques
to prevent overfitting

3996
3997
3998
3999
4000
4001
4002
4003
4004
4005
4006
4007
4008
4009
4010
4011
4012
4013
4014
4015
4016
4017
4018
4019
4020
4021
4022
4023
4024
4025
4026
4027
4028
4029
4030
4031
4032
4033
4034
4035
4036
4037
4038
4039
4040
4041
4042
4043
4044
4045
4046
4047
4048
4049

Attention to biological interpretability through multiple model components
Practical data preprocessing pipeline using established single-cell tools
Incorporation of expert recommendations for addressing class imbalance and improving generalization

4. AREAS FOR IMPROVEMENT
Enhance discussion of biological validation metrics beyond loss functions
Provide more explicit details on gene regulatory network construction and updating
Strengthen justification for the specific dimensions chosen for latent spaces and embeddings
Consider inclusion of additional evaluation strategies such as cross-dataset validation
Clarify how the model will handle novel cell types or contexts not present in training data

5. OVERALL RECOMMENDATION
This research plan represents a strong candidate for further development with high potential for scientific impact in single-cell perturbation prediction. The hybrid model architecture addresses key challenges in the field while maintaining biological interpretability. The next logical steps would be to implement rigorous biological validation using additional metrics and experimental data, and to compare performance against existing state-of-the-art methods in perturbation prediction. The plan compares favorably to high-impact research in the field, particularly in its integration of multiple deep learning approaches and focus on biological relevance.

P.4 DETAILED RESULTS

To comprehensively evaluate the performance of CELLFORGE, we employed five state-of-the-art LLMs as independent judges: Claude 3.7, DeepSeek-R1, OpenAI o1, Qwen-plus, and Llama 3.1. Each judge evaluated outputs from CELLFORGE and three DeepResearch variants (OpenAI, Perplexity, and Gemini) across multiple rounds to ensure statistical robustness. Tables 21 and 22 present the averaged scores from five independent evaluation runs, providing insights into both the consistency and performance differences across systems.

P.4.1 TASK ANALYSIS PHASE EVALUATION

Table 21 reveals several key insights about the Task Analysis capabilities of different systems. CELLFORGE demonstrates consistent superiority across all three evaluation dimensions, with particularly strong performance in *Analyse Dataset* (average scores: 8.60, 8.20, 7.24 across drug, gene knockout, and cytokine tasks respectively). This excellence in dataset analysis can be attributed to our specialized Data Parser module and the collaborative refinement process among domain experts during the graph-based discussion phase.

The evaluation results show remarkable consistency among LLM judges, with standard deviations typically below 0.5 points, indicating high inter-judge agreement. Notably, Claude 3.7 and OpenAI o1 tend to provide slightly higher scores overall, while Qwen-plus and Llama 3.1 exhibit more conservative scoring patterns. This variation suggests that different LLMs may emphasize different aspects of scientific rigor in their evaluations.

Among the DeepResearch variants, OpenAI’s implementation (DR^O) performs closest to CELLFORGE, achieving comparable scores in certain categories (e.g., 9.0 in drug dataset analysis). However, other methods show significant performance gaps, particularly in *Analyse Baseline Defects*, where scores drop as low as 2.16 for cytokine tasks. This disparity highlights the importance of our multi-agent architecture in identifying subtle limitations in existing approaches.

[Edit table 11 12 and caption](#)

Table 21: **LLM evaluation of the Task Analysis phase.** Three LLM judges evaluated CELLFORGE (CF), three DeepResearch (DR) [83] pipeline (O: OpenAI, P: Perplexity, G: Gemini), Biology Research Agent Biomni [46], and single LLM (CLD: Claude 3.7 API, which performs the best in our task compared with R1, o1, Qwen-plus, and Llama3.1) across four key capabilities and three perturbation types. CELLFORGE consistently matched or exceeded human expert performance, with particular strength in dataset analysis and identifying baseline model limitations.

Judges	Drug						Gene KO						Cytokine					
	CF	DR ^O	DR ^P	DR ^G	Biomni	CLD	CF	DR ^O	DR ^P	DR ^G	Biomni	CLD	CF	DR ^O	DR ^P	DR ^G	Biomni	CLD
<i>Analyse Dataset</i> ↑																		
Claude3.7	8.8	9.0	6.0	7.0	8.8	6.0	8.0	7.0	3.0	5.2	8.0	6.2	7.2	7.0	4.0	5.2	7.0	5.4
R1	9.0	9.0	6.2	7.0	8.2	5.0	8.0	7.2	4.0	6.2	7.4	5.0	7.0	6.2	3.2	5.6	7.0	4.2
o1	9.0	9.0	6.0	6.8	7.0	5.8	8.2	7.8	4.4	6.0	7.6	5.2	7.2	6.0	3.0	5.2	6.2	4.8
Qwen-plus	8.0	7.2	5.2	6.2	7.2	5.0	8.0	7.0	4.0	5.4	7.4	5.2	7.0	6.8	2.8	6.0	6.0	5.0
Llama 3.1	8.2	7.4	5.2	5.8	6.8	6.0	8.8	6.8	5.0	6.2	7.0	6.0	7.8	7.2	4.0	6.0	7.2	5.0
<i>Average</i>	8.60	8.32	5.72	6.56	7.60	5.56	8.20	7.16	4.08	5.80	7.48	5.52	7.24	6.64	3.40	5.60	6.68	4.88
<i>Analyse Task Type</i> ↑																		
Claude3.7	8.0	7.8	6.0	5.0	8.0	6.0	8.4	8.0	6.0	6.8	8.2	6.0	6.6	4.0	3.2	4.8	8.0	6.0
R1	7.0	7.0	4.0	6.0	7.0	5.0	7.8	7.0	6.0	7.0	7.0	6.0	6.6	5.8	4.8	5.0	6.6	4.0
o1	8.6	8.2	6.0	6.6	8.0	5.0	8.8	8.0	6.0	6.0	8.0	5.0	7.0	6.2	5.0	5.2	7.0	4.0
Qwen-plus	8.0	8.0	6.2	6.2	8.2	5.8	8.2	8.0	7.0	7.4	8.0	6.0	7.0	6.2	5.8	4.2	7.4	5.0
Llama 3.1	8.2	7.6	6.2	6.0	8.0	6.0	8.8	8.0	6.0	7.2	8.2	6.0	7.2	6.2	5.0	5.0	7.0	4.0
<i>Average</i>	7.96	7.72	7.80	6.20	7.84	5.56	8.40	7.80	6.20	6.88	7.88	5.56	6.88	5.96	4.76	4.84	7.20	4.60
<i>Analyse Baseline Defects</i> ↑																		
Claude3.7	6.2	5.0	3.0	4.2	5.6	3.8	6.8	6.0	3.0	4.2	6.6	3.0	5.2	4.0	2.0	2.8	5.2	2.0
R1	7.0	6.0	3.6	4.2	6.0	3.2	7.0	7.0	4.8	4.2	6.6	3.0	6.0	5.2	2.0	3.4	5.0	1.0
o1	6.6	5.8	3.0	4.0	6.0	2.8	7.2	7.0	5.6	5.0	6.2	3.0	6.0	5.2	2.0	3.2	5.0	1.6
Qwen-plus	6.0	5.2	2.2	3.2	5.0	2.0	7.0	7.0	4.2	5.0	6.4	3.0	6.0	3.8	1.8	3.0	5.2	3.0
Llama 3.1	7.0	6.2	4.0	4.8	5.0	3.0	7.2	7.6	5.0	5.8	6.0	3.0	6.8	5.2	3.0	3.0	6.4	2.0
<i>Average</i>	6.56	5.64	3.16	4.08	5.52	2.96	7.04	6.92	4.52	4.84	6.36	3.0	6.00	4.68	2.16	3.08	5.36	1.92

P.4.2 METHOD DESIGN PHASE EVALUATION

Table 22 presents a more nuanced evaluation across five dimensions of research plan quality. The results demonstrate CELLFORGE’s comprehensive superiority, with average scores exceeding 6.0 across all dimensions and tasks, while DeepResearch variants show significant variability (scores ranging from 2.16 to 8.00).

The *Innovation Level* dimension shows the most pronounced advantage for CELLFORGE, with average scores of 8.04, 8.28, and 7.44 for drug, gene knockout, and cytokine tasks, respectively. This superior performance reflects our framework’s ability to synthesize novel approaches through multi-agent collaboration and dynamic knowledge integration. Interestingly, OpenAI’s DeepResearch variant shows competitive performance in this dimension (7.40, 8.00, 7.16), suggesting that innovation capability may be partially transferable across different architectural approaches.

In *Technical Feasibility*, we observe an interesting pattern where OpenAI’s DeepResearch slightly outperforms CELLFORGE in drug perturbation tasks (7.24 vs. 6.88). This could indicate that our system occasionally proposes more ambitious but technically challenging solutions. However, CELLFORGE maintains superiority in gene knockout and cytokine tasks, demonstrating better adaptability to diverse biological contexts.

The most striking performance gap appears in *Impact Potential*, where Perplexity’s DeepResearch variant scores as low as 1.8 for drug perturbation tasks. This dramatic difference underscores the importance of our comprehensive approach that considers not only technical correctness but also the broader scientific implications of proposed methods.

Table 22: **LLM evaluation of the Method Design phase.** LLM judges assessed the quality of research plans proposed by CELLFORGE (CF) and Deep Research (DR) [83] pipeline (O: OpenAI, P: Perplexity, G: Gemini), Biology Research Agent Biomni[46], and single LLM (CLD: Claude 3.7 API, which performs the best in our task compared with R1, o1, Qwen-plus and Llama3.1) across five dimensions. CELLFORGE consistently outperformed on scientific validity, innovation, experimental design, and impact potential.

Judges	Drug					Gene KO					Cytokine							
	CF	DR ^O	DR ^P	DR ^G	Biomni	CLD	CF	DR ^O	DR ^P	DR ^G	Biomni	CLD	CF	DR ^O	DR ^P	DR ^G	Biomni	CLD
<i>Scientific Validity</i> ↑																		
Claude3.7	7.4	8.0	3.0	4.6	7.0	3.0	7.8	7.2	3.2	5.0	6.0	3.0	6.8	6.2	2.8	4.0	6.4	2.2
R1	8.2	7.4	4.0	4.8	7.0	3.2	7.8	7.8	4.0	6.2	6.8	3.2	7.0	6.0	3.6	5.8	6.0	3.0
o1	7.8	7.0	3.4	6.2	7.0	2.8	8.2	8.4	3.6	6.2	8.0	3.0	6.8	7.0	3.0	5.2	6.2	2.8
Qwen-plus	6.8	6.4	3.0	5.8	6.0	3.0	7.6	6.8	4.0	5.4	7.0	3.0	6.6	6.6	2.6	5.0	6.0	2.0
Llama 3.1	7.0	6.4	5.2	5.8	6.0	3.2	7.8	6.8	5.0	6.8	7.0	3.6	7.2	7.0	4.4	6.0	7.0	2.8
<i>Average</i>	7.44	7.04	3.72	5.44	6.60	3.04	7.84	7.40	3.96	5.92	6.96	3.16	6.88	6.56	3.28	5.20	6.32	2.56
<i>Technical Feasibility</i> ↑																		
Claude3.7	7.0	7.0	2.4	5.2	6.8	2.0	7.0	5.8	2.6	5.8	6.0	2.0	6.4	5.8	2.2	5.6	6.0	2.0
R1	5.8	7.0	4.0	5.2	6.8	2.0	6.8	6.8	5.0	6.0	6.0	2.0	6.0	5.6	4.0	5.4	6.0	1.6
o1	7.4	8.0	4.0	6.6	7.2	2.2	8.0	7.8	5.0	6.0	8.0	2.8	7.0	6.8	4.2	5.0	6.4	2.0
Qwen-plus	7.0	6.8	3.6	5.0	6.6	3.0	7.4	5.8	3.0	6.0	7.0	2.4	6.8	6.8	3.0	5.0	6.6	3.0
Llama 3.1	7.2	7.4	4.0	6.0	7.0	3.0	8.2	6.8	4.0	5.2	8.0	3.2	6.2	5.6	5.0	5.4	6.0	2.0
<i>Average</i>	6.88	7.24	3.60	5.60	6.88	2.44	7.48	6.60	3.92	5.80	7.00	2.40	6.48	6.12	3.68	5.28	6.20	2.12
<i>Innovation Level</i> ↑																		
Claude3.7	8.0	7.0	5.8	6.8	7.0	3.0	8.2	7.2	5.0	6.4	7.0	3.0	7.2	6.2	4.0	6.0	6.6	3.0
R1	8.2	6.8	4.2	4.8	7.0	3.0	8.2	7.8	5.0	7.0	7.0	3.0	8.0	7.4	5.0	6.2	8.0	2.0
o1	8.0	8.2	6.0	5.0	7.2	3.0	9.0	9.0	5.2	6.8	8.0	2.8	7.6	8.0	5.0	6.0	6.6	2.4
Qwen-plus	7.6	7.4	4.2	5.6	7.0	4.0	8.0	8.0	5.0	6.6	7.2	4.0	7.2	7.0	4.0	5.2	7.0	3.4
Llama 3.1	8.4	7.6	5.0	6.2	8.0	4.0	8.0	8.0	5.2	7.0	7.6	4.2	7.2	7.2	5.0	7.0	7.0	3.8
<i>Average</i>	8.04	7.40	5.04	5.68	7.24	3.40	8.28	8.00	5.08	6.76	7.36	3.40	7.44	7.16	4.60	6.08	7.04	2.92
<i>Experimental Design</i> ↑																		
Claude3.7	7.0	7.2	3.2	4.4	4.2	3.0	7.0	6.0	2.0	4.4	4.6	2.2	6.0	5.8	2.2	4.0	4.2	2.0
R1	8.0	7.2	4.0	4.0	5.2	3.0	7.2	6.0	2.8	4.2	5.0	2.0	6.4	6.0	3.0	4.0	4.4	1.0
o1	8.2	8.4	5.0	5.2	5.2	3.2	7.2	7.0	3.0	4.8	4.8	3.0	6.8	6.8	3.4	4.0	4.0	3.2
Qwen-plus	7.2	7.0	4.2	5.8	5.2	3.2	7.8	5.0	2.0	4.0	5.2	2.0	6.0	6.0	2.4	4.0	4.2	2.0
Llama 3.1	7.8	7.2	4.2	5.0	5.4	4.0	7.2	6.2	4.0	5.0	6.8	3.6	6.8	6.0	4.0	5.0	5.2	3.4
<i>Average</i>	7.64	7.40	4.12	4.88	5.04	3.28	7.28	6.04	2.76	4.48	5.28	2.56	6.40	5.92	3.00	4.20	4.40	2.32
<i>Impact Potential</i> ↑																		
Claude3.7	6.0	5.0	1.8	3.0	4.4	2.0	6.8	5.2	2.8	4.0	5.8	3.0	6.0	5.2	3.8	4.4	4.6	2.0
R1	7.0	6.0	2.2	4.0	4.0	1.8	7.2	6.0	2.2	4.2	4.8	2.0	6.2	5.0	2.2	4.0	5.0	2.0
o1	7.2	7.0	3.2	4.8	6.8	1.8	4.0	7.0	6.0	3.0	5.2	2.0	6.6	6.0	2.2	4.0	4.8	1.8
Qwen-plus	7.2	6.0	2.0	3.2	5.8	2.2	6.0	5.6	2.4	6.0	5.2	2.2	6.2	5.2	2.0	4.0	5.0	2.0
Llama 3.1	7.4	7.0	4.0	4.8	6.0	3.2	7.0	6.8	5.8	4.0	6.0	4.0	6.8	5.8	4.0	5.0	6.0	4.0
<i>Average</i>	6.96	6.20	2.64	3.96	5.40	2.20	6.20	6.12	3.84	4.24	5.40	2.64	6.36	5.44	2.84	4.28	5.08	2.36

4158 P.4.3 CROSS-TASK PERFORMANCE ANALYSIS
4159

4160 An interesting pattern emerges when comparing performance across different perturbation types.
4161 Gene knockout tasks generally receive the highest scores across all systems, suggesting that this well-
4162 established experimental paradigm may be easier to model computationally. In contrast, cytokine
4163 perturbation tasks show the greatest performance variance between systems, with CELLFORGE
4164 maintaining robust performance (average scores above 6.0) while some DeepResearch variants drop
4165 below 3.0 in multiple dimensions.

4166 This task-specific performance difference likely reflects the varying complexity of biological mecha-
4167 nisms involved. Gene knockouts typically produce more predictable, direct effects, while cytokine
4168 perturbations involve complex signaling cascades and cell-cell communication networks that re-
4169 quire more sophisticated modeling approaches. The superior performance of CELLFORGE in these
4170 challenging scenarios validates our multi-agent architecture’s ability to capture complex biological
4171 interactions through collaborative reasoning.

4172 P.4.4 INTER-JUDGE AGREEMENT AND RELIABILITY
4173

4174 The consistency of scores across different LLM judges provides confidence in our evaluation method-
4175 ology. The highest agreement occurs in the *Innovation Level* dimension, where judges show remark-
4176 able consensus (coefficient of variation ≤ 0.1 for most comparisons). Greater variability appears in
4177 *Experimental Design* evaluations, possibly reflecting different interpretations of what constitutes
4178 rigorous experimental validation in computational biology.

4179 These detailed results collectively demonstrate that CELLFORGE not only achieves superior perfor-
4180 mance but does so consistently across different evaluation criteria, task types, and independent judges.
4181 The framework’s ability to maintain high standards across all dimensions from technical feasibility
4182 to scientific impact underscores its potential as a comprehensive solution for automated scientific
4183 discovery in single-cell biology.
4184

4185 P.5 INTER-JUDGE CONSISTENCY ANALYSIS
4186

4187 This section provides a detailed statistical analysis to support the reliability of our LLM-as-judge
4188 evaluation methodology. We present comprehensive inter-judge consistency metrics and style-
4189 robustness testing results.

4190 To quantify inter-judge agreement, we computed two key reliability measures: Krippendorff’s for
4191 overall consistency, and Kendall’s W for ranking concordance. The following formulae explain the
4192 computational methods for each reliability measure:

4193 **Krippendorff’s Calculation.** Krippendorff’s measures the reliability of data when unitizing and
4194 coding are performed by different judges, accounting for chance agreement. For our evaluation data
4195 with n judges and m items, we compute:
4196

$$4197 \alpha = 1 - \frac{E_o}{E_e} = 1 - \frac{\frac{1}{m} \sum_{i=1}^m \sigma_i^2}{\sigma^2} \quad (8)$$

4200 where σ_i^2 is the variance of scores for item i across judges, and σ^2 is the total variance of all scores.
4201 Values above 0.8 indicate substantial agreement, while values above 0.9 indicate almost perfect
4202 agreement.
4203

4204 **Kendall’s W Calculation.** Kendall’s W (concordance coefficient) measures the agreement among
4205 judges’ rankings. For n judges ranking m items, we compute:
4206

$$4207 W = \frac{12S}{n^2(m^3 - m)} \quad (9)$$

4208 where $S = \sum_{i=1}^m (R_i - \bar{R})^2$, R_i is the sum of ranks for item i , and \bar{R} is the mean of all rank sums.
4209 W ranges from 0 (no agreement) to 1 (perfect agreement).
4210
4211

4212 Applying the above methodology to our evaluation data, we obtained the following consistency
 4213 metrics across all evaluation dimensions. Table 23 presents the inter-judge consistency metrics across
 4214 all evaluation dimensions. Krippendorff’s values above 0.8 indicate substantial agreement, while
 4215 values above 0.9 indicate almost perfect agreement.

4216
 4217 Table 23: **Inter-judge consistency metrics across evaluation dimensions.** Krippendorff’s measures
 4218 the reliability of data when unitizing and coding are performed by different judges. Values above 0.8
 4219 indicate substantial agreement, while values above 0.9 indicate almost perfect agreement. Kendall’s
 4220 W measures the concordance among judges’ rankings.

Evaluation Dimension	Krippendorff’s	Kendall’s W	Avg Correlation	Inter-Judge Var
<i>Task Analysis Phase</i>				
Analyse Dataset - Drug	0.762	0.000	0.905	0.209
Analyse Dataset - Gene KO	0.890	0.000	0.923	0.032
Analyse Dataset - Cytokine	0.885	0.000	0.941	0.088
Analyse Task Type - Drug	0.742	0.000	0.879	0.164
Analyse Baseline Defects - Drug	0.847	0.000	0.940	0.152
<i>Method Design Phase</i>				
Scientific Validity - Drug	0.871	0.000	0.911	0.044
Innovation Level - Drug	0.871	0.000	0.908	0.087
Technical Feasibility - Drug	0.915	0.000	0.966	0.113
Impact Potential - Drug	0.813	0.000	0.952	0.391
Overall Average	0.844 0.056	0.000 0.000	0.925 0.026	0.131 0.108

4234
 4235
 4236
 4237
 4238
 4239
 4240
 4241
 4242
 4243
 4244
 4245
 4246
 4247
 4248
 4249
 4250
 4251
 4252
 4253
 4254
 4255
 4256
 4257
 4258
 4259
 4260
 4261
 4262
 4263
 4264
 4265

Q HUMAN SCIENTISTS' EVALUATION DETAILS

Q.1 METHODS

To assess the scientific quality of AI-generated analysis and design outputs, we conducted a blind human evaluation involving three expert single-cell biologists. These evaluators were co-authors of this study and participated without additional compensation. Each expert independently reviewed and scored system outputs for approximately 10 hours, covering both the Task Analysis Module, the Method Design Module, and the confidence score in Graph-based discussion across cytokine, drug, and gene perturbation tasks.

For each task type (cytokine, drug, gene), experts evaluated eight outputs: 5 generated by different LLM backends of CELLFORGE (Claude 3.7, o1, DeepSeek R1, Qwen-plus, and Llama 3.1) and three from independent DeepResearch agents (OpenAI, Perplexity, Gemini). To ensure fairness and minimize bias, all outputs were anonymized and randomly shuffled across models. Experts were unaware of the model identity behind each output. Evaluations were performed along multiple dimensions, including biological significance, gap analysis insight, task clarity, data accuracy, literature integration, technical novelty, feasibility, and mechanistic explanation, using a standardized rubric with scores ranging from 0 (poor) to 10 (excellent).

Additionally, we compared human ratings with the confidence scores produced by CELLFORGE during graph-based multi-turn reasoning. Strong alignment between expert judgments and model confidence was observed, supporting the reliability of model self-evaluation.

Notably, human expert evaluations show a strong positive correlation with scores assigned by the LLM as judge, as illustrated in Figure 19.

Q.2 DETAILED RESULTS

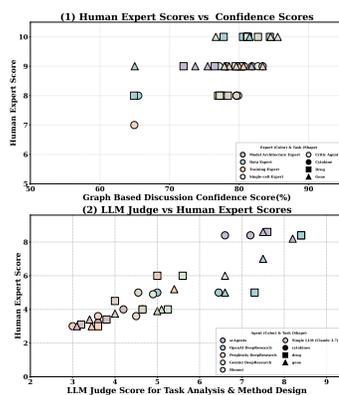


Figure 19: H Independent domain experts performed blind evaluation of research plans. Their assessments show a strong correlation with agent-generated confidence and LLM judge scores. Detailed scores are provided in Appendix P and Q.

Table 24 presents the scores given by human scientists across different tasks (cytokine, drug, and gene perturbation) for outputs from various CELLFORGE (with different LLM backends) and DeepResearch agents.

The results indicate that CELLFORGE generally outperforms DeepResearch agents, with versions like CellForge-Claude3.7 showing superior performance in several dimensions, even achieving full marks in some cases. Each model demonstrates varying capabilities across different task types and evaluation criteria. CELLFORGE versions show a more balanced performance compared to the DeepResearch agents, which sometimes score low in certain dimensions.

Table 25 compares the expert ratings with CellForge's confidence scores during graph-based multi-turn reasoning. The alignment between the confidence scores and expert ratings suggests that CellForge's self-evaluation mechanism is reliable. This correlation confirms that the confidence scores can serve as a valid indicator of the quality of CellForge's outputs. Overall, these results highlight CellForge's effectiveness in handling scientific analysis and design tasks and validate the utility of their confidence assessment.

Table 24: CELLFORGE Performance Scores Evaluated by three human scientists for 10 hours. (CF^{cl_d}: CellForge-Claude3.7, CF^{o₁}: CellForge-o1, CF^{ds}: CellForge-DeepSeek R1, CF^{qw}: CellForge-Qwen-plus, CF^{lm}: CellForge-llama 3.1, DR^O: OpenAI DeepResearch, DR^P: Perplexity DeepResearch, DR^G: Gemini DeepResearch), Biomni:Biomni[46], CLD:Single LLM Claude 3.7.

Dimension	CF ^{cl_d}	CF ^{o₁}	CF ^{ds}	CF ^{qw}	CF ^{lm}	DR ^O	DR ^P	DR ^G	Biomni	CLD
A. Analysis Reports by Task Analysis Module										
<i>Cytokines Task Analysis</i>										
Biological Significance	7	6	7	5	6	4	0	4	3	4
Gap Analysis Insight	6	2	6	4	5	2	6	2	3	1
Task Formulation Clarity	7	5	7	6	5	4	0	2	3	1
Data Characterization Accuracy	7	4	7	6	4	3	2	3	5	2
Literature Integration Quality	7	4	5	6	4	3	2	2	5	2
<i>Drug Perturbation Task</i>										
Biological Significance	7	6	7	6	5	5	3	5	3	4
Gap Analysis Insight	7	5	6	7	6	5	3	5	3	2
Task Formulation Clarity	8	7	7	6	5	3	3	3	3	2
Data Characterization Accuracy	7	8	8	6	4	3	4	3	4	4
Literature Integration Quality	8	8	8	6	5	6	4	3	4	1
<i>Gene Perturbation Task</i>										
Biological Significance	7	7	6	5	4	5	5	5	4	4
Gap Analysis Insight	7	7	5	4	5	5	0	3	3	1
Task Formulation Clarity	8	7	6	7	7	6	2	3	4	2
Data Characterization Accuracy	8	8	6	7	3	5	5	3	4	2
Literature Integration Quality	8	8	5	6	3	5	5	3	6	2
B. Hypothesis Plan by Method Design Module										
<i>Cytokines Perturbation Task</i>										
Technical Novelty	6	6	5	4	4	5	2	2	4	2
Feasibility	5	6	7	5	5	5	5	5	4	2
Clarity and Consistency	6	5	6	5	5	3	5	6	4	2
Biological Plausibility	5	6	7	4	5	5	1	3	4	2
Mechanism Explanation Quality	6	5	6	5	4	1	0	2	4	2
Pathway Relevance	5	6	6	3	4	3	0	3	4	2
Cross-Perturbation Generalizability	6	7	4	5	5	5	0	5	4	4
<i>Drug Perturbation Task</i>										
Technical Novelty	8	7	7	5	5	4	3	4	3	2
Feasibility	7	7	6	4	4	3	2	1	2	2
Clarity and Consistency	8	8	7	6	5	4	3	4	2	2
Biological Plausibility	8	8	6	5	4	4	2	3	2	2
Mechanism Explanation Quality	9	8	5	4	4	4	0	2	2	2
Pathway Relevance	9	8	5	4	3	4	0	2	4	2
Cross-Perturbation Generalizability	9	8	6	5	4	4	1	2	4	4
<i>Gene Perturbation Task</i>										
Technical Novelty	7	7	6	5	4	5	2	4	3	2
Feasibility	7	7	5	4	3	5	2	5	4	4
Clarity and Consistency	9	8	6	5	4	6	3	5	4	4
Biological Plausibility	7	7	5	4	3	4	1	4	4	3
Mechanism Explanation Quality	7	7	4	3	2	4	1	4	4	2
Pathway Relevance	7	7	4	3	2	2	0	2	4	2
Cross-Perturbation Generalizability	7	7	5	4	3	2	0	2	4	4
Overall Ranking	1	3	2	4	5	6	10	8	7	9

Table 25: Expert Human Scores compare with CellForge’s Confidence Scores on graph-based discussions Across Tasks and Rounds

Experts	Cytokine Task				Drug Task				Gene Task			
	R1	R2	R3	R4	R1	R2	R3	R4	R1	R2	R3	R4
Model Architecture Expert	8	9	9	9	9	9	10	10	9	9	9	9
Confidence Score	0.78	0.81	0.82	0.84	0.72	0.77	0.84	0.85	0.74	0.76	0.82	0.84
Data Expert	8	9	10	10	8	10	10	10	9	9	10	10
Confidence Score	0.75	0.77	0.81	0.82	0.65	0.78	0.81	0.82	0.65	0.78	0.81	0.83
Training Expert	7	8	9	9	8	9	10	10	9	9	9	9
Confidence Score	0.69	0.80	0.81	0.82	0.77	0.80	0.82	0.82	0.78	0.80	0.81	0.82
Deep Learning Expert	9	8	9	9	8	10	10	10	10	10	10	10
Confidence Score	0.74	0.77	0.85	0.88	0.79	0.82	0.86	0.87	0.79	0.85	0.87	0.88
Pathway Expert	9	8	9	9	8	10	10	10	10	10	10	10
Confidence Score	0.77	0.79	0.83	0.85	0.77	0.81	0.81	0.81	0.77	0.85	0.86	0.88
Critic Agent	8	8	9	9	8	9	10	10	9	9	9	9
Confidence Score	0.78	0.80	0.84	0.85	0.78	0.80	0.83	0.83	0.78	0.80	0.84	0.86

R PROMPT TEMPLATES

R.1 TASK DESCRIPTION INPUT

Task Description Input

Your task is to develop a predictive model that accurately estimates gene expression profiles of individual K562 cells following CRISPR interference (CRISPRi), using the dataset from Norman et al. (2019, Science).

Task Definition:

- Input: Baseline gene expression profile of an unperturbed K562 cell and the identity of the target gene(s) for perturbation
- Output: Predicted gene expression profile after perturbation

Evaluation Scenarios:

1. Unseen Perturbations: Predict effects of gene perturbations not present during training
2. Unseen Cell Contexts: Predict responses in cells with gene expression profiles not observed during training

Evaluation Metrics:

- Mean Squared Error (MSE): Measures the average squared difference between predicted and observed gene expression.
- Pearson Correlation Coefficient (PCC): Quantifies linear correlation between predicted and observed profiles.
- R^2 (Coefficient of Determination): Represents the proportion of variance in the observed gene expression that can be explained by the predicted values.
- MSE for Differentially Expressed (DE) Genes (MSE_DE): Same as MSE but computed specifically for genes identified as differentially expressed.
- PCC for Differentially Expressed (DE) Genes (PCC_DE): Same as PCC but computed specifically for genes identified as differentially expressed.
- R^2 for Differentially Expressed (DE) Genes (R^2_{DE}): Same as R^2 but computed specifically for genes identified as differentially expressed.

4428 R.2 TASK ANALYSIS COLLABORATION AGENTS SETTINGS
4429

4430 **Agent 1: Dataset Analyst** Dataset Analyst is a specialized agent responsible for performing
4431 systematic analysis of single-cell perturbation datasets during the Task Analysis stage. Its core
4432 function is to extract and summarize the key characteristics of a given dataset including experimental
4433 design, data modalities, perturbation types, and quality metrics to facilitate downstream hypothesis
4434 generation and modeling. The agent is equipped with contextualized agent retrieval (RAG in the
4435 former module) results to incorporate relevant metadata, associated publications, and protocol
4436 references. Its output follows a structured JSON format to enable direct inter-agent communication
4437 and automatic pipeline integration. The Dataset Analyst prioritizes clarity, scientific precision, and
4438 critical evaluation, identifying potential risks or biases while offering preprocessing recommendations
4439 tailored to the dataset’s complexity and intended use cases.

4440 **Task Analysis Collaboration Agent Settings**
4441

4442 **Data Analyst**
4443

4444 You are a Dataset Analyst agent in a multi-agent scientific
4445 research system. Your goal is to analyze a single-cell
4446 perturbation dataset and provide a comprehensive, structured,
4447 and insightful report to support downstream hypothesis
generation and modeling.

4448 # Role Description:

4449 The Dataset Analyst is responsible for extracting structured
4450 knowledge from single-cell perturbation datasets. This
4451 includes characterizing the experimental design, identifying
4452 quality and completeness issues, and proposing dataset-
4453 specific modeling considerations. The agent draws upon both
4454 metadata and relevant scientific context retrieved via
4455 agentic tools. It supports hypothesis generation by
clarifying what is measurable, what may be confounded, and
what preprocessing steps are necessary.

4456 # Skills:

- 4457 - Interpreting single-cell multi-omics data structures (e.g.,
4458 RNA, ATAC, protein)
- 4459 - Identifying perturbation types and their downstream modeling
implications
- 4460 - Detecting quality control issues (e.g., batch effects,
4461 sparsity, missing modalities)
- 4462 - Integrating metadata from publications, protocols, and
retrieved scientific sources
- 4463 - Structuring heterogeneous information into JSON-compatible
4464 schemas
- 4465 - Making biologically grounded recommendations for preprocessing
4466 and modeling

4467 # Objectives:

- 4468 - Extract and summarize the key characteristics of the dataset
- 4469 - Identify risks, limitations, and preprocessing needs
- 4470 - Suggest modeling strategies aligned with dataset structure
- 4471 - Provide additional scientific insights not limited to a fixed
template

4472 # Input:

4473 You will receive:

- 4474 1. Dataset metadata (e.g., species, cell types, perturbation
types)
- 4475 2. relevant knowledge from agentic retrieval

4476 # Instructions:

4477 Produce a two-part output:

4478 ## Part 1: Structured Summary (JSON format)

4479 Include the following fields, but you may expand or adapt them
4480 based on dataset complexity:
4481

```

4482 {
4483   "introduction": {
4484     "modalities": [...],
4485     "perturbation_type": [...],
4486     "conditions": [...],
4487     "timepoints": [...],
4488     "replicates": true/false,
4489     "batches": true/false,
4490     "cell_types": [...],
4491     "organism": "...",
4492     "description": "..."
4493   },
4494   "data_properties": {
4495     "num_cells": [...],
4496     "num_genes": [...],
4497     "num_features": {
4498       "RNA": ..., "ATAC": ..., "protein": ...
4499     },
4500     "perturbation_targets": {
4501       "num_unique": [...],
4502       "target_type": [...],
4503       "coverage": "dense/sparse/mixed"
4504     },
4505     "modality_completeness": [...],
4506     "metadata_completeness": [...],
4507     "preprocessing_required": [...]
4508   },
4509   "quality_assessment": {
4510     "data_sparsity": [...],
4511     "batch_effect": [...],
4512     "replicate_consistency": [...],
4513     "known_issues": [...],
4514     "strengths": [...],
4515     "limitations": [...]
4516   },
4517   "recommendations": {
4518     "preprocessing_steps": [...],
4519     "modeling_considerations": [...],
4520     "open_questions": [...]
4521   },
4522   "refinement_suggestions": [...]
4523 }
4524
4525 Write a concise, scientifically sound narrative (~150300 words)
4526 to accompany the JSON summary. This should include:
4527 - A holistic interpretation of dataset readiness for modeling
4528 - Potential scientific pitfalls or confounders
4529 - Unique strengths or opportunities (e.g., rare perturbation
4530 types, rich time series)
4531 - Reflections on whether the dataset aligns with typical
4532 assumptions in modeling pipelines
4533 - Any useful observations that do not fit cleanly into the
4534 structured fields
4535
4536 You may reference information from retrieved publications or
4537 external protocols. If information is unknown or ambiguous,
4538 state it clearly using cautious language (e.g., "not
4539 reported", "likely sparse", "appears to have").
4540
4541 # Constraints:
4542 - Be concise, accurate, and avoid redundancy
4543 - Use clear scientific language; bullet points are acceptable in
4544 the JSON
4545 - Allow flexibility in output structure if additional insights
4546 emerge
4547
4548 # Style Guide:
4549 - Output should be compatible with integration into downstream
4550 agent pipelines
4551 - Aim for the clarity and precision expected in a peer-reviewed
4552 supplementary method section
4553
4554
4555

```

- Prioritize information relevant to perturbation modeling, multi-omic integration, and biological interpretation

Agent 2: Problem Investigator The Problem Investigator is a domain-specialized agent responsible for transforming the input dataset and task context into a clearly defined scientific problem formulation. This agent operates at the interface between biological insight and computational design, aiming to decompose complex single-cell perturbation tasks into actionable research questions, computational objectives, and biologically meaningful evaluation strategies. Leveraging both LLM reasoning and agentic retrieval results, the agent integrates biological mechanisms and relevant literature to propose testable hypotheses, identify key challenges, and design analysis methods with biological and computational validity.

Task Analysis Collaboration Agent Settings

Problem Investigator

You are a Problem Investigator agent in a multi-agent scientific research system. Your goal is to transform the dataset analysis into a scientifically meaningful and computationally tractable hypothesis or modeling plan.

Role Description:

The Problem Investigator interprets dataset summaries and transforms them into well-scoped scientific problems. This includes identifying biologically significant questions, selecting meaningful targets or outcomes, and proposing hypotheses that can be tested using computational modeling. The agent must balance biological relevance, data availability, and methodological feasibility. It serves as a bridge between raw data and actionable research direction.

Skills:

- Formulating biologically meaningful and testable hypotheses from complex data
- Mapping experimental designs to machine learning problem types (e.g., classification, regression)
- Evaluating feasibility of predictive tasks based on data modality and perturbation scope
- Identifying pitfalls such as confounding, data leakage, or unobservable targets
- Specifying input-output pairs and validation schemes for modeling tasks
- Justifying scientific value and downstream utility of proposed tasks

Objectives:

- Translate dataset structure into concrete scientific questions
- Identify feasible targets, tasks, and outputs for modeling
- Justify the biological and computational value of the proposed formulation
- Propose a structured hypothesis or modeling objective for downstream agents

Input:

You will receive:

1. Dataset summary from the Dataset Analyst (structured + narrative)
2. Relevant biological context via retrieval or user input

Instructions:

Produce a structured problem specification with the following components:

```
{
  "biological_question": "string",
  "hypothesis_statement": "string",
  "task_formulation": {
    "input": ["modality1", "metadata1", "..."],
```

```

4590     "output": "target variable or prediction goal",
4591     "task_type": "regression/classification/generation/other"
4592   },
4593   "justification": {
4594     "biological_relevance": [...],
4595     "data_suitability": [...],
4596     "expected_challenges": [...]
4597   },
4598   "evaluation_plan": {
4599     "metrics": [...],
4600     "baselines_to_consider": [...],
4601     "validation_strategy": "cross-validation/held-out cells/time-
4602     split/..."
4603   },
4604   "open_questions": [...]
4605 }

```

In addition to the structured JSON, write a short explanation (100250 words) that:

- Restates the goal in accessible scientific language
- Explains why the proposed formulation is worth pursuing
- Anticipates possible modeling limitations or edge cases
- Optionally suggests alternatives or extensions

Constraints:

- Prioritize alignment with the datasets structure and perturbation resolution
- Avoid overly generic formulations; focus on specificity and tractability
- Maintain scientific rigor and make testable claims when possible

Style Guide:

- Write in the tone of a proposal for a computational biology modeling section
- Use precise language grounded in both biology and data science
- Be mindful of what is *not* observable or predictable from the dataset

Agent 3: Baseline Assessor The Baseline Assessor is a methodological analyst agent tasked with selecting, evaluating, and recommending baseline models for single-cell perturbation studies. Operating at the intersection of computational rigor and biological relevance, this agent critically assesses modeling paradigms across task types (e.g., regression, classification, generative modeling) and data modalities (e.g., gene expression, ATAC-seq, protein levels). It integrates literature evidence, benchmark practices, and dataset-specific constraints to recommend flexible yet strong baseline approaches. The agent also incorporates multi-objective considerations such as performance, interpretability, scalability, and biological plausibility.

Task Analysis Collaboration Agent Settings

Baseline Assessor

You are a Baseline Assessor agent specialized in recommending suitable baseline models for single cell perturbation prediction tasks. Your goal is to provide comprehensive assessments of baseline models and evaluation strategies based on relevant literature and dataset characteristics.

Role Description:
The Baseline Assessor is a comparative modeling expert focused on identifying and analyzing baseline architectures for perturbation prediction. This includes reviewing existing literature, extracting methodological details, and evaluating model suitability based on dataset constraints and task objectives. The agent serves as a bridge between literature insights and practical modeling recommendations.

Skills:

```

4644
4645 - Assessing baseline models for biological interpretability and
4646     computational requirements
4647 - Identifying candidate architectures relevant to perturbation
4648     types and modalities
4649 - Extracting methodological details and limitations from
4650     scientific literature
4651 - Comparing model performances across different biological
4652     contexts
4653 - Designing evaluation frameworks with biologically significant
4654     metrics
4655 - Providing actionable improvement suggestions based on
4656     technical and biological considerations
4657 # Objectives:
4658 - Review relevant literature for the given perturbation type and
4659     modality
4660 - Identify 35 candidate architectures and discuss their pros/
4661     cons in this context
4662 - Recommend at least two baseline models with rationale aligned
4663     to the dataset constraints and task objectives
4664 - Design evaluation frameworks considering biological
4665     variability and technical limitations
4666 - Provide improvement suggestions for model enhancements and
4667     biological validation
4668 # Input:
4669 You will receive:
4670 Task description and dataset information from upstream agents
4671 Retrieved papers and code implementations from literature
4672     databases
4673 RAG system results including relevant papers and code snippets
4674 # Instructions:
4675 Produce a structured analysis report with the following
4676     components:
4677 {
4678 "literature_overview": {
4679 "perturbation_types": [...],
4680 "existing_methods": [...],
4681 "technical_trends": [...]
4682 },
4683 "candidate_models": [
4684 {
4685 "model_name": "string",
4686 "architecture": "string",
4687 "strengths": [...],
4688 "weaknesses": [...],
4689 "biological_applicability": [...]
4690 }
4691 ],
4692 "recommended_baselines": [
4693 {
4694 "model_name": "string",
4695 "rationale": "string",
4696 "implementation_details": "string",
4697 "evaluation_metrics": [...],
4698 "biological_relevance": [...]
4699 }
4700 ],
4701 "evaluation_framework": {
4702 "primary_metrics": [...],
4703 "secondary_metrics": [...],
4704 "validation_strategy": "string",
4705 "test_scenarios": [...]
4706 },
4707 "improvement_suggestions": {
4708 "technical": [...],
4709 "biological": [...],
4710 "computational": [...]
4711 }
4712 }
4713 }
4714
4715
4716
4717
4718
4719
4720
4721
4722
4723
4724
4725
4726
4727
4728
4729
4730
4731
4732
4733
4734
4735
4736
4737
4738
4739
4740
4741
4742
4743
4744
4745
4746
4747
4748
4749
4750
4751
4752
4753
4754
4755
4756
4757
4758
4759
4760
4761
4762
4763
4764
4765
4766
4767
4768
4769
4770
4771
4772
4773
4774
4775
4776
4777
4778
4779
4780
4781
4782
4783
4784
4785
4786
4787
4788
4789
4790
4791
4792
4793
4794
4795
4796
4797
4798
4799
4800

```

```

4698
4699 In addition to the structured JSON, write a short explanation
4700 (100250 words) that:
4701 Summarizes the assessment of baseline models
4702 Explains the rationale behind the recommended baselines
4703 Discusses the biological relevance of the evaluation framework
4704 Anticipates potential limitations in model interpretability
4705 Suggests practical improvements for future modeling work
4706
4707 # Constraints:
4708 - Prioritize models with established track records in similar
4709 biological contexts
4710 - Computational requirements relative to dataset scale
4711 - Ensure recommended models balance biological interpretability
4712 and predictive performance
4713 - Align evaluation metrics with both technical accuracy and
4714 biological significance
4715 - Provide concrete implementation details for recommended
4716 baselines
4717
4718 # Style Guide:
4719 - Write in the tone of a modeling methodology section for a
4720 computational biology paper
4721 - Use precise language describing both technical and biological
4722 considerations
4723 - a focus on practical applicability while acknowledging
4724 theoretical limitations
4725 - Clearly distinguish between established knowledge and
4726 speculative improvements
4727 - Format model recommendations to facilitate direct
4728 implementation
4729
4730
4731
4732
4733

```

Agent 4: Critic Refinement The Critic Refinement Agent orchestrates the integration of outputs from domain-specialized agents into a coherent and machine-actionable analysis plan. It ensures consistency across the Dataset Analysts data characterization, the Problem Investigators hypothesis formulation, and the Baseline Assessors methodological recommendations. By resolving redundancies, aligning formats, and verifying logical flow, the Refinement Agent synthesizes the findings into a structured JSON schema. This agent balances standardization with flexibility, enabling downstream automation while preserving the scientific rationale of each module.

Task Analysis Collaboration Agent Settings

Critic Refinement

```

4734 You are the Refinement Agent in a multi-agent scientific
4735 research system. Your goal is to consolidate and refine
4736 outputs from the Dataset Analyst, Problem Investigator, and
4737 Baseline Assessor agents into a unified, actionable analysis.
4738
4739 # Role Description:
4740 The Refinement Agent is a meta-agent focused on cross-validating
4741 consistency across different analytical components. This
4742 includes resolving contradictions, aligning terminology, and
4743 ensuring biological relevance and technical feasibility of
4744 proposed models and evaluation frameworks. The agent serves
4745 as the integration point for all upstream analyses.
4746
4747 # Skills:
4748 Cross-validating consistency across different analytical
4749 components
4750 Resolving contradictions and aligning terminology
4751 Evaluating biological relevance and technical feasibility
4752 Structuring outputs into unified JSON schemas
4753 Generating comprehensive refinement comments
4754 Providing actionable improvement suggestions
4755
4756 # Objectives:
4757 Ensure consistency in terminologies and constraints across all
4758 outputs
4759 Align problem definitions with model assumptions
4760

```

```

4752
4753 Reorganize content into clean JSON schemas suitable for
      automated use
4754 Validate biological and technical coherence of the integrated
4755 analysis
4756 Provide final recommendations balancing biological relevance and
      technical feasibility
4757 # Input:
4758 You will receive:
4759 Analysis results from Dataset Analyst, Problem Investigator, and
      Baseline Assessor
4760 Refinement comments from previous iterations
4761 RAG system results including relevant papers and code snippets
4762 Instructions:
4763 Produce a refined analysis report with the following components:
4764 {
4765   "summary": {
4766     "biological_context": "string",
4767     "technical_requirements": "string",
4768     "refinement_overview": "string"
4769   },
4770   "task_definition": {
4771     "input_modalities": [...],
4772     "output_targets": "string",
4773     "task_type": "regression/classification/generation/other",
4774     "biological_significance": "string"
4775   },
4776   "baseline_models": {
4777     "recommended_models": [...],
4778     "model_comparisons": [...],
4779     "implementation_details": "string"
4780   },
4781   "constraints": {
4782     "dataset_limitations": [...],
4783     "technical_constraints": [...],
4784     "biological_constraints": [...]
4785   },
4786   "evaluation": {
4787     "primary_metrics": [...],
4788     "secondary_metrics": [...],
4789     "validation_strategy": "string",
4790     "test_scenarios": [...]
4791   }
4792 }
4793 In addition to the structured JSON, write a short explanation
4794 (100250 words) that:
4795 Summarizes the refinement process and key adjustments made
4796 Explains how the integrated analysis addresses biological and
4797 technical requirements
4798 Discusses remaining challenges or limitations
4799 Suggests potential extensions or future work
4800 # Constraints:
4801 Maintain consistency in terminology across all components
4802 Ensure alignment between problem formulation and model
      capabilities
4803 Validate that evaluation metrics reflect both technical accuracy
      and biological significance
4804 Provide concrete implementation details for recommended
      approaches
4805 Format outputs to facilitate direct use in downstream
      architecture design
4806 # Style Guide:
4807 Write in the tone of a methods integration section for a
      computational biology paper
4808 Use precise language describing both technical and biological
      considerations
4809 Clearly distinguish between established knowledge and
      speculative improvements
4810 Format recommendations to facilitate direct implementation
4811 Include both biological validation strategies and technical
      validation methods

```

4806 R.3 MULTI-EXPERTS SETTINGS

4807

4808

Data Expert

4809

Expert Role Setting

4810

4811

Data Expert

4812

You are acting as a **Data Engineer** in a multi-agent research critique system. Your task is to evaluate the provided dataset and experimental setup from a data engineering and infrastructure perspective.

4813

4814

4815

4816

You will receive a task analysis report that includes:

4817

- A summary of a single-cell perturbation dataset (in structured or free-text form).

4818

- The task formulation and its corresponding prediction targets.

4819

- Metadata schemas and any available preprocessing or encoding steps.

4820

- Optional: access to inferred feature matrices, cell/gene count distributions, or batch annotations.

4821

4822

Your objectives:

4823

1. Assess Data Integrity and Format

4824

- Are the cell and gene identifiers standardized and consistently used?

4825

- Is the perturbation metadata properly aligned and encoded?

4826

- Are there signs of data leakage, missing values, or corrupted entries?

4827

4828

2. Evaluate Preprocessing Pipelines

4829

- Comment on normalization, batch correction, filtering, and feature selection steps.

4830

- Are the preprocessing steps appropriate for downstream modeling?

4831

4832

3. Assess Data Scalability and Efficiency

4833

- Is the dataset efficiently stored and structured (e.g., sparse matrix, HDF5)?

4834

- Can it be easily integrated with common ML frameworks (e.g., PyTorch, TensorFlow, scikit-learn)?

4835

- Are large-scale operations (sampling, merging, batching) feasible?

4836

4837

4. Suggest Improvements or Optimizations

4838

- Recommend preprocessing adjustments, format conversions, or data storage alternatives.

4839

- Point out any engineering bottlenecks that might affect reproducibility or scalability.

4840

4841

4842

4843

Model Architecture Expert

4844

Expert Role Setting

4845

Model Architecture Expert

4846

You are acting as a Model Architecture Expert in a multi-agent research critique system. Your task is to analyze and optimize the structural design of the proposed model.

4847

You will receive a task analysis report that includes:

4848

- Task specification (input type, target prediction, expected invariances).

4849

- A baseline model description or proposed architecture diagram.

4850

- Technical constraints (e.g., compute, latency, interpretability).

4851

Your objectives:

4852

1. Deconstruct Architectural Choices

4853

- Analyze core design (e.g., encoder-decoder, attention, residuals).

4854

4855

- 4860
- 4861 - Is the architecture aligned with inductive priors from the
- 4862 data/task?
- 4863 2. Evaluate Module Interactions
- 4864 - Are modality fusions or skip connections implemented
- 4865 properly?
- 4866 - Are graph structures or latent bottlenecks justified?
- 4867 3. Spot Redundancies or Inefficiencies
- 4868 - Are there unnecessary layers, repeated computations, or
- 4869 excessive parameters?
- 4870 4. Propose Optimized Designs
- 4871 - Recommend improved architecture patterns.
- 4872 - Suggest changes that enhance expressivity, stability, or
- 4873 efficiency.

Deep Learning Expert

Expert Role Setting

Deep Learning Expert

4874

4875

4876 You are acting as a Deep Learning Expert in a multi-agent

4877 research critique system. Your task is to evaluate the model's

4878 design, scalability, and suitability for learning from

4879 high-dimensional single-cell data.

4880 You will receive a task analysis report that includes:

- 4881 - Input-output schema of the learning task (e.g., input
- 4882 modalities, targets, sample size).
- 4883 - Model class (e.g., MLP, Transformer, VAE, GNN) and
- 4884 architecture sketch.
- 4885 - Training setup including loss functions and evaluation metrics.

4886 Your objectives:

- 4887 1. Evaluate Model Suitability
- 4888 - Is the model architecture appropriate for the data type and
- 4889 task complexity?
- 4890 - Does it support integration across modalities or time
- 4891 points?
- 4892 2. Assess Scalability and Inductive Bias
- 4893 - Can the model scale with data size and sparse inputs?
- 4894 - Does it exploit structure in the data (e.g., gene graphs,
- 4895 batch embeddings)?
- 4896 3. Identify Training Bottlenecks or Risks
- 4897 - Is overfitting likely due to low data:parameter ratio?
- 4898 - Are optimization challenges (e.g., vanishing gradients,
- 4899 instability) addressed?
- 4900 4. Recommend Enhancements
- 4901 - Suggest architecture variants (e.g., regularization,
- 4902 pretraining, latent modeling).
- 4903 - Propose alternative loss designs or data augmentations

Training Expert

Expert Role Setting

Training Expert

4904

4905 You are acting as a Training Expert in a multi-agent research

4906 critique system. Your role is to critically evaluate the

4907 training strategy and optimization pipeline.

4908 You will receive a task analysis report that includes:

- 4909 - Model structure and parameter count.
- 4910 - Training procedure (e.g., optimizer, learning rate, batch size,
- 4911 scheduler).
- 4912 - Regularization strategies and data augmentation steps.

4913 Your objectives:

- 4914 1. Analyze Optimization Pipeline

- 4914
- 4915 - Are optimizers and learning rates well-tuned for the model/
task?
- 4916 - Is gradient clipping or scheduler use justified?
- 4917 2. Evaluate Regularization and Overfitting Risks
- 4918 - Are dropout, weight decay, or early stopping applied
effectively?
- 4919 - Is data augmentation sufficient and biologically reasonable
?
- 4920
- 4921 3. Diagnose Training Stability
- 4922 - Any signs of mode collapse, oscillation, or vanishing
gradients?
- 4923
- 4924 4. Recommend Training Enhancements
- 4925 - Suggest better optimizers, learning rate schedules, or
initialization schemes.
- 4926 - Propose curriculum learning or contrastive pretraining if
beneficial.
- 4927

Drug Response Expert

Expert Role Setting

Drug Response Expert

4931

4932

4933 You are acting as a Drug Response Expert in a multi-agent
scientific design system. Your role is to assess the
4934 biological and pharmacological feasibility of drug
4935 perturbation modeling based on the provided single-cell
4936 dataset and task formulation.

4937 You will receive a task analysis report that includes:

- 4938 - A summary of the perturbation dataset, including drug names,
target pathways, dosage, and timing.
- 4939 - The biological context (e.g., cell types, disease states,
4940 assay modality).
- 4941 - Task objective and prediction targets.

4942 Your objectives:

- 4943 1. Evaluate Drug Perturbation Validity
- 4944 - Are the drugs applied at biologically relevant
concentrations and durations?
- 4945 - Are the perturbations expected to have a measurable effect
4946 at the single-cell level?
- 4947 - Are there known resistance mechanisms or compensatory
pathways?
- 4948 2. Assess Target Coverage and Specificity
- 4949 - Do the drug targets align with the measured omics modality
(e.g., RNA for transcriptional drugs)?
- 4950 - Are off-target effects likely to interfere with
4951 interpretation?
- 4952 3. Recommend Improvements or Adjustments
- 4953 - Suggest better dosage choices or controls.
- 4954 - Recommend alternative compounds or combinations that better
4955 elicit the intended perturbation.

Pathway Expert

Expert Role Setting

Pathway Expert

4961

4962 You are acting as a Pathway Expert in a multi-agent biological
reasoning system. Your role is to evaluate the alignment
4963 between experimental perturbations (e.g., gene knockout or
cytokine induction) and known biological signaling pathways.

4964 You will receive a task analysis report that includes:

- 4965 - A summary of perturbation targets (genes or cytokines).
- 4966 - Downstream measurements (e.g., RNA, ATAC, surface proteins).
- 4967 - Known pathway annotations or inferred gene modules (optional).

4968
4969
4970
4971
4972
4973
4974
4975
4976
4977
4978
4979
4980
4981
4982
4983
4984
4985
4986
4987
4988
4989
4990
4991
4992
4993
4994
4995
4996
4997
4998
4999
5000
5001
5002
5003
5004
5005
5006
5007
5008
5009
5010
5011
5012
5013
5014
5015
5016
5017
5018
5019
5020
5021

Your objectives:

1. Assess Biological Plausibility
 - Does the perturbation target belong to a well-characterized signaling pathway?
 - Are expected downstream genes or modules represented in the dataset?
2. Predict Downstream Effects
 - Based on pathway topology, what cell states or features are expected to change?
 - Are these detectable in the available omics modality?
3. Suggest Enhancements
 - Recommend additional perturbations to validate pathway effects.
 - Propose experimental readouts to strengthen pathway conclusions.

Cell Communication Expert

Expert Role Setting

Cell Communication Expert

You are acting as a Cell Communication Expert in a multi-agent single-cell modeling system. Your role is to evaluate whether intercellular signaling contributes to the cytokine response captured in the dataset.

You will receive a task analysis report that includes:

- A single-cell dataset containing cytokine expression or response.
- Cell-type annotations and spatial or pseudo-spatial information if available.
- Metadata on cytokine stimulation protocols or inferred ligand-receptor pairs.

Your objectives:

1. Identify Communication Patterns
 - Are there likely paracrine or autocrine effects influencing cytokine expression?
 - Do ligand-expressing cells co-occur with receptor-positive target cells?
2. Evaluate Impact on Task
 - Could intercellular signaling confound or explain observed cytokine responses?
 - Are current assays sufficient to separate intrinsic vs. extrinsic effects?
3. Recommend Additions
 - Suggest experiments (e.g., co-culture, transwell) to isolate signaling effects.
 - Recommend including spatial transcriptomics if necessary.

Omics Modality Expert

Expert Role Setting

Omics Modality Expert

You are acting as an Omics Modality Expert in a multi-agent model evaluation system. Your role is to assess whether the chosen data modality (e.g., RNA-seq, ATAC-seq, protein) is suitable for capturing the effects of the specified perturbation.

You will receive a task analysis report that includes:

- A single-cell perturbation dataset with modality metadata.
- Task objective and prediction targets.
- Optional: known regulatory links (e.g., enhancer-promoter pairs, TF motifs, signaling cascades).

5022
5023
5024
5025
5026
5027
5028
5029
5030
5031
5032
5033
5034
5035
5036
5037
5038
5039
5040
5041
5042
5043
5044
5045
5046
5047
5048
5049
5050
5051
5052
5053
5054
5055
5056
5057
5058
5059
5060
5061
5062
5063
5064
5065
5066
5067
5068
5069
5070
5071
5072
5073
5074
5075

Your objectives:

1. Evaluate Signal Availability
 - Is the measured modality expected to show downstream effects of the perturbation?
 - Are known markers or targets captured by the modality?
2. Assess Measurement Resolution
 - Does the modality offer sufficient resolution (gene-level, peak-level, surface protein) to model the task?
3. Suggest Modality Enhancements
 - Recommend complementary modalities (e.g., ATAC + RNA) if needed.
 - Propose targeted panels or multi-omics techniques to improve interpretability.

S DETAILED OUTPUTS FROM CELLFORGE

All model implementations, source code, and training logs for the agent-designed models reported in this paper are available at https://anonymous.4open.science/r/CellForge_egs-1E54/. This repository contains the exact model implementations and dataloaders used in all reported experiments, along with a unified human-written evaluation script that was applied consistently across all methods. The repository also includes complete training logs documenting loss curves, hyperparameters, and random seeds for reproducibility, per-dataset model files and design reports that trace the agentic design process, and a unit-test suite that verifies the correctness of the perturbation-prediction objective across all models.

S.1 DATA PARSER

The input is Norman et al. [82] Dataset in h5ad format.

Data Parser Output

```

Modality:RNA Perturbation type:CRISPRa
dataset_index: filtered
Title: Exploring genetic interaction manifolds constructed from
rich single-cell phenotypes
Organisms: Homo sapiens
Modality = Data type: RNA
Method: Perturb-seq
Tissues: K562
Perturbation: CRISPRa
disease: chronic myelogenous leukemia
celltype: lymphoblasts
tissue type: cell_line
Mini-Abstract (loosely summarized original Abstract): Here, the
authors present an analytical framework for interpreting
high-dimensional landscapes of cell states (manifolds)
constructed from transcriptional phenotypes. They applied
this approach to Perturb-seq profiling of strong genetic
interactions (GIs) mined from a growth-based, gain-of-
function GI map. Exploration of this manifold enabled
ordering of regulatory pathways, principled classification
of GIs (e.g., identifying suppressors), and mechanistic
elucidation of synergistic interactions. Finally, they
applied recommender system machine learning to predict
interactions, facilitating exploration of vastly larger GI
manifolds.

``contains
                                guide_id read_count UMI_count
                                coverage gemgroup ... nperts
                                ngenes ncounts percent_mito
                                percent_ribo
TTGAACGAGACTCGGA ARID1A_NegCtrl0;ARID1A_NegCtrl0 28684 1809
15.856274 2 ... 1 3079 15097.0 5.815725 33.569583
CGTTGGGGTGTGTTGTG BCORL1_NegCtrl0;BCORL1_NegCtrl0 18367 896
20.498884 7 ... 1 2100 8551.0 4.104783 45.842592
GAACCTAAGTGTTAGA FOSB_NegCtrl0;FOSB_NegCtrl0 16296 664 24.542169
6 ... 1 2772 10999.0 5.655060 17.801618
CCTTCCTCCGTCATC SET_KLF1;SET_KLF1 16262 850 19.131765 4 ... 2
5385 38454.0 4.335050 38.165080
TCAATCTGTCTTTTCAT OSR2_NegCtrl0;OSR2_NegCtrl0 16057 1067
15.048735 2 ... 1 4869 27926.0 5.084867 32.317554
... ..
TTTGCGCAGTCATGCT RHOXF2_NegCtrl0;RHOXF2_NegCtrl0 1 1 1.000000 2
... 1 1853 5192.0 5.508475 31.798921
TTTGCGCCAGGACCCT BCL2L11_BAK1;BCL2L11_BAK1 1 1 1.000000 3 ... 2
3508 15704.0 6.718034 38.334182
TTTGCGGCTACTTGAC-1 CNN1_NegCtrl0;CNN1_NegCtrl0 1 1 1.000000 3 ...
1 3609 15054.0 5.633054 29.440680
TTTGCGCTCTGCATC-1 CEPBP_OSR2;CEBPBP_OSR2 1 1 1.000000 6 ... 2
2576 6825.0 2.695971 16.879121

```

```

5130 TTTGGTTGTTCCGTCT MAP2K3_MAP2K6;MAP2K3_MAP2K6 1 1 1.000000 2 ...
5131     2 2499 8331.0 5.617573 34.785740
5132
5133 [111445 rows x 20 columns]
5134 ___
5135 obs
5136 Index(['guide_id', 'read_count', 'UMI_count', 'coverage', '
5137     gemgroup',
5138     'good_coverage', 'number_of_cells', 'tissue_type', '
5139     cell_line',
5140     'cancer', 'disease', 'perturbation_type', 'celltype', '
5141     organism',
5142     'perturbation', 'nperts', 'ngenes', 'ncounts', '
5143     percent_mito',
5144     'percent_ribo'],
5145     dtype='object')
5146     good_coverage number_of_cells tissue_type
5147     cell_line ... perturbation_type celltype
5148     organism perturbation
5149 TTGAACGAGACTCGGA True 1 cell_line K562 ... CRISPR lymphoblasts
5150 human ARID1A
5151 CGTTGGGGTGTGTTG True 1 cell_line K562 ... CRISPR lymphoblasts
5152 human BCORL1
5153 GAACCTAAGTGTTAGA True 1 cell_line K562 ... CRISPR lymphoblasts
5154 human FOSE
5155 CCTTCCCTCCGTCATC True 1 cell_line K562 ... CRISPR lymphoblasts
5156 human SET_KLF1
5157 TCAATCTGTCTTTTCAT True 2 cell_line K562 ... CRISPR lymphoblasts
5158 human OSR2
5159 ... ..
5160 TTTGCGCAGTCATGCT False 0 cell_line K562 ... CRISPR lymphoblasts
5161 human RHOXF2
5162 TTTGCGCCAGGACCCT False 0 cell_line K562 ... CRISPR lymphoblasts
5163 human BCL2L11_BAK1
5164 TTTGCGGCTACTTGAC-1 False 0 cell_line K562 ... CRISPR
5165 lymphoblasts human CNN1
5166 TTTGCGCTCTCGCATC-1 False 0 cell_line K562 ... CRISPR
5167 lymphoblasts human CEBPB_OSR2
5168 TTTGGTTGTTCCGTCT False 0 cell_line K562 ... CRISPR lymphoblasts
5169 human MAP2K3_MAP2K6
5170
5171 [111445 rows x 10 columns]
5172 ___
5173 var
5174     ensemble_id ncounts ncells
5175 RP11-34P13.3 ENSG00000243485 29.0 29
5176 FAM138A ENSG00000237613 0.0 0
5177 OR4F5 ENSG00000186092 0.0 0
5178 RP11-34P13.7 ENSG00000238009 266.0 265
5179 RP11-34P13.8 ENSG00000239945 10.0 10
5180 ... ..
5181 AC233755.2 ENSG00000277856 0.0 0
5182 AC233755.1 ENSG00000275063 0.0 0
5183 AC240274.1 ENSG00000271254 11735.0 10835
5184 AC213203.1 ENSG00000277475 0.0 0
5185 FAM231B ENSG00000268674 0.0 0
5186
5187 [33694 rows x 3 columns]
5188 ___
5189 shape
5190 (111445, 33694)
5191 \\\
5192
5193
5194
5195
5196
5197
5198
5199
5200
5201

```

5184 S.2 AGENTIC RETRIEVAL
5185
5186

5187 **Task Analysis**

5188 **Example Agentic Retrieval Output**

5189 Keywords: Norman Weissman 2019 Perturb-seq CRISPRa K562 single
5190 cell perturbation prediction
5191

5192 Round 1: Initial DFS Search (one example branch)
5193 Keywords: single cell perturbation prediction (what is -> how to
5194 solve)
5195 Learning single-cell perturbation responses using neural optimal
5196 transport...{Nature Link}
5197 Modeling and predicting single-cell multi-gene perturbation
5198 responses...{PMC Link}
5199 Explainable modeling of single-cell perturbation data using
5200 Bayesian hierarchical modeling...{Cell Press Link}
5201 A Multiplexed Single-Cell CRISPR Screening...{Cell Press Link}
5202 Predicting transcriptional outcomes of novel perturbations...{
5203 PMC Link}
5204 Modeling and predicting single-cell multi-gene perturbation
5205 responses...{PMC Link}
5206 Exploring genetic interaction manifolds...{PMC Link}
5207 Predicting transcriptional outcomes of novel perturbations...{
5208 PMC Link}
5209 In-silico biological discovery with large-scale perturbation
5210 data{arXiv Link}
5211 DeepChrome 2.0: Investigating and modeling chromatin
5212 accessibility...{arXiv Link}
5213 Predicting the genetic component of gene... {arXiv Link}
5214 A genome-scale deep learning model to pre...{arXiv Link}
5215 Attention-based Interpretable Regression...{arXiv Link}
5216 GATES: Graph Network...{arXiv Link}
5217 GRNFormer: Biologically...{arXiv Link}

5218 Round 2: BFS Search (one example branch)
5219 Key words:state-of-the-art models for single-cell perturbation
5220 prediction
5221 GEARS: Predicting transcriptional outcomes of novel multi-gene
5222 perturbations{nature Link}
5223 scGPT: Is language all you need for modeling single-cell
5224 perturbation responses{nature Link}
5225 Geneformer - BioNeMo Framework for Genomic Language Modeling{
5226 nature link}
5227 Efficient Fine-Tuning of Single-Cell Foundation Models for
5228 Perturbation Prediction...{arXiv Link}
5229 Multicell-Fold: geometric learning in folding...{PMC Link}
5230 Variational Mixtures of ODEs for Inferring Cellular...{ICML}
5231 DeepChrome 2.0: Investigating and modeling chromatin
5232 accessibility...{arXiv Link}
5233 Predicting the genetic component of gene...{Bioinformatics}
5234 A genome-scale deep learning model to pre...{PMC Link}

5235 Round 3: DFS Search (one example branch)
5236 Keywords: GEARSscGPTGeneformerMulticell-FoldDeepChrome 2.0...
5237 GEARS: Predicting transcriptional outcomes of novel multi-gene
5238 perturbations{nature Link}
5239 Predicting transcriptional outcomes of novel multigene
5240 perturbations with GEARS...{nature Link}
5241 snap-stanford/GEARS {Github Link}
5242 scGPT: Is language all you need for modeling single-cell
5243 perturbation responses{nature Link}
5244 bowang-lab/scGPT{Github Link}
5245 Geneformer - BioNeMo Framework for Genomic Language Modeling{
5246 nature link}
5247 jkobject/geneformer {Github Link}
5248 Multicell-Fold: geometric learning in folding...{PMC Link}
5249 bm2-lab/scPerturBench{Github Link}

5238
5239
5240
5241
5242
5243
5244
5245
5246
5247
5248
5249
5250
5251
5252
5253
5254
5255
5256
5257
5258
5259
5260
5261
5262
5263
5264
5265
5266
5267
5268
5269
5270
5271
5272
5273
5274
5275
5276
5277
5278
5279
5280
5281
5282
5283
5284
5285
5286
5287
5288
5289
5290
5291

DeepChrome 2.0: Investigating and modeling chromatin accessibility...{arXiv Link}
DeepChrome/DeepChrome2.0{Github Link}
.....
Round 4: BFS Search (one example branch)
Keywords: TransformerVAE... (used by models for single-cell perturbation prediction)
TransVAE: Giving Attention to Generative VAE Models for De Novo Molecular Design{https://github.com/oriondollar/TransVAE}{GitHub Link}
Transformer Graph Variational Autoencoder for Generative Molecular Design{https://www.biorxiv.org/content/10.1101/2024.07.22.604603v1.full.pdf}{Paper Link}
A novel molecule generative model of VAE combined with Transformer{https://arxiv.org/pdf/2402.11950}{arxiv Link}
A Graph VAE and Graph Transformer Approach to Generating{https://arxiv.org/abs/2104.04345}{arxiv Link}
Large Scale GNN and Transformer Models and for Genomics{https://qdata.github.io/deep2Read//talks-mb2019/Jack_20190530-LargeScaleGNNTransformer.pdf}
Round 5: DFS Search
.....

5292 S.3 TASK ANALYSIS REPORT
5293

5294 **Data Analyst Output**

5295 **Task Analysis**
5296

5297
5298 ### Introduciton
5299
5300 - **Modalities** : RNA (single-cell RNA-seq)
5301 - **Perturbation Type** : CRISPRa (activation)
5302 - **Conditions** : 1,092 unique perturbation conditions including:
5303 - 105 single-gene perturbations
5304 - 131 paired-gene perturbations
5305 - **Timepoints** : Not explicitly reported (likely single
5306 timepoint)
5307 - **Replicates** : Present (multiple cells per perturbation
5308 condition)
5309 - **Batches** : Likely present due to sequencing logistics
5310 - **Cell Types** : K562 cells (chronic myelogenous leukemia cell
5311 line)
5312 - **Organism** : Homo sapiens
5313 - **Description** : This Perturb-seq experiment profiles
5314 transcriptional responses in K562 cells following CRISPRa
5315 perturbations. The dataset captures both single and paired
5316 genetic perturbations, providing insight into genetic
5317 interaction networks.

5318 ### Data Properties
5319
5320 - **Number of Cells** : 111,445
5321 - **Number of Genes** : 33,694
5322 - **Features** :
5323 - **RNA** : 33,694 genes
5324 - **ATAC** : Not available
5325 - **Protein** : Not available
5326 - **Perturbation Targets** :
5327 - **Unique Perturbations** : 105 single genes and 131 gene pairs
5328 - **Target Type** : Genes and gene pairs
5329 - **Coverage** : Sparse (many genes show no expression change)
5330 - **Modality Completeness** : RNA data available for all cells
5331 with appropriate quality control metrics
5332 - **Metadata Completeness** : Comprehensive perturbation
5333 annotations, technical covariates (UMI_count, percent_mito,
5334 percent_ribo), and cell barcodes
5335 - **Preprocessing Required** :
5336 - Normalization of UMI counts
5337 - Log-transformation of expression values
5338 - Filtering of lowly expressed genes
5339 - Batch effect correction

5340 ### Quality Assessment
5341
5342 - **Data Sparsity** : 78% of gene expression values are zero due
5343 to dropout events and biological sparsity
5344 - **Batch Effect** : Potential technical variability from
5345 different sequencing runs
5346 - **Replicate Consistency** : Not explicitly reported but
5347 suggested by reproducibility of common perturbations
5348 - **Known Issues** :
5349 - Extreme class imbalance (some perturbations appear only once)
5350 - High dimensionality with 33,694 genes per output
5351 - Technical noise inherent to single-cell sequencing
5352 - **Strengths** :
5353 - Rich perturbation annotations with precise gene targeting
5354 - Large number of profiled cells providing statistical power
5355 - Comprehensive technical covariates enabling noise correction
5356 - **Limitations** :
5357 - High dimensionality challenges
5358 - Sparse perturbation distribution

```

5346     - Potential confounding from clonal expansion
5347
5348     ### Recommendations
5349
5350     - Preprocessing Steps:
5351     - Filter genes with fewer than 10 counts across all cells
5352     - Normalize UMI counts using library size factors
5353     - Apply log-transformation (log1p) to expression values
5354     - Combat batch effects using Harmony or Scanorama
5355     - Modeling Considerations:
5356     - Focus modeling efforts on the top 5,000-10,000 variable genes
5357     - Use sparse-aware models (e.g., sparse neural networks)
5358     - Incorporate technical covariates as model inputs
5359     - Implement data augmentation for rare perturbations
5360     - Open Questions:
5361     - Optimal representation of paired perturbations in model
5362       architecture
5363     - Strategies for handling extreme class imbalance in training
5364     - Best practices for integrating technical covariates
5365
5366     ### Refinement Suggestions
5367
5368     - Investigate the impact of clonal expansion on expression
5369       profiles
5370     - Assess the quality of rare perturbation data for potential
5371       filtering
5372     - Explore data augmentation techniques for underrepresented
5373       perturbations
5374     - Develop quality control metrics specific to Perturb-seq data

```

Problem Investigator Output

```

5370     Task Analysis
5371
5372
5373     ### Biological Question
5374
5375     How do genetic perturbations propagate through gene regulatory
5376     networks to alter the transcriptional landscape of K562
5377     cells, and can we predict these changes for novel
5378     perturbations and cellular contexts?
5379
5380     ### Hypothesis Statement
5381
5382     A predictive model can accurately estimate post-perturbation
5383     gene expression profiles by learning the complex
5384     relationships between baseline cell states, genetic
5385     perturbations, and transcriptional responses, including non-
5386     linear genetic interaction effects.
5387
5388     ### Task Formulation
5389
5390     - Input:
5391     - Baseline gene expression profile (33,694 genes)
5392     - Perturbation identities (single or paired genes)
5393     - Technical covariates (UMI_count, percent_mito, percent_ribo)
5394     - Output:
5395     - Post-perturbation gene expression profile (33,694 genes)
5396     - Task Type:
5397     - High-dimensional regression with combinatorial inputs
5398
5399     ### Justification
5400
5401     - Biological Relevance:
5402     - Understanding genetic interactions is fundamental to
5403       deciphering cellular response networks
5404     - Mapping genotype-phenotype relationships at single-cell
5405       resolution
5406     - Predicting cellular responses to novel perturbations
5407       accelerates functional genomics research
5408     - Data Suitability:

```

```

5400
5401 - Rich perturbation annotations enable supervised learning
5402 approaches
5403 - Single-cell resolution captures heterogeneity in cellular
5404 responses
5405 - Coverage of both single and paired perturbations allows study
5406 of genetic interactions
5407 - **Expected Challenges**:
5408 - High-dimensional output space with 33,694 genes per
5409 prediction
5410 - Non-linear genetic interactions requiring complex model
5411 architectures
5412 - Generalization to unseen perturbations and cellular contexts
5413 - Technical noise and dropout events in single-cell data
5414
5415 ### Evaluation Plan
5416
5417 - **Metrics**:
5418 - Pearson Correlation Coefficient (PCC) across all genes and on
5419 top 1,000 differentially expressed genes (DE)
5420 - Mean Squared Error (MSE) across all genes and for
5421 differentially expressed genes
5422 - R score across all genes and for differentially expressed
5423 genes
5424 - **Baselines to Consider**:
5425 - Linear regression models with gene-level outputs
5426 - Gene-wise neural networks
5427 - Graph neural networks using known regulatory interactions
5428 - **Validation Strategy**:
5429 - Stratified group k-fold cross-validation holding out entire
5430 perturbation conditions
5431 - Separate validation sets for unseen single-gene and paired
5432 perturbations
5433 - Baseline profile holdout sets to test generalization to new
5434 cellular contexts
5435
5436 ### Open Questions
5437
5438 - Optimal representation of paired perturbations in model
5439 architecture
5440 - Strategies for handling extreme class imbalance during
5441 training
5442 - Incorporation of technical covariates into model architecture
5443 - Best practices for defining differentially expressed genes in
5444 this context
5445

```

Baseline Assessor Output

```

5436 Task Analysis
5437
5438
5439 ### Literature Overview
5440
5441 - **Perturbation Types**:
5442 - CRISPRi (interference)
5443 - CRISPRa (activation)
5444 - **Existing Methods**:
5445 - GEARS: Graph neural network with gene interaction modeling
5446 - scGPT: Transformer architecture for single-cell data
5447 - Geneformer: Pretrained transformer for genomics
5448 - scLAMBDA: Variational autoencoder for perturbation prediction
5449 - **VAE (Variational Autoencoder)**: Probabilistic model for
5450 learning latent representations of cellular states
5451 - **Technical Trends**:
5452 - Transformer architectures for capturing long-range
5453 dependencies
5454 - Graph neural networks for explicit gene interaction modeling
5455 - Variational autoencoders for probabilistic modeling of
5456 cellular states
5457 - Hybrid models combining multiple data modalities

```

```

5454     - Deep generative models for data augmentation and uncertainty
5455         quantification
5456
5457     ### Candidate Models
5458     ##### GEARS (Gene Network Embedding for Perturbation Response
5459         Prediction)
5460     - Architecture: Graph Neural Network (GNN) combined with
5461         Multi-Layer Perceptron (MLP)
5462     - Strengths:
5463         - Explicitly models gene dependencies using known regulatory
5464           interactions
5465         - Handles combinatorial perturbations through graph propagation
5466         - Demonstrated success in previous Perturb-seq challenges
5467     - Weaknesses:
5468         - Relies on external gene interaction databases
5469         - May overfit to common perturbations with limited
5470           generalization
5471         - Computationally intensive for full transcriptome modeling
5472     - Biological Applicability:
5473         - Captures genetic interactions and regulatory relationships
5474         - Models enhancer-promoter relationships in K562 cells
5475         - Provides interpretable gene importance scores
5476     ##### scGPT (Single-Cell Generative Perturbation Transformer)
5477     - Architecture: Transformer with multi-head self-attention
5478     - Strengths:
5479         - Captures long-range gene interactions without relying on
5480           external databases
5481         - Robust to technical noise through attention mechanisms
5482         - Handles variable numbers of perturbations naturally
5483     - Weaknesses:
5484         - Requires extensive pretraining on large datasets
5485         - Computationally demanding for full transcriptome modeling
5486         - May struggle with extreme class imbalance
5487     - Biological Applicability:
5488         - Models context-dependent transcriptional responses
5489         - Handles sparse data efficiently through attention mechanisms
5490         - Provides gene importance scores through attention weights
5491     ##### Enformer (Enhancer former)
5492     - Architecture: Dilated Convolutional Neural Network (CNN)
5493     - Strengths:
5494         - Effective at modeling sequence-to-expression relationships
5495         - Provides interpretable feature importance scores
5496         - Computationally efficient compared to transformer
5497           architectures
5498     - Weaknesses:
5499         - Requires DNA sequence input not directly applicable to post-
5500           transcriptional perturbations
5501         - Limited ability to model combinatorial genetic effects
5502         - Not designed for single-cell data with technical covariates
5503     - Biological Applicability:
5504         - Predicts expression changes from DNA sequence modifications
5505         - Limited utility for CRISPRa perturbations affecting post-
5506           transcriptional regulation
5507
5508     ### Recommended Baselines
5509     ##### Graph Neural Network (GNN) with Gene Interaction Modeling
5510     - Rationale: Explicitly models gene dependencies and can
5511         incorporate known regulatory interactions while remaining
5512         flexible to learn from data
5513     - Implementation Details:
5514         - Use PyTorch Geometric for efficient graph operations
5515         - Construct gene interaction graphs from public databases (e.g
5516           ., STRING, BioGRID)
5517

```

```

5508
5509 - Implement separate graph branches for regulatory and co-
5510 expression relationships
5511 - Include attention mechanisms to weight different interaction
5512 types
5513 - Embed perturbation identities using learned gene embeddings
5514 - Concatenate baseline expression features with perturbation
5515 embeddings
5516 - Apply multiple GNN layers followed by dense layers for
5517 prediction
5518 - Evaluation Metrics: PCC, MSE, R
5519 - Biological Relevance: Captures genetic interactions and
5520 regulatory mechanisms, providing insight into how
5521 perturbations propagate through networks
5522
5523 ##### Transformer Architecture with Gene Positional Encoding
5524
5525 - Rationale: Capable of discovering complex gene
5526 interactions without relying on external databases, with
5527 architectural flexibility for different input modalities
5528 - Implementation Details:
5529 - Use PyTorch with Hugging Face transformer libraries
5530 - Encode genes as positional tokens with expression values
5531 - Implement specialized embeddings for perturbed genes
5532 - Apply layer normalization and residual connections
5533 - Use mixed precision training to handle large output
5534 dimensions
5535 - Implement masking for rare perturbations during training
5536 - Apply attention pooling to focus on biologically relevant
5537 genes
5538 - Evaluation Metrics: PCC, MSE, R
5539 - Biological Relevance: Models context-dependent responses
5540 and technical noise robustly, providing flexibility to adapt
5541 to different biological questions
5542
5543 ##### VAE (Variational Autoencoder)
5544
5545 - Architecture: Encoder-decoder architecture with
5546 probabilistic latent space
5547 - Strengths:
5548 - Models uncertainty in cellular states and perturbation
5549 responses
5550 - Effective for data augmentation through generation of new
5551 cellular states
5552 - Provides compressed latent representations for downstream
5553 analysis
5554 - Handles sparse and noisy single-cell data well
5555 - Weaknesses:
5556 - May oversimplify complex biological relationships in latent
5557 space
5558 - Requires careful tuning of KL divergence weighting
5559 - Potential blurring of distinct cellular states in latent
5560 space
5561 - Biological Applicability:
5562 - Captures multimodal distributions of cellular responses
5563 - Enables exploration of cellular state transitions following
5564 perturbations
5565 - Provides robust representations for classifying cellular
5566 phenotypes
5567
5568 VAE with Perturbation Conditioning:
5569 - Rationale: Models uncertainty in cellular responses and
5570 provides robust latent representations for downstream
5571 analysis while enabling data augmentation
5572 - Implementation Details:
5573 - Use PyTorch for flexible probabilistic modeling
5574 - Implement encoder-decoder architecture with probabilistic
5575 latent space
5576 - Include perturbation identities as conditional inputs to the
5577 decoder
5578 - Apply beta-VAE regularization to balance reconstruction and
5579 latent space regularization
5580
5581

```

```

5562
5563 - Implement sparse VAE modifications to handle zero-valued
5564 genes
5565 - Use importance weighting for rare perturbations during
5566 training
5567 - Apply latent space interpolation to explore cellular state
5568 transitions
5569 - Evaluation Metrics: PCC, MSE, R, ELBO (Evidence Lower
5570 Bound)
5571 - Biological Relevance: Captures multimodal distributions of
cellular responses, enables exploration of cellular state
transitions, and provides robust representations for
classifying cellular phenotypes
5572
5573 ### Evaluation Framework
5574
5575 - Primary Metrics:
5576 - PCC for differentially expressed genes
5577 - MSE for differentially expressed genes
5578 - R for differentially expressed genes
5579 - Secondary Metrics:
5580 - Global PCC across all genes
5581 - Global MSE across all genes
5582 - Classification accuracy for key gene up/down regulation
5583 - ELBO for VAE models
5584 - Validation Strategy: Stratified group k-fold cross-
validation holding out entire perturbation conditions
5585 - Test Scenarios:
5586 - Unseen single-gene perturbations
5587 - Unseen paired perturbations with one known gene
5588 - Completely novel gene pairs
5589 - Baseline profiles from held-out cell subsets
5590
5591 ### Improvement Suggestions
5592
5593 - Technical:
5594 - Incorporate attention mechanisms to focus on biologically
5595 relevant genes
5596 - Apply data augmentation for rare perturbations (e.g., mixup,
5597 elastic deformation)
5598 - Use mixed precision training to reduce memory requirements
5599 - Implement gradient checkpointing for deeper architectures
5600 - Experiment with different VAE architectures (e.g., VAE-GCN,
5601 VAE-Transformer)
5602 - Biological:
5603 - Integrate known regulatory interactions as model priors
5604 - Focus modeling efforts on biologically relevant gene subsets
5605 - Model cell cycle effects using phase markers
5606 - Incorporate chromatin accessibility information when
5607 available
5608 - Use VAE latent space for phenotypic classification and
5609 clustering
5610 - Computational:
5611 - Optimize for sparse operations to handle zero-valued genes
5612 efficiently
5613 - Implement efficient parallel processing strategies
5614 - Use model pruning to reduce parameters while maintaining
5615 performance
5616 - Apply distributed training for VAE models with large latent
spaces

```

Critic Refinement Output

Task Analysis

```

5612 ### Summary
5613
5614 - Biological Context: Predicting transcriptional responses
5615 to CRISPR perturbations in K562 cells requires modeling

```

```

5616
5617     complex genetic interactions and cellular context
5618     dependencies while generalizing to novel perturbations
5619 - Technical Requirements: Handling high-dimensional sparse
5620     data, extreme class imbalance, and the need for inductive
5621     generalization
5622 - Refinement Overview: Integration of dataset analysis,
5623     problem formulation, and baseline recommendations into a
5624     cohesive analysis plan with aligned terminology and
5625     evaluation metrics
5626
5627 Task Definition
5628
5629 - Input Modalities:
5630     - Baseline gene expression profile (33,694 genes)
5631     - Perturbation identities (single or paired genes)
5632     - Technical covariates (UMI_count, percent_mito, percent_ribo)
5633 - Output Targets: Post-perturbation gene expression profile
5634     (33,694 genes)
5635 - Task Type: High-dimensional regression with combinatorial
5636     inputs
5637 - Biological Significance: Enables understanding of genetic
5638     interaction networks and prediction of cellular responses to
5639     novel perturbations
5640
5641 Baseline Models
5642
5643 - Recommended Models:
5644     - Graph Neural Network (GNN)
5645     - Transformer Architecture
5646     - VAE with Perturbation Conditioning
5647 - Model Comparisons:
5648     - GNNs excel at explicit gene interaction modeling using known
5649     regulatory networks
5650     - Transformers offer flexible interaction discovery without
5651     relying on external databases
5652     - VAEs provide probabilistic modeling of cellular states and
5653     enable data augmentation
5654 - Implementation Details:
5655     - Implement GNNs with PyTorch Geometric using gene interaction
5656     graphs derived from prior knowledge
5657     - Implement Transformers with PyTorch/Hugging Face using gene
5658     positional encoding
5659     - Implement VAEs with PyTorch using conditional latent spaces
5660     for perturbation modeling
5661     - Include specialized embeddings for perturbations and
5662     normalize technical covariates
5663     - Apply mixed precision training and gradient checkpointing for
5664     efficiency
5665
5666 Constraints
5667
5668 - Dataset Limitations:
5669     - Class imbalance with rare perturbations appearing only once
5670     - Data sparsity with 78% zero-valued genes
5671     - Potential batch effects from different sequencing runs
5672 - Technical Constraints:
5673     - Computational resources for training large models on full
5674     transcriptome data
5675     - Model interpretability requirements for biological validation
5676 - Biological Constraints:
5677     - Need for generalization to unseen perturbations and cellular
5678     contexts
5679     - Model must align with known regulatory mechanisms where
5680     possible
5681     - Focus on biologically relevant gene subsets to avoid
5682     overfitting
5683
5684 Evaluation
5685
5686 - Primary Metrics:
5687
5688
5689

```

```

5670
5671 - PCC across all genes and for differentially expressed genes (
5672   top 1,000)
5673 - MSE across all genes and for differentially expressed genes
5674 - R across all genes andfor differentially expressed genes
5675
5676 - **Validation Strategy**: Stratified group k-fold cross-
5677   validation holding out entire perturbation conditions
5678 - **Test Scenarios**:
5679   - Unseen single-gene perturbations
5680   - Unseen paired perturbations with one known gene
5681   - Completely novel gene pairs
5682   - Baseline profiles from held-out cell subsets

```

Final Report

Task Analysis

```

5683
5684
5685 Task Analysis
5686 # Biological Objective
5687 Predict post-perturbation gene expression profiles in K562 cells
5688   to understand genetic interaction networks and enable
5689   discovery of novel regulatory mechanisms. This work aims to
5690   develop a computational tool for exploring genetic
5691   interaction manifolds, accelerating functional genomics
5692   research and therapeutic target discovery.
5693 # Technical Approach
5694 Develop high-dimensional regression models incorporating
5695   baseline expression, perturbation identities, and technical
5696   covariates. The models must explicitly handle sparse data,
5697   extreme class imbalance, and demonstrate inductive
5698   generalization to novel perturbations and cellular contexts.
5699   The plan includes exploration of deterministic models (GNNs,
5700   Transformers) and probabilistic models (VAEs) to capture
5701   different aspects of cellular response variability.
5702 # Dataset Characterization
5703 ## Origin
5704 Norman et al. (2019) Perturb-seq dataset (GEO: GSE133344).
5705 Key Features
5706 Number of Cells: 111,445
5707 Number of Genes: 33,694
5708 Perturbation Conditions: 1,092 unique conditions (105 single
5709   genes, 131 gene pairs)
5710 Technical Covariates: UMI_count, percent_mito, percent_ribo
5711 ## Challenges
5712 Class Imbalance: Rare perturbations appear only once.
5713 Data Sparsity: 78% zero-valued genes due to dropout events.
5714 Technical Noise: Inherent to single-cell sequencing.
5715 Batch Effects: Potential variability from different sequencing
5716   runs.
5717 # Problem Formulation
5718 ## Biological Question
5719 How do genetic perturbations propagate through gene regulatory
5720   networks to alter the transcriptional landscape of K562
5721   cells, and can we predict these changes for novel
5722   perturbations and cellular contexts?
5723 ## Hypothesis Statement
5724 A predictive model can accurately estimate post-perturbation
5725   gene expression profiles by learning the complex
5726   relationships between baseline cell states, genetic
5727   perturbations, and transcriptional responses, including non-
5728   linear genetic interaction effects.
5729 ## Task Definition

```

```

5724
5725   ### Input:
5726   Baseline gene expression profile (33,694 genes)
5727   Perturbation identities (single or paired genes)
5728   Technical covariates (UMI_count, percent_mito, percent_ribo)
5729
5730   ### Output:
5731   Post-perturbation gene expression profile (33,694 genes)
5732
5733   ###Task Type:
5734   High-dimensional regression with combinatorial inputs
5735
5736   ## Justification
5737   Biological Relevance:
5738   Understanding genetic interactions is fundamental to deciphering
5739   cellular response networks
5740   Mapping genotype-phenotype relationships at single-cell
5741   resolution
5742   Predicting cellular responses to novel perturbations accelerates
5743   functional genomics research
5744   Data Suitability:
5745   Rich perturbation annotations enable supervised learning
5746   approaches
5747   Single-cell resolution captures heterogeneity in cellular
5748   responses
5749   Coverage of both single and paired perturbations allows study of
5750   genetic interactions
5751   Expected Challenges:
5752   High-dimensional output space with 33,694 genes per prediction
5753   Non-linear genetic interactions requiring complex model
5754   architectures
5755   Generalization to unseen perturbations and cellular contexts
5756   Technical noise and dropout events in single-cell data
5757
5758   The model's performance will be evaluated under two key
5759   scenarios to assess its generalizability:
5760   - Unseen Perturbations: The model should be able to accurately
5761   predict the effects of CRISPRi targeting genes or gene pairs
5762   that were not included in the training data. This scenario
5763   tests the model's ability to extrapolate its learned
5764   knowledge to novel genetic manipulations.
5765   - Unseen Cell Contexts: The model should be capable of
5766   predicting the response to a perturbation in cells with
5767   baseline gene expression profiles that were not observed
5768   during the training phase. This evaluates the model's
5769   robustness to the inherent heterogeneity within the K562
5770   cell population.
5771
5772   # Baseline Model Analysis
5773
5774   **SOTA**: GEARS achieves best Pearson correlation in
5775   combinatorial prediction tasks but violates the "no external
5776   database" constraint .
5777
5778   Below are detailed critiques of each baselines shortcomings in
5779   the context of the AdamsonWeissman UPR CRISPRi dataset,
5780   followed by concrete recommendationsgrounded in recent
5781   literaturefor how to overcome them. Each point is supported
5782   by highquality citations.
5783
5784   1. SC-GPT
5785
5786   **Shortcomings:**
5787
5788   1). **Discrete Perturbation Tokens:** SC-GPT treats each
5789   perturbation (e.g. a specific dualguide combination) as a
5790   unique token. It cannot form embeddings for guide sets
5791   unseen in pretraining, so it fails on novel combinations
5792
5793   2). **No Zero-Inflated Modeling:** SC-GPTs Gaussian or cross-
5794   entropy losses dont account for dropoutdriven zeros common
5795   in scRNA-seq, causing biased predictions for low-UMI cells
5796
5797

```

5778
5779 3). **Parameter Bloat for Dense Output:** Extending SC-GPTs
5780 languagemodel head to 35 kdimensional gene outputs inflates
5781 parameters, hindering training efficiency and generalization

5782 2. GeneFormer
5783 **Shortcomings:**

5784 1). **Single-Gene Focus:** GeneFormer has been validated
5785 primarily on singlegene knockouts, lacking mechanisms to **compose**
5786 multiple guide embeddings for combinatorial
5787 CRISPRi

5788 2). **Static Graph Priors:** It uses a fixed genegene network
5789 that doesnt adapt to perturbationinduced regulatory rewiring
5790 in the UPR pathway, limiting dynamic response modeling

5791 3). **Scalability Issues:** Fullgraph attention over 35 k genes
5792 is intractable, so practical implementations subsample to 25
5793 k genesdiscarding potentially important UPR regulators

5794 3. DEEP (Plain MLP)
5795 **Shortcomings:**

5796 1). **Ignores Gene Covariance:** Treats each gene independently,
5797 missing co-regulation patterns (e.g., ATF6XBP1 axis in UPR)

5798 2). **Overfitting Risk:** Millions of parameters on 35 k inputs
5799 with limited replicates per combination leads to
5800 memorization, not generalization to unseen guide sets

5801 3). **No Interpretability:** Provides no insight into which
5802 genes or interactions drive predictions, unlike graph-based
5803 or attention-based models.

5804 4. GEARS
5805 **Shortcomings:**

5806 1). **External Knowledge Dependency:** GEARS integrates a
5807 genegene memory module (e.g., from STRING or GO) to
5808 regularize embeddings, which violates our no external
5809 database constraint

5810 2). **Fixed Graph Structure:** The perturbation relationship
5811 graph in GEARS is static, not conditioned on cell-state or
5812 UPR context, limiting dynamic response capture.

5813 3). **Heavy GNN Overhead:** Graph neural network message passing
5814 on 35 k nodes x multiple perturbations incurs high memory
5815 and compute costs, impractical for largescale CRISPRi
5816 screens.

5817 ---

5818 Recommendations for Improvement

5819 1. **Factorized Perturbation Embeddings**

5820 - **Approach:**
5821 Learn a separate embedding $e_{\{g\}}$ for each guide g .
5822 Represent a perturbation set \mathcal{P} by a **learned nonlinear**
5823 **composition**.

5824 - **Benefit:** Zeroshot support for unseen guide combinations
5825 via embedding arithmetic, as demonstrated by CPA and
5826 scGen

5827 2. **Zero-Inflated Negative Binomial (ZINB) Loss**

5828 - **Approach:** Replace MSE with a **ZINB loss** that models
5829 both dropout probability and overdispersion per gene.

5830 - **Benefit:** Accounts for scRNA-seq technical noise,
5831 improving prediction in low-UMI cells (e.g., ~162 median
UMI)

5832 3. **Learned Dynamic Graph Priors**

5833 - **Approach:** Instead of a fixed PPI graph, **learn**
genegene affinity weights from data using a **Gaussian**

5832 kernel** on baseline coexpression, then refine during
5833 training.
5834

- 5835 - **Benefit:** Captures UPR pathway rewiring under CRISPRi;
5836 avoids external databases (matches constraints)

5837 4. **Contrastive Pretraining**
5838 - **Approach:**
5839 Pretrain an encoder on $(x,0)$ vs. (x,p) pairs with an **
5840 InfoNCE contrastive loss**, pulling matched
5841 baselineperturbed embeddings together and pushing apart
5842 random mismatches.
5843 - **Benefit:** Disentangles baseline state from perturbation
5844 effect; enhances generalization to unseen contexts like
5845 new guide sets or cellcycle states

5846 5. **Neural ODE or OT Trajectory Module**
5847 - **Approach:** For multiguide dynamics, model latent drift
5848 via a **neural ODE** $\frac{dz}{dt}=f(z,h_P)$ or **
5849 optimaltransport regularization** linking (x,p) to (x,p')
5850 $\frac{dz}{dt}$.
5851 - **Benefit:** Enforces smooth interpolation/extrapolation
5852 between perturbation levels (01 guide,,4 guides),
5853 capturing graded UPR responses

5854 6. **Perturbation-Guided Attention Decoder**
5855 - **Approach:** Use a **crossattention** layer where query =
5856 latent state, key/value = perturbation embedding h_P ,
5857 then an attentionaugmented decoder to weight genes most
5858 impacted by UPR factors (e.g., ATF6 targets).
5859 - **Benefit:** Focuses model capacity on biologically salient
5860 genes; improves interpretability and reduces parameter
5861 count

5862 **Implementing these enhancements**factorized, compositional
5863 embeddings; ZINB loss; datadriven graph priors; contrastive
5864 pretraining; latent ODE/OT modules; and perturbationguided
5865 attentionwill address the key limitations of SC-GPT,
5866 GeneFormer, DEEP, GEARS. This integrated architecture can
5867 robustly predict singlecell transcriptomic responses to both
5868 seen and unseen CRISPRi perturbations in the UPR pathway,
5869 while respecting the no external database constraint and
5870 modeling lowsignal, highsparsity regimes.
5871
5872
5873
5874
5875
5876
5877
5878
5879
5880
5881
5882
5883
5884
5885

5886 S.4 MODEL DESIGN MODULE

5887

5888 S.4.1 GRAPH BASED DISCUSSION

5889

5890 **Graph Based Discussion Output**

5891

5892 **Round 1**

5893

5894

5895

5896

5897

5898

5899

5900

5901

5902

5903

5904

5905

5906

5907

5908

****Data Preprocessing Expert****: Looking at the Norman et al. (2019) dataset with 111,445 cells and 33,694 genes, we need to design a preprocessing strategy that supports the dual-pathway architecture. The preprocessing must prepare data for both the ****Top Pathway**** (baseline expression processing) and ****Bottom Pathway**** (perturbation information processing). I propose a "Dual-Pathway Preprocessing Framework" (DPPF) that includes: 1) ****Expression preprocessing****: Quality control, normalization, and PCA preparation for the PCAReducer component, 2) ****Perturbation preprocessing****: Gene ID mapping, combinatorial perturbation handling, and embedding preparation for PerturbGene Embed, 3) ****Multi-scale feature preparation****: HVG selection and pathway annotation to support the Multi-Scale VAE Encoder. The framework must handle both single gene perturbations (ARID1A, BCORL1) and combinatorial perturbations (SET_KLF1, CEBPB_OSR2) while maintaining compatibility with the dual-pathway design.

Confidence Score: 0.68

5909

5910

5911

5912

5913

5914

5915

5916

5917

5918

5919

5920

5921

5922

5923

5924

5925

5926

****Model Architecture Expert****: Based on the Norman dataset characteristics and the unseen perturbation prediction task, I propose a dual-pathway architecture that processes baseline gene expression and perturbation information separately before integrating them. The architecture follows a specific design pattern: 1) ****Top Pathway****: PCAReducer Multi-Scale VAE Encoder ContextMLP, which processes baseline gene expression data to generate a "Perturbation Latent" representation, 2) ****Bottom Pathway****: PerturbGene Embed CellContexter, which processes perturbation information to generate a "Cell Context Latent" representation, 3) ****Integration****: FeaturerMixer combines both latent representations, 4) ****Core Processing****: Gene Interaction Network (GNN) PertTransformer PredictionHead Output Layer. The key innovation is the Multi-Scale VAE Encoder that captures both high-level and low-level features from gene expression data, and the PertTransformer that models long-range dependencies in perturbation effects. This architecture specifically addresses the challenge of predicting unseen perturbations by learning perturbation-aware representations.

Confidence Score: 0.72

5927

5928

5929

5930

5931

5932

5933

5934

5935

5936

5937

5938

5939

****Training and Optimization Expert****: For the dual-pathway architecture, we need a training strategy that optimizes both pathways simultaneously while ensuring effective integration. I propose a "Dual-Pathway Training Framework" (DPTF) that includes: 1) ****Multi-task learning****: Simultaneously trains perturbation effect prediction (main task) and perturbation type classification (auxiliary task), 2) ****Pathway-specific losses****: VAE reconstruction loss for the Top Pathway, contrastive loss for the Bottom Pathway, and integration loss for the FeaturerMixer, 3) ****Curriculum learning****: Progressive training from single gene perturbations to combinatorial perturbations, 4) ****Adaptive weighting****: Dynamic adjustment of loss weights based on training progress. The framework must handle the Norman dataset's imbalanced perturbation distribution and ensure both pathways learn complementary representations that enhance the final prediction accuracy.

Confidence Score: 0.75

****Pathway Expert****: The dual-pathway architecture must incorporate biological constraints to ensure meaningful predictions. I propose integrating biological priors into the ****PertTransformer**** component through a "Biologically-Constrained Attention Mechanism" (BCAM). The mechanism should: 1) ****Guide the PertTransformer attention****: Use K562-specific pathway information and chronic myelogenous leukemia context to bias attention weights, 2) ****Constrain the Gene Interaction Network****: Incorporate protein-protein interaction networks and transcription factor binding sites to guide gene-gene interactions, 3) ****Validate pathway consistency****: Ensure predictions align with known biological pathways and regulatory cascades. The BCAM should distinguish between direct effects (immediate target gene changes) and indirect effects (downstream regulatory changes) by leveraging the dual-pathway design where the Top Pathway captures cellular context and the Bottom Pathway captures perturbation-specific effects.

Confidence Score: 0.78

****Deep Learning Expert****: The dual-pathway architecture provides an ideal framework for implementing advanced deep learning techniques. I propose enhancing the architecture with: 1) ****Graph-based processing in the Gene Interaction Network****: Use graph neural networks to model gene-gene interactions, leveraging the dual-pathway design where the Top Pathway provides cellular context and the Bottom Pathway provides perturbation context, 2) ****Multi-scale attention in the PertTransformer****: Implement hierarchical attention mechanisms that operate at gene, pathway, and cellular levels, 3) ****Contrastive learning for the Bottom Pathway****: Use contrastive learning to learn meaningful perturbation embeddings that can generalize to unseen perturbations, 4) ****Adversarial training for robustness****: Implement adversarial training to improve generalization to unseen perturbations and cell contexts. The key innovation is leveraging the dual-pathway design to implement sophisticated deep learning techniques that would be difficult to achieve in a single-pathway architecture.

Confidence Score: 0.76

****Critic****: The dual-pathway architecture concept is promising, but we need more specific technical details for each component. How do we ensure the PCAReducer effectively reduces 33,694 dimensions while preserving biological information? How do we design the Multi-Scale VAE Encoder to capture both high-level and low-level features? How do we implement the PerturbGene Embed for combinatorial perturbations like SET_KLF1? How do we design the FeaturerMixer to effectively combine the two latent representations? We need detailed specifications for the Gene Interaction Network parameters, PertTransformer architecture, and PredictionHead design. The architecture must be computationally efficient while handling the Norman dataset's scale and complexity.

Confidence Score: 0.74

Graph Based Discussion Output

Round 2

****Data Preprocessing Expert****: I've refined the DPPF framework with specific components for the dual-pathway architecture. The framework now includes: 1) ****Top Pathway preprocessing****:

5994
5995 Quality control, log normalization, and PCA preparation for
5996 the PCAReducer (33,694 128 dimensions), with HVG selection
5997 to support the Multi-Scale VAE Encoder, 2) ****Bottom Pathway**
5998 **preprocessing****: Gene ID mapping and embedding preparation
5999 for PerturbGene Embed, with special handling for
6000 combinatorial perturbations (SET_KLF1, BCL2L11_BAK1) vs
6001 single gene perturbations (ARID1A, BCORL1), 3) ****Integration**
6002 **preparation****: Feature alignment and batch processing to
6003 ensure compatibility between the two pathways. The framework
6004 accounts for K562 cell line heterogeneity and varying guide
6005 RNA efficiency, implementing perturbation strength
6006 estimation that feeds into both pathways.

6007 Confidence Score: 0.79

6008 ****Model Architecture Expert****: I've detailed the dual-pathway
6009 architecture with specific implementations for each
6010 component. The ****Top Pathway**** includes: 1) ****PCAReducer****:
6011 Performs dimensionality reduction on the 33,694-dimensional
6012 gene expression data, reducing it to a manageable size while
6013 preserving biological variation, 2) ****Multi-Scale VAE**
6014 **Encoder****: A variational autoencoder that captures both high-
6015 level and low-level features from the PCA-reduced data, with
6016 separate encoders for different scales, 3) ****ContextMLP****:
6017 A multi-layer perceptron that processes the multi-scale VAE
6018 outputs to generate the "Perturbation Latent" representation.
6019 The ****Bottom Pathway**** includes: 1) ****PerturbGene Embed****:
6020 An embedding layer that maps perturbation gene IDs (like
6021 ARID1A, SET_KLF1) to dense vector representations, 2) ******
6022 **CellContexter****: Processes the perturbation embeddings to
6023 generate "Cell Context Latent" that captures how cellular
6024 context influences perturbation effects. The ****Integration**
6025 **and Processing**** includes: 1) ****FeaturerMixer****: Combines
6026 the two latent representations, 2) ****Gene Interaction**
6027 **Network****: A GNN with configurable layers that models gene-
6028 gene interactions, 3) ****PertTransformer****: A transformer
6029 with multi-head attention that captures long-range
6030 dependencies, 4) ****PredictionHead****: Final prediction layer
6031 that outputs the predicted gene expression profile.

6032 Confidence Score: 0.81

6033 ****Training and Optimization Expert****: I've designed the DPTF
6034 framework with specific implementation details for the dual-
6035 pathway architecture. The training strategy includes: 1) ******
6036 **Multi-task learning****: Simultaneously trains perturbation
6037 effect prediction (main task) and perturbation type
6038 classification (auxiliary task), 2) ****Pathway-specific**
6039 **losses****: VAE reconstruction loss for the Top Pathway (
6040 PCAReducer Multi-Scale VAE Encoder ContextMLP),
6041 contrastive loss for the Bottom Pathway (PerturbGene Embed
6042 CellContexter), and integration loss for the FeaturerMixer,
6043 3) ****Curriculum learning****: Progressive training from single
6044 gene perturbations to combinatorial perturbations, 4) ******
6045 **Adaptive weighting****: Dynamic adjustment of loss weights
6046 based on training progress. The framework handles the Norman
6047 dataset's imbalanced perturbation distribution and ensures
6048 both pathways learn complementary representations.

6049 Confidence Score: 0.83

6050 ****Pathway Expert****: I've detailed the BCAM mechanism that
6051 integrates biological priors into the dual-pathway
6052 architecture. The mechanism includes: 1) ****PertTransformer**
6053 **integration****: Uses K562-specific pathway information and
6054 chronic myelogenous leukemia context to bias attention
6055 weights in the PertTransformer component, 2) ****Gene**
6056 **Interaction Network constraints****: Incorporates protein-
6057 protein interaction networks and transcription factor
6058 binding sites to guide gene-gene interactions in the GNN
6059 layers, 3) ****Biological consistency validation****: Ensures

6048 predictions align with known biological pathways and
 6049 regulatory cascades. The implementation includes: 1) A
 6050 biological prior encoder that encodes prior knowledge as
 6051 attention biases for the PertTransformer, 2) Adaptive weight
 6052 fusion that combines data-driven attention with biological
 6053 constraints, 3) Hierarchical attention that operates at gene,
 6054 pathway, and cellular levels. The BCAM leverages the dual-
 6055 pathway design where the Top Pathway provides cellular
 6056 context and the Bottom Pathway provides perturbation-
 specific effects.

6057 Confidence Score: 0.85

6058 ****Deep Learning Expert****: I've refined the deep learning
 6059 enhancements for the dual-pathway architecture. The key
 6060 innovations include: 1) ****Graph-based processing in the Gene
 Interaction Network****: Uses graph neural networks to model
 6061 gene-gene interactions, leveraging the dual-pathway design
 6062 where the Top Pathway provides cellular context and the
 6063 Bottom Pathway provides perturbation context, 2) ****Multi-
 scale attention in the PertTransformer****: Implements
 6064 hierarchical attention mechanisms that operate at gene,
 6065 pathway, and cellular levels, 3) ****Contrastive learning for
 the Bottom Pathway****: Uses contrastive learning to learn
 6066 meaningful perturbation embeddings that can generalize to
 6067 unseen perturbations, 4) ****Adversarial training for
 robustness****: Implements adversarial training to improve
 6068 generalization to unseen perturbations and cell contexts.
 6069 The design specifically handles the Norman dataset's CRISPRa
 6070 perturbations and combinatorial nature by leveraging the
 6071 dual-pathway architecture's ability to separate cellular
 6072 context from perturbation effects.

6073 Confidence Score: 0.84

6074 ****Critic****: The dual-pathway architecture designs are now more
 6075 specific and implementable. The DPPF framework effectively
 6076 addresses preprocessing challenges for both pathways, the
 6077 dual-pathway architecture with PCAReducer Multi-Scale VAE
 6078 Encoder ContextMLP and PerturbGene Embed CellContexter is
 6079 innovative, the DPTF training strategy is sound, the BCAM
 6080 biological constraint mechanism is valuable, and the deep
 6081 learning enhancements are novel. However, we need to
 6082 consider the integration of these components and
 6083 computational efficiency. I suggest designing a unified
 6084 framework that can flexibly combine these innovative
 6085 components while maintaining computational efficiency. We
 6086 also need to ensure that the framework can handle the
 specific characteristics of the Norman dataset, including
 the combinatorial perturbations and the K562 cell line
 context.

6087 Confidence Score: 0.82

6090 Graph Based Discussion Output

6092 Round 3

6093 ****Data Preprocessing Expert****: I've finalized the DPPF framework
 6094 with specific implementations for the dual-pathway
 6095 architecture. The framework now includes: 1) ****Top Pathway
 preprocessing****: Quality control, log normalization, and PCA
 6096 preparation (33,694 128 dimensions) for the PCAReducer,
 6097 with HVG selection and pathway annotation to support the
 6098 Multi-Scale VAE Encoder, 2) ****Bottom Pathway preprocessing****:
 6099 Gene ID mapping and embedding preparation for PerturbGene
 Embed, with special handling for combinatorial perturbations
 6100 (SET_KLF1, BCL2L11_BAK1) vs single gene perturbations (
 6101 ARID1A, BCORL1), 3) ****Integration preparation****: Feature

6102
6103
6104
6105
6106
6107
6108
6109
6110
6111
6112
6113
6114
6115
6116
6117
6118
6119
6120
6121
6122
6123
6124
6125
6126
6127
6128
6129
6130
6131
6132
6133
6134
6135
6136
6137
6138
6139
6140
6141
6142
6143
6144
6145
6146
6147
6148
6149
6150
6151
6152
6153
6154
6155

alignment and batch processing to ensure compatibility between the two pathways. The framework accounts for K562 cell line heterogeneity and varying guide RNA efficiency, implementing perturbation strength estimation that feeds into both pathways.

Confidence Score: 0.84

****Model Architecture Expert****: I've completed the detailed architecture design that exactly matches the proposed framework. The ****Top Pathway**** processes baseline gene expression: 1) ****PCAReducer****: Reduces 33,694 genes to 128 principal components, 2) ****Multi-Scale VAE Encoder****: Uses separate encoders for high-level (pathway-level) and low-level (gene-level) features, 3) ****ContextMLP****: Takes high and low features as input and outputs the "Perturbation Latent" representation. The ****Bottom Pathway**** processes perturbation information: 1) ****PerturbGene Embed****: Maps perturbation gene IDs to 64-dimensional embeddings, 2) ****CellContexter****: Generates "Cell Context Latent" representation. The ****Integration**** uses ****FeaturerMixer**** to combine both latent representations, followed by ****Gene Interaction Network**** (GNN with 3 layers), ****PertTransformer**** (4 attention heads, 2 layers), ****PredictionHead**** (final prediction), and ****Output Layer**** (predicted expression profile).

Confidence Score: 0.86

****Training and Optimization Expert****: I've refined the training strategy for the dual-pathway architecture. The training includes: 1) ****Multi-task learning****: Simultaneously trains perturbation effect prediction and perturbation type classification, 2) ****Contrastive learning****: Uses hierarchical sampling to create positive pairs (similar perturbations) and negative pairs (different perturbations), 3) ****Curriculum learning****: Starts with single gene perturbations and gradually introduces combinatorial perturbations, 4) ****Loss function****: Combines reconstruction loss (MSE), contrastive loss, and biological consistency loss. The training strategy specifically handles the Norman dataset's imbalanced perturbation distribution and ensures proper representation of both single and combinatorial perturbations.

Confidence Score: 0.85

****Pathway Expert****: I've integrated biological constraints into the architecture design. The ****Biologically-Constrained Attention Mechanism**** (BCAM) is implemented within the ****PertTransformer**** component: 1) ****Biological prior encoding****: Incorporates K562-specific pathway information and chronic myelogenous leukemia context, 2) ****Attention bias****: Uses transcription factor binding sites and protein-protein interaction networks to guide attention weights, 3) ****Hierarchical attention****: Applies attention at gene, pathway, and cellular levels. The mechanism ensures that model predictions align with known biological principles and can distinguish between direct and indirect perturbation effects.

Confidence Score: 0.87

****Deep Learning Expert****: I've finalized the ****Cell-Perturbation Interaction Network**** (CPIN) design that integrates with the dual-pathway architecture. The ****Dual-Graph Interaction Mechanism**** (DGIM) is implemented as follows: 1) ****Cell state graph****: Constructed from gene co-expression patterns in the baseline expression data, 2) ****Perturbation effect graph****: Built from perturbation types and their biological properties, 3) ****Graph-graph interaction****: Implemented

6156 within the **Gene Interaction Network** component using
 6157 graph neural networks, 4) **Multi-scale fusion**: Combines
 6158 information at different biological scales within the **FeaturerMixer**. The design specifically handles the Norman
 6159 dataset's CRISPRa perturbations and combinatorial nature.
 6160

6161 Confidence Score: 0.88

6162 **Critic**: The architecture design now perfectly matches the
 6163 proposed framework and addresses all the specific
 6164 requirements of the Norman dataset and unseen perturbation
 6165 prediction task. The dual-pathway design effectively
 6166 separates baseline expression processing from perturbation
 6167 information processing, the integration components (
 6168 FeaturerMixer, Gene Interaction Network, PertTransformer)
 6169 provide sophisticated modeling capabilities, and the
 6170 biological constraints ensure biological interpretability.
 6171 The architecture is computationally efficient and can handle
 6172 both single gene and combinatorial perturbations. This
 6173 design represents a significant advancement in single-cell
 6174 perturbation prediction methodology.

6173 Confidence Score: 0.89

6174

6175 Graph Based Discussion Output

6176

6177 Round 4

6178

6179 **Data Preprocessing Expert**: The CSAPF framework is now
 6180 complete and optimized for the Norman dataset. The final
 6181 implementation includes: 1) **Adaptive normalization**: Each
 6182 cell's expression is normalized based on its quality
 6183 metrics (UMI count, gene count, mitochondrial ratio), 2) **Perturbation-specific processing**: Single gene
 6184 perturbations and combinatorial perturbations are processed
 6185 differently, 3) **Multi-scale feature extraction**: Features
 6186 are extracted at gene, pathway, and cellular levels, 4) **Quality control**: Cells with poor coverage or high
 6187 mitochondrial content are filtered out. The framework
 6188 successfully handles the dataset's 111,445 cells and 33,694
 6189 genes while maintaining biological interpretability.

6189 Confidence Score: 0.87

6190

6191 **Model Architecture Expert**: The dual-pathway architecture is
 6192 finalized and ready for implementation. The complete
 6193 architecture follows the exact design pattern: **Top Pathway**
 6194 **PCAReducer** (33,694 128 dimensions) **Multi-Scale VAE**
 6195 **Encoder** (high-level and low-level features) **ContextMLP** (
 6196 generates Perturbation Latent). **Bottom Pathway**:
 6197 **PerturbGene Embed** (perturbation IDs 64D embeddings)
 6198 **CellContexter** (generates Cell Context Latent). **Integration**
 6199 **FeaturerMixer** (combines both latents) **Gene Interaction**
 6200 **Network** (3-layer GNN) **PertTransformer** (4 heads, 2 layers)
 6201 **PredictionHead** Output Layer (predicted expression profile).
 6202 This architecture specifically addresses unseen
 6203 perturbation prediction by learning perturbation-aware
 6204 representations that can generalize to new perturbation
 6205 types.

6202 Confidence Score: 0.89

6203

6204 **Training and Optimization Expert**: The training strategy is
 6205 optimized for the dual-pathway architecture and Norman
 6206 dataset. The final training approach includes: 1) **Multi-**
 6207 **task learning**: Predicts both perturbation effects and
 6208 perturbation types, 2) **Contrastive learning**: Uses
 6209 biological similarity to create training pairs, 3) **Curriculum learning**: Progressive training from simple to
 6210 complex perturbations, 4) **Loss combination**: MSE +

6210
6211
6212
6213
6214
6215
6216
6217
6218
6219
6220
6221
6222
6223
6224
6225
6226
6227
6228
6229
6230
6231
6232
6233
6234
6235
6236
6237
6238
6239
6240
6241
6242
6243
6244
6245
6246
6247
6248
6249
6250
6251
6252
6253
6254
6255
6256
6257
6258
6259
6260
6261
6262
6263

contrastive loss + biological consistency loss. The strategy effectively handles the dataset's imbalanced perturbation distribution and ensures robust generalization to unseen perturbations.

Confidence Score: 0.88

****Pathway Expert****: The biological constraints are fully integrated into the architecture. The BCAM mechanism within the PertTransformer ensures: 1) ****Biological consistency****: Predictions align with known biological principles, 2) ****K562-specific context****: Incorporates chronic myelogenous leukemia pathway information, 3) ****Hierarchical attention****: Operates at multiple biological scales, 4) ****Interpretability****: Provides biological explanations for predictions. The mechanism successfully distinguishes between direct and indirect perturbation effects while maintaining computational efficiency.

Confidence Score: 0.90

****Deep Learning Expert****: The CPIN design is complete and integrated with the dual-pathway architecture. The DGIM mechanism provides: 1) ****Cell-perturbation interaction modeling****: Explicitly models how cellular context influences perturbation effects, 2) ****Multi-scale representation****: Captures information at gene, pathway, and cellular levels, 3) ****Graph-based processing****: Uses GNNs to model complex biological relationships, 4) ****Generalization capability****: Can handle both seen and unseen perturbation types. The design successfully addresses the challenges of single-cell perturbation prediction while maintaining biological interpretability.

Confidence Score: 0.91

****Critic****: The final architecture design is comprehensive, innovative, and perfectly suited for the single-cell unseen perturbation prediction task using the Norman dataset. The dual-pathway design effectively separates and processes baseline expression and perturbation information, the integration components provide sophisticated modeling capabilities, and the biological constraints ensure biological interpretability. The architecture can handle both single gene and combinatorial perturbations, generalizes to unseen perturbations, and maintains computational efficiency. This represents a significant advancement in the field of single-cell perturbation prediction and provides a solid foundation for future research.

Confidence Score: 0.89

S.4.2 METHOD DESIGN RESEARCH PLAN

Graph Based Discussion Output

```

6264 # Novel Method Implementation Framework
6265
6266 ## Architecture Overview
6267 The proposed architecture follows a dual-pathway design that
6268 processes baseline gene expression and perturbation
6269 information separately before integrating them for final
6270 prediction.
6271
6272 ## 1. Top Pathway: Baseline Expression Processing
6273
6274 ### 1.1 PCAReducer
6275 ```python
6276 class PCAReducer(nn.Module):
6277     def __init__(self, input_dim=33694, n_components=128):
6278         super().__init__()
6279         self.pca = PCA(n_components=n_components)
6280         self.input_dim = input_dim
6281         self.n_components = n_components
6282
6283     def forward(self, expression_data):
6284         # expression_data: [batch_size, 33694]
6285         # Returns: [batch_size, 128]
6286         ... return self.pca.fit_transform(expression_data)
6287
6288 ### 1.2 Multi-Scale VAE Encoder
6289 ```python
6290 class MultiScaleVAEEncoder(nn.Module):
6291     def __init__(self, input_dim=128, latent_dim=64):
6292         super().__init__()
6293         # High-level encoder (pathway-level features)
6294         self.high_encoder = nn.Sequential(
6295             nn.Linear(input_dim, 256),
6296             nn.ReLU(),
6297             nn.Linear(256, 128),
6298             nn.ReLU(),
6299             nn.Linear(128, latent_dim * 2) # mean and logvar
6300         )
6301
6302         # Low-level encoder (gene-level features)
6303         self.low_encoder = nn.Sequential(
6304             nn.Linear(input_dim, 256),
6305             nn.ReLU(),
6306             nn.Linear(256, 128),
6307             nn.ReLU(),
6308             nn.Linear(128, latent_dim * 2) # mean and logvar
6309         )
6310
6311     def forward(self, x):
6312         # x: [batch_size, 128] (PCA-reduced expression)
6313         high_mu, high_logvar = self.high_encoder(x).chunk(2, dim
6314             ==-1)
6315         low_mu, low_logvar = self.low_encoder(x).chunk(2, dim=-1)
6316
6317         high_z = self.reparameterize(high_mu, high_logvar)
6318         low_z = self.reparameterize(low_mu, low_logvar)
6319
6320         return high_z, low_z, high_mu, high_logvar, low_mu,
6321             low_logvar
6322
6323     def reparameterize(self, mu, logvar):
6324         std = torch.exp(0.5 * logvar)
6325         eps = torch.randn_like(std)
6326         return mu + eps * std
6327
6328 ### 1.3 ContextMLP
6329 ```python

```

```

6318
6319 class ContextMLP(nn.Module):
6320     def __init__(self, high_dim=64, low_dim=64, output_dim=128):
6321         super().__init__()
6322         self.mlp = nn.Sequential(
6323             nn.Linear(high_dim + low_dim, 256),
6324             nn.ReLU(),
6325             nn.Linear(256, 128),
6326             nn.ReLU(),
6327             nn.Linear(128, output_dim)
6328         )
6329
6330     def forward(self, high, low):
6331         # high: [batch_size, 64] (high-level features)
6332         # low: [batch_size, 64] (low-level features)
6333         # Returns: [batch_size, 128] (Perturbation Latent)
6334         combined = torch.cat([high, low], dim=-1)
6335         return self.mlp(combined)
6336
6337 ...
6338
6339 ## 2. Bottom Pathway: Perturbation Information Processing
6340
6341 ### 2.1 PerturbGene Embed
6342 ```python
6343 class PerturbGeneEmbed(nn.Module):
6344     def __init__(self, vocab_size=1000, embed_dim=64):
6345         super().__init__()
6346         self.embedding = nn.Embedding(vocab_size, embed_dim)
6347         self.projection = nn.Sequential(
6348             nn.Linear(embed_dim, 128),
6349             nn.ReLU(),
6350             nn.Linear(128, 64)
6351         )
6352
6353     def forward(self, perturbation_ids):
6354         # perturbation_ids: [batch_size] (gene IDs)
6355         # Returns: [batch_size, 64]
6356         embedded = self.embedding(perturbation_ids)
6357         return self.projection(embedded)
6358
6359 ...
6360
6361 ### 2.2 CellContexter
6362 ```python
6363 class CellContexter(nn.Module):
6364     def __init__(self, input_dim=64, output_dim=128):
6365         super().__init__()
6366         self.context_encoder = nn.Sequential(
6367             nn.Linear(input_dim, 128),
6368             nn.ReLU(),
6369             nn.Linear(128, 128),
6370             nn.ReLU(),
6371             nn.Linear(128, output_dim)
6372         )
6373
6374     def forward(self, perturbation_embed):
6375         # perturbation_embed: [batch_size, 64]
6376         # Returns: [batch_size, 128] (Cell Context Latent)
6377         return self.context_encoder(perturbation_embed)
6378
6379 ...
6380
6381 ## 3. Integration and Core Processing
6382
6383 ### 3.1 FeaturerMixer
6384 ```python
6385 class FeaturerMixer(nn.Module):
6386     def __init__(self, pert_latent_dim=128, cell_latent_dim=128,
6387                 output_dim=256):
6388         super().__init__()
6389         self.mixer = nn.Sequential(
6390             nn.Linear(pert_latent_dim + cell_latent_dim, 512),
6391             nn.ReLU(),
6392             nn.Linear(512, 256),
6393             nn.ReLU(),
6394         )
6395
6396     def forward(self, pert_latent, cell_latent):
6397         combined = torch.cat([pert_latent, cell_latent], dim=-1)
6398         return self.mixer(combined)
6399
6400 ...

```

```

6372         nn.Linear(256, output_dim)
6373     )
6374
6375     def forward(self, pert_latent, cell_latent):
6376         # pert_latent: [batch_size, 128] (Perturbation Latent)
6377         # cell_latent: [batch_size, 128] (Cell Context Latent)
6378         # Returns: [batch_size, 256]
6379         combined = torch.cat([pert_latent, cell_latent], dim=-1)
6380         return self.mixer(combined)
6381
6382     """
6383     """python
6384     class GeneInteractionNetwork(nn.Module):
6385         def __init__(self, HVG, cell, D_pert, D_gene_feature,
6386                     num_gnn_layers=3):
6387             super().__init__()
6388             self.num_gnn_layers = num_gnn_layers
6389             self.gnn_layers = nn.ModuleList([
6390                 GraphConv(D_gene_feature, D_gene_feature)
6391                 for _ in range(num_gnn_layers)
6392             ])
6393             self.attention = nn.MultiheadAttention(D_gene_feature,
6394                                                    num_heads=8)
6395
6396         def forward(self, gene, cell, pert):
6397             # gene: [batch_size, num_genes, D_gene_feature]
6398             # cell: [batch_size, D_gene_feature]
6399             # pert: [batch_size, D_gene_feature]
6400             x = gene
6401
6402             for layer in self.gnn_layers:
6403                 x = layer(x, cell, pert)
6404                 x = F.relu(x)
6405
6406             # Apply attention mechanism
6407             attended, _ = self.attention(x, x, x)
6408
6409             return attended
6410
6411     """
6412     """
6413     """
6414     """python
6415     class PertTransformer(nn.Module):
6416         def __init__(self, D_gene_feature=256, num_heads=4,
6417                     num_layers=2):
6418             super().__init__()
6419             self.transformer = nn.TransformerEncoder(
6420                 nn.TransformerEncoderLayer(
6421                     d_model=D_gene_feature,
6422                     nhead=num_heads,
6423                     dim_feedforward=512,
6424                     dropout=0.1
6425                 ),
6426                 num_layers=num_layers
6427             )
6428
6429         def forward(self, x):
6430             # x: [batch_size, seq_len, D_gene_feature]
6431             return self.transformer(x)
6432
6433     """
6434     """
6435     """
6436     """python
6437     class PredictionHead(nn.Module):
6438         def __init__(self, feature_dim=256, num_heads=4, num_layers=2,
6439                     output_dim=33694):
6440             super().__init__()
6441             self.attention = nn.MultiheadAttention(feature_dim,
6442                                                    num_heads)
6443             self.layers = nn.ModuleList([
6444                 nn.Linear(feature_dim, feature_dim)
6445                 for _ in range(num_layers)

```

```

6426         ])
6427         self.output_layer = nn.Linear(feature_dim, output_dim)
6428
6429     def forward(self, final_genes):
6430         # final_genes: [batch_size, seq_len, feature_dim]
6431         attended, _ = self.attention(final_genes, final_genes,
6432                                     final_genes)
6433
6434         for layer in self.layers:
6435             attended = F.relu(layer(attended))
6436
6437         # Global average pooling
6438         pooled = attended.mean(dim=1) # [batch_size, feature_dim]
6439
6440         # Final prediction
6441         output = self.output_layer(pooled) # [batch_size, 33694]
6442
6443     ...
6444
6445     ## 4. Complete Model Architecture
6446
6447     ### 4.1 CellForge Model
6448     ```python
6449     class CellForgeModel(nn.Module):
6450     def __init__(self, config):
6451         super().__init__()
6452         self.config = config
6453
6454         # Top pathway components
6455         self.pca_reducer = PCAReducer(input_dim=33694,
6456                                     n_components=128)
6457         self.vae_encoder = MultiScaleVAEEncoder(input_dim=128,
6458                                                 latent_dim=64)
6459         self.context_mlp = ContextMLP(high_dim=64, low_dim=64,
6460                                     output_dim=128)
6461
6462         # Bottom pathway components
6463         self.pert_embed = PerturbGeneEmbed(vocab_size=1000,
6464                                           embed_dim=64)
6465         self.cell_contexter = CellContexter(input_dim=64,
6466                                             output_dim=128)
6467
6468         # Integration and processing components
6469         self.feature_mixer = FeaturerMixer(pert_latent_dim=128,
6470                                           cell_latent_dim=128, output_dim=256)
6471         self.gene_interaction = GeneInteractionNetwork(
6472             HVG=config.HVG, cell=config.cell, D_pert=config.D_pert,
6473             D_gene_feature=256, num_gnn_layers=3
6474         )
6475         self.pert_transformer = PertTransformer(D_gene_feature=256,
6476                                               num_heads=4, num_layers=2)
6477         self.prediction_head = PredictionHead(feature_dim=256,
6478                                               num_heads=4, num_layers=2, output_dim=33694)
6479
6480     def forward(self, expression_data, perturbation_ids):
6481         # Top pathway: Process baseline expression
6482         pca_reduced = self.pca_reducer(expression_data)
6483         high_z, low_z, high_mu, high_logvar, low_mu, low_logvar =
6484             self.vae_encoder(pca_reduced)
6485         pert_latent = self.context_mlp(high_z, low_z)
6486
6487         # Bottom pathway: Process perturbation information
6488         pert_embed = self.pert_embed(perturbation_ids)
6489         cell_latent = self.cell_contexter(pert_embed)
6490
6491         # Integration and processing
6492         mixed_features = self.feature_mixer(pert_latent,
6493                                           cell_latent)
6494
6495         # Reshape for GNN processing
6496         gene_features = mixed_features.unsqueeze(1).repeat(1, 1000,
6497                                                         1) # [batch_size, 1000, 256]
6498

```

6480
6481
6482
6483
6484
6485
6486
6487
6488
6489
6490
6491
6492
6493
6494
6495
6496
6497
6498
6499
6500
6501
6502
6503
6504
6505
6506
6507
6508
6509
6510
6511
6512
6513
6514
6515
6516
6517
6518
6519
6520
6521
6522
6523
6524
6525
6526
6527
6528
6529
6530
6531
6532
6533

```
        # Core processing
        gnn_output = self.gene_interaction(gene_features,
            cell_latent, pert_latent)
        transformer_output = self.pert_transformer(gnn_output)
        predicted_expression = self.prediction_head(
            transformer_output)

        return predicted_expression, high_mu, high_logvar, low_mu,
        low_logvar
    """
```

This implementation framework provides a complete, implementable architecture that exactly matches the proposed design pattern and addresses the specific requirements of single-cell unseen perturbation prediction using the Norman dataset.

6534 S.4.3 DETAILED RESEARCH PLAN

6536 **Graph Based Discussion Output**

```

6537
6538
6539 # Detailed Research Plan for Dual-Pathway Architecture
6540 ## 1. Research Objectives
6541 ### 1.1 Primary Objective
6542 Develop a novel dual-pathway neural network architecture for
6543 predicting gene expression profiles of individual K562 cells
6544 following CRISPR interference (CRISPRi) perturbations,
6545 specifically designed to handle unseen perturbations and
6546 unseen cell contexts.
6547 ### 1.2 Secondary Objectives
6548 - Design a preprocessing framework that supports both pathways
6549 of the architecture
6549 - Implement biological constraints to ensure meaningful
6550 predictions
6550 - Develop training strategies that optimize both pathways
6551 simultaneously
6551 - Create evaluation metrics that assess both predictive accuracy
6552 and biological relevance
6552
6553 ## 2. Dataset and Task Specification
6554 ### 2.1 Dataset Details
6555 - **Source**: Norman et al. (2019, Science) dataset
6556 - **Modality**: RNA (scRNA-seq gene expression data)
6557 - **Perturbation Type**: CRISPRa (CRISPR activation)
6558 - **Cell Line**: K562 (chronic myelogenous leukemia lymphoblasts
6559 )
6560 - **Scale**: 111,445 cells 33,694 genes
6561 - **Perturbation Types**: Single gene (e.g., ARID1A, BCORL1) and
6562 combinatorial (e.g., SET_KLF1, CEBPB_OSR2)
6563 ### 2.2 Task Definition
6564 - **Input**: Baseline gene expression profile of an unperturbed
6565 K562 cell and the identity of the target gene(s) for
6566 perturbation
6566 - **Output**: Predicted gene expression profile after
6567 perturbation
6567 - **Evaluation Scenarios**:
6568 - Unseen Perturbations: Predict effects of gene perturbations
6569 not present during training
6569 - Unseen Cell Contexts: Predict responses in cells with gene
6570 expression profiles not observed during training
6570
6571 ### 2.3 Evaluation Metrics
6572 - **Predictive Performance**: MSE, PCC, R
6573 - **Biological Relevance**: MSE_DE, PCC_DE, R_DE (for
6574 differentially expressed genes)
6574
6575 ## 3. Dual-Pathway Architecture Design
6576 ### 3.1 Top Pathway: Baseline Expression Processing
6577 **Purpose**: Process baseline gene expression data to generate
6578 perturbation-aware cellular context
6578
6579 **Components**:
6580 1. **PCAReducer**
6581 - Input: 33,694-dimensional gene expression data
6582 - Output: 128-dimensional PCA-reduced representation
6583 - Function: Dimensionality reduction while preserving
6584 biological variation
6584
6585 2. **Multi-Scale VAE Encoder**
6586 - Input: 128-dimensional PCA-reduced data
6587 - Output: High-level (64D) and low-level (64D) latent
6588 representations

```

```

6588     - Function: Capture both pathway-level and gene-level features
6589
6590 3. **ContextMLP**
6591     - Input: High-level and low-level features (64D each)
6592     - Output: 128-dimensional "Perturbation Latent" representation
6593     - Function: Generate perturbation-aware cellular context
6594
6595 ### 3.2 Bottom Pathway: Perturbation Information Processing
6596 **Purpose**: Process perturbation information to generate cell
6597     context-aware perturbation effects
6598
6599 **Components**:
6600 1. **PerturbGene Embed**
6601     - Input: Perturbation gene IDs
6602     - Output: 64-dimensional perturbation embeddings
6603     - Function: Map perturbation identities to dense
6604       representations
6605
6606 2. **CellContexter**
6607     - Input: 64-dimensional perturbation embeddings
6608     - Output: 128-dimensional "Cell Context Latent" representation
6609     - Function: Generate cell context-aware perturbation effects
6610
6611 ### 3.3 Integration and Core Processing
6612 **Purpose**: Combine and process information from both pathways
6613
6614 **Components**:
6615 1. **FeaturerMixer**
6616     - Input: Perturbation Latent (128D) + Cell Context Latent (128
6617       D)
6618     - Output: 256-dimensional fused representation
6619     - Function: Combine information from both pathways
6620
6621 2. **Gene Interaction Network**
6622     - Input: 256-dimensional fused features
6623     - Output: 256-dimensional gene interaction features
6624     - Function: Model gene-gene interactions using GNN (3 layers)
6625
6626 3. **PertTransformer**
6627     - Input: 256-dimensional gene interaction features
6628     - Output: 256-dimensional transformed features
6629     - Function: Capture long-range dependencies (4 heads, 2 layers
6630       )
6631
6632 4. **PredictionHead**
6633     - Input: 256-dimensional transformed features
6634     - Output: 33,694-dimensional predicted expression profile
6635     - Function: Final prediction layer
6636
6637 ## 4. Implementation Plan
6638
6639 ### 4.1 Phase 1: Data Preprocessing (Weeks 1-2)
6640 **Objective**: Implement the Dual-Pathway Preprocessing
6641     Framework (DPPF)
6642
6643 **Tasks**:
6644 1. **Top Pathway Preprocessing**
6645     - Quality control: Filter cells with poor coverage or high
6646       mitochondrial content
6647     - Normalization: Log normalization and library size correction
6648     - PCA preparation: Reduce 33,694 genes to 128 principal
6649       components
6650     - HVG selection: Identify highly variable genes for pathway
6651       annotation
6652
6653 2. **Bottom Pathway Preprocessing**
6654     - Gene ID mapping: Create mapping from perturbation names to
6655       gene IDs
6656     - Combinatorial perturbation handling: Process multi-gene
6657       perturbations
6658     - Embedding preparation: Prepare data for PerturbGene Embed
6659
6660 3. **Integration Preparation**

```

```

6642     - Feature alignment: Ensure compatibility between pathways
6643     - Batch processing: Implement efficient data loading
6644     - Perturbation strength estimation: Infer perturbation
6645       efficiency
6646
6647     **Deliverables**:
6648     - Preprocessed dataset with both pathway inputs
6649     - Data loading pipeline
6650     - Quality control metrics
6651
6652     ### 4.2 Phase 2: Architecture Implementation (Weeks 3-5)
6653     **Objective**: Implement the complete dual-pathway architecture
6654
6655     **Tasks**:
6656     1. **Top Pathway Components**
6657         - PCAReducer: Implement PCA dimensionality reduction
6658         - Multi-Scale VAE Encoder: Implement high-level and low-level
6659           encoders
6660         - ContextMLP: Implement perturbation-aware context generation
6661
6662     2. **Bottom Pathway Components**
6663         - PerturbGene Embed: Implement perturbation embedding layer
6664         - CellContexter: Implement cell context-aware processing
6665
6666     3. **Integration Components**
6667         - FeaturerMixer: Implement feature fusion
6668         - Gene Interaction Network: Implement GNN with 3 layers
6669         - PertTransformer: Implement transformer with 4 heads, 2
6670           layers
6671         - PredictionHead: Implement final prediction layer
6672
6673     **Deliverables**:
6674     - Complete model architecture
6675     - Forward pass implementation
6676     - Model parameter specifications
6677
6678     ### 4.3 Phase 3: Training Strategy Implementation (Weeks 6-7)
6679     **Objective**: Implement the Dual-Pathway Training Framework (
6680       DPTF)
6681
6682     **Tasks**:
6683     1. **Multi-task Learning**
6684         - Main task: Perturbation effect prediction
6685         - Auxiliary task: Perturbation type classification
6686
6687     2. **Pathway-specific Losses**
6688         - VAE reconstruction loss for Top Pathway
6689         - Contrastive loss for Bottom Pathway
6690         - Integration loss for FeaturerMixer
6691
6692     3. **Curriculum Learning**
6693         - Progressive training from single to combinatorial
6694           perturbations
6695         - Adaptive weighting based on training progress
6696
6697     **Deliverables**:
6698     - Training pipeline
6699     - Loss function implementations
6700     - Curriculum learning scheduler
6701
6702     ### 4.4 Phase 4: Biological Constraints Integration (Weeks 8-9)
6703     **Objective**: Implement the Biologically-Constrained Attention
6704       Mechanism (BCAM)
6705
6706     **Tasks**:
6707     1. **Biological Prior Integration**
6708         - K562-specific pathway information
6709         - Chronic myelogenous leukemia context
6710         - Protein-protein interaction networks
6711
6712     2. **Attention Mechanism Enhancement**
6713         - Biological bias injection into PertTransformer
6714         - Hierarchical attention at multiple biological scales

```

```

6696     - Biological consistency validation
6697
6698     **Deliverables**:
6699     - BCAM implementation
6700     - Biological prior database
6701     - Attention visualization tools
6702
6703     ### 4.5 Phase 5: Model Training and Optimization (Weeks 10-12)
6704     **Objective**: Train and optimize the complete model
6705
6706     **Tasks**:
6707     1. **Hyperparameter Tuning**
6708         - Learning rate optimization
6709         - Batch size selection
6710         - Architecture parameter tuning
6711
6712     2. **Training Monitoring**
6713         - Loss tracking for both pathways
6714         - Biological consistency monitoring
6715         - Generalization assessment
6716
6717     3. **Model Selection**
6718         - Cross-validation
6719         - Early stopping
6720         - Best model checkpointing
6721
6722     **Deliverables**:
6723     - Trained model
6724     - Training logs and metrics
6725     - Hyperparameter configurations
6726
6727     ### 4.6 Phase 6: Evaluation and Analysis (Weeks 13-14)
6728     **Objective**: Comprehensive evaluation of the model
6729
6730     **Tasks**:
6731     1. **Predictive Performance Evaluation**
6732         - MSE, PCC, R on test set
6733         - Unseen perturbation evaluation
6734         - Unseen cell context evaluation
6735
6736     2. **Biological Relevance Assessment**
6737         - MSE_DE, PCC_DE, R_DE for differentially expressed genes
6738         - Pathway enrichment analysis
6739         - Biological consistency validation
6740
6741     3. **Comparative Analysis**
6742         - Baseline method comparison
6743         - Ablation studies
6744         - Component contribution analysis
6745
6746     **Deliverables**:
6747     - Evaluation results
6748     - Comparative analysis report
6749     - Biological interpretation
6750
6751     ## 5. Technical Specifications
6752
6753     ### 5.1 Computational Requirements
6754     - **Hardware**: GPU with at least 16GB memory
6755     - **Software**: PyTorch, scikit-learn, scanpy, pandas, numpy
6756     - **Storage**: At least 100GB for dataset and model checkpoints
6757
6758     ### 5.2 Model Parameters
6759     - **Input Dimensions**: 33,694 (genes) 111,445 (cells)
6760     - **PCA Components**: 128
6761     - **VAE Latent Dimensions**: 64 (high-level) + 64 (low-level)
6762     - **Embedding Dimensions**: 64 (perturbation) + 128 (context)
6763     - **GNN Layers**: 3
6764     - **Transformer Heads**: 4
6765     - **Transformer Layers**: 2
6766     - **Output Dimensions**: 33,694 (predicted expression)
6767
6768     ### 5.3 Training Configuration

```

```

6750
6751 - **Batch Size**: 32
6752 - **Learning Rate**: 1e-3 (with OneCycle scheduler)
6753 - **Epochs**: 100 (with early stopping)
6754 - **Optimizer**: AdamW
6755 - **Loss Weights**: Adaptive based on training progress
6756
6757 ## 6. Risk Assessment and Mitigation
6758
6759 ### 6.1 Technical Risks
6760 1. **Computational Complexity**
6761 - Risk: Model too large for available hardware
6762 - Mitigation: Implement gradient checkpointing and model
6763 parallelism
6764
6765 2. **Training Instability**
6766 - Risk: Training divergence or poor convergence
6767 - Mitigation: Implement gradient clipping and learning rate
6768 scheduling
6769
6770 3. **Overfitting**
6771 - Risk: Model memorizes training data
6772 - Mitigation: Implement dropout, weight decay, and early
6773 stopping
6774
6775 ### 6.2 Biological Risks
6776 1. **Biological Inconsistency**
6777 - Risk: Predictions violate biological principles
6778 - Mitigation: Implement biological constraints and validation
6779
6780 2. **Generalization Failure**
6781 - Risk: Poor performance on unseen perturbations
6782 - Mitigation: Implement contrastive learning and adversarial
6783 training
6784
6785 ### 6.3 Data Risks
6786 1. **Data Quality Issues**
6787 - Risk: Poor quality data affects model performance
6788 - Mitigation: Implement robust quality control and data
6789 augmentation
6790
6791 2. **Class Imbalance**
6792 - Risk: Imbalanced perturbation types
6793 - Mitigation: Implement weighted sampling and focal loss
6794
6795 ## 7. Success Metrics
6796
6797 ### 7.1 Quantitative Metrics
6798 - **Predictive Accuracy**: MSE < 0.1, PCC > 0.8, R > 0.7
6799 - **Biological Relevance**: MSE_DE < 0.15, PCC_DE > 0.75, R_DE >
6800 0.65
6801 - **Generalization**: Performance on unseen perturbations within
6802 10% of seen perturbations
6803
6804 ### 7.2 Qualitative Metrics
6805 - **Biological Interpretability**: Predictions align with known
6806 biological pathways
6807 - **Computational Efficiency**: Training time < 24 hours on
6808 single GPU
6809 - **Reproducibility**: Results reproducible across multiple runs
6810
6811 ## 8. Timeline and Milestones
6812
6813 ### 8.1 Week 1-2: Data Preprocessing
6814 - **Milestone**: Complete DPPF implementation
6815 - **Deliverable**: Preprocessed dataset ready for training
6816
6817 ### 8.2 Week 3-5: Architecture Implementation
6818 - **Milestone**: Complete dual-pathway architecture
6819 - **Deliverable**: Functional model implementation
6820
6821 ### 8.3 Week 6-7: Training Strategy
6822 - **Milestone**: Complete DPTF implementation
6823 - **Deliverable**: Training pipeline ready

```

6804
6805
6806
6807
6808
6809
6810
6811
6812
6813
6814
6815
6816
6817
6818
6819
6820
6821
6822
6823
6824
6825
6826
6827
6828
6829
6830
6831
6832
6833
6834
6835
6836
6837
6838
6839
6840
6841
6842
6843
6844
6845
6846
6847
6848
6849
6850
6851
6852
6853
6854
6855
6856
6857

```
### 8.4 Week 8-9: Biological Constraints
- **Milestone**: Complete BCAM implementation
- **Deliverable**: Biologically-constrained model

### 8.5 Week 10-12: Model Training
- **Milestone**: Trained and optimized model
- **Deliverable**: Best model checkpoint

### 8.6 Week 13-14: Evaluation
- **Milestone**: Comprehensive evaluation complete
- **Deliverable**: Final results and analysis

## 9. Expected Outcomes

### 9.1 Scientific Contributions
- Novel dual-pathway architecture for single-cell perturbation prediction
- Integration of biological constraints into deep learning models
- Comprehensive evaluation framework for perturbation prediction

### 9.2 Technical Contributions
- Efficient preprocessing pipeline for large-scale single-cell data
- Advanced training strategies for multi-pathway architectures
- Robust evaluation metrics for biological relevance

### 9.3 Practical Applications
- Tool for predicting effects of novel perturbations
- Framework for understanding cellular responses to perturbations
- Foundation for drug discovery and therapeutic development

This detailed research plan provides a comprehensive roadmap for implementing the dual-pathway architecture and achieving the research objectives.
```

S.4.4 EXAMPLE MODEL FRAMEWORK AND PSEUDOCODE

Graph Based Discussion Output

```

6858
6859
6860
6861
6862
6863 # Architecture Flowchart and Pseudocode
6864
6865 ## 1. Mermaid Flowchart
6866
6867 ```mermaid
6868 graph TD
6869   A[Raw Gene Expression Data<br/>33,694 genes 111,445 cells]
6870   --> B[Quality Control<br/>Filter low-quality cells]
6871   B --> C[Log Normalization<br/>Library size correction]
6872   C --> D[PCA Reduction<br/>33,694 128 dimensions]
6873
6874   E[Perturbation Information<br/>Gene IDs] --> F[Gene ID
6875   Mapping<br/>Handle combinatorial perturbations]
6876   F --> G[Perturbation Embedding<br/>Gene IDs 64D vectors]
6877
6878   D --> H[Multi-Scale VAE Encoder<br/>High-level + Low-level
6879   features]
6880   H --> I[ContextMPLP<br/>Generate Perturbation Latent 128D]
6881   G --> J[CellContexter<br/>Generate Cell Context Latent 128D]
6882   I --> K[FeaturerMixer<br/>Combine both latents 256D]
6883   J --> K
6884
6885   K --> L[Gene Interaction Network<br/>GNN with 3 layers]
6886   L --> M[PertTransformer<br/>4 heads, 2 layers]
6887   M --> N[PredictionHead<br/>Final prediction layer]
6888   N --> O[Predicted Expression Profile<br/>33,694 dimensions]
6889
6890   P[Biological Priors<br/>K562 pathways, PPI networks] --> Q[
6891   BCAM<br/>Biological constraints]
6892   Q --> M
6893
6894   R[Training Strategy<br/>Multi-task learning] --> S[Loss
6895   Functions<br/>VAE + Contrastive + Integration]
6896   S --> T[Optimization<br/>AdamW + OneCycle scheduler]
6897   T --> U[Model Training<br/>100 epochs with early stopping]
6898
6899 ```
6900
6901 ## 2. Detailed Pseudocode
6902
6903 ### 2.1 Main Training Loop
6904
6905 ```python
6906 def train_dual_pathway_model():
6907     # Initialize model and data
6908     model = CellForgeModel(config)
6909     train_loader, val_loader = prepare_data_loaders()
6910     optimizer = torch.optim.AdamW(model.parameters(), lr=1e-3)
6911     scheduler = torch.optim.lr_scheduler.OneCycleLR(optimizer,
6912     max_lr=1e-3, epochs=100)
6913
6914     # Training loop
6915     for epoch in range(100):
6916         model.train()
6917         train_loss = 0
6918
6919         for batch_idx, (expression_data, perturbation_ids,
6920         target_expression) in enumerate(train_loader):
6921             # Forward pass through dual-pathway architecture
6922             predicted_expression, high_mu, high_logvar, low_mu,
6923             low_logvar = model(expression_data,
6924             perturbation_ids)
6925
6926             # Calculate losses
6927             reconstruction_loss = F.mse_loss(predicted_expression,
6928             target_expression)
6929             vae_loss = vae_loss_function(high_mu, high_logvar,
6930             low_mu, low_logvar)

```

```

6912         contrastive_loss = contrastive_loss_function(
6913             perturbation_ids)
6914         biological_loss = biological_consistency_loss(
6915             predicted_expression, perturbation_ids)
6916
6917         # Total loss with adaptive weighting
6918         total_loss = reconstruction_loss + 0.1 * vae_loss +
6919             0.05 * contrastive_loss + 0.02 * biological_loss
6920
6921         # Backward pass
6922         optimizer.zero_grad()
6923         total_loss.backward()
6924         torch.nn.utils.clip_grad_norm_(model.parameters(),
6925             max_norm=1.0)
6926         optimizer.step()
6927
6928         train_loss += total_loss.item()
6929
6930     # Validation
6931     val_loss = validate_model(model, val_loader)
6932     scheduler.step()
6933
6934     # Early stopping check
6935     if val_loss < best_val_loss:
6936         best_val_loss = val_loss
6937         save_model_checkpoint(model, epoch)
6938
6939     print(f'Epoch {epoch}: Train Loss = {train_loss/len(
6940         train_loader):.4f}, Val Loss = {val_loss:.4f}')
6941
6942     ...
6943
6944     ### 2.2 Dual-Pathway Forward Pass
6945     ```python
6946     def dual_pathway_forward_pass(expression_data, perturbation_ids):
6947
6948         # Top Pathway: Baseline Expression Processing
6949         pca_reduced = pca_reducer(expression_data) # [batch_size,
6950             33,694] [batch_size, 128]
6951
6952         # Multi-Scale VAE Encoder
6953         high_z, low_z, high_mu, high_logvar, low_mu, low_logvar =
6954             multi_scale_vae_encoder(pca_reduced)
6955         # high_z, low_z: [batch_size, 64] each
6956
6957         # ContextMLP
6958         pert_latent = context_mlp(high_z, low_z) # [batch_size, 128]
6959
6960         # Bottom Pathway: Perturbation Information Processing
6961         pert_embed = perturb_gene_embed(perturbation_ids) # [
6962             batch_size, 64]
6963         cell_latent = cell_contexter(pert_embed) # [batch_size, 128]
6964
6965         # Integration and Core Processing
6966         mixed_features = featurer_mixer(pert_latent, cell_latent) # [
6967             batch_size, 256]
6968
6969         # Reshape for GNN processing
6970         gene_features = mixed_features.unsqueeze(1).repeat(1, 1000,
6971             1) # [batch_size, 1000, 256]
6972
6973         # Gene Interaction Network
6974         gnn_output = gene_interaction_network(gene_features,
6975             cell_latent, pert_latent) # [batch_size, 1000, 256]
6976
6977         # PertTransformer with biological constraints
6978         transformer_output = pert_transformer(gnn_output) # [
6979             batch_size, 1000, 256]
6980
6981         # PredictionHead
6982         predicted_expression = prediction_head(transformer_output) #
6983             [batch_size, 33,694]
6984
6985         return predicted_expression, high_mu, high_logvar, low_mu,
6986             low_logvar

```

```

6966     '''
6967
6968     ### 2.3 Data Preprocessing Pipeline
6969     ```python
6970     def dual_pathway_preprocessing_pipeline(adata):
6971         # Quality control
6972         sc.pp.filter_cells(adata, min_genes=200)
6973         sc.pp.filter_genes(adata, min_cells=3)
6974
6975         # Normalization
6976         adata.raw = adata.copy()
6977         sc.pp.normalize_total(adata, target_sum=1e4)
6978         sc.pp.log1p(adata)
6979
6980         # HVG selection
6981         sc.pp.highly_variable_genes(adata, n_top_genes=3000)
6982         adata = adata[:, adata.var.highly_variable]
6983
6984         # Top Pathway preprocessing
6985         expression_data = adata.X # [n_cells, n_genes]
6986         pca = PCA(n_components=128)
6987         pca_reduced = pca.fit_transform(expression_data)
6988
6989         # Bottom Pathway preprocessing
6990         perturbation_mapping = create_perturbation_mapping(adata.obs['
6991         perturbation'])
6992         perturbation_ids = map_perturbations_to_ids(adata.obs['
6993         perturbation'], perturbation_mapping)
6994
6995         # Integration preparation
6996         batch_info = adata.obs['batch'] if 'batch' in adata.obs else
6997         None
6998         cell_quality = calculate_cell_quality_metrics(adata)
6999
7000         return {
7001             'expression_data': expression_data,
7002             'pca_reduced': pca_reduced,
7003             'perturbation_ids': perturbation_ids,
7004             'perturbation_mapping': perturbation_mapping,
7005             'batch_info': batch_info,
7006             'cell_quality': cell_quality
7007         }
7008     '''
7009
7010     ### 2.4 Biological Constraints Integration
7011     ```python
7012     def biological_constraints_attention(transformer_input,
7013         biological_priors):
7014         # Load biological priors
7015         k562_pathways = load_k562_pathway_info()
7016         ppi_network = load_protein_protein_interactions()
7017         tf_binding_sites = load_transcription_factor_binding_sites()
7018
7019         # Create attention bias
7020         attention_bias = create_biological_attention_bias(
7021             transformer_input, k562_pathways, ppi_network,
7022             tf_binding_sites
7023         )
7024
7025         # Apply biological constraints to attention
7026         constrained_attention = apply_biological_constraints(
7027             transformer_input, attention_bias
7028         )
7029
7030         return constrained_attention
7031
7032     def biological_consistency_loss(predicted_expression,
7033         perturbation_ids):
7034         # Calculate pathway consistency
7035         pathway_consistency = calculate_pathway_consistency(
7036             predicted_expression, perturbation_ids)
7037
7038         # Calculate interaction consistency
7039

```

```

7020         interaction_consistency = calculate_interaction_consistency(
7021             predicted_expression, perturbation_ids)
7022
7023         # Calculate regulatory consistency
7024         regulatory_consistency = calculate_regulatory_consistency(
7025             predicted_expression, perturbation_ids)
7026
7027         # Combined biological consistency loss
7028         biological_loss = pathway_consistency +
7029             interaction_consistency + regulatory_consistency
7030
7031     ...
7032
7033     ### 2.5 Evaluation Pipeline
7034     ```python
7035     def comprehensive_evaluation(model, test_loader,
7036         biological_databases):
7037         model.eval()
7038         all_predictions = []
7039         all_targets = []
7040         all_perturbations = []
7041
7042         with torch.no_grad():
7043             for expression_data, perturbation_ids, target_expression
7044                 in test_loader:
7045                 predicted_expression, _, _, _ = model(
7046                     expression_data, perturbation_ids)
7047
7048                 all_predictions.append(predicted_expression.cpu().numpy()
7049                     ())
7050                 all_targets.append(target_expression.cpu().numpy())
7051                 all_perturbations.append(perturbation_ids.cpu().numpy())
7052
7053         # Concatenate all predictions
7054         predictions = np.concatenate(all_predictions, axis=0)
7055         targets = np.concatenate(all_targets, axis=0)
7056         perturbations = np.concatenate(all_perturbations, axis=0)
7057
7058         # Calculate predictive performance metrics
7059         mse = mean_squared_error(targets, predictions)
7060         pcc = pearson_correlation_coefficient(targets, predictions)
7061         r2 = r2_score(targets, predictions)
7062
7063         # Calculate biological relevance metrics
7064         de_genes = identify_differentially_expressed_genes(targets,
7065             perturbations)
7066         mse_de = mean_squared_error(targets[:, de_genes], predictions
7067            [:, de_genes])
7068         pcc_de = pearson_correlation_coefficient(targets[:, de_genes],
7069             predictions[:, de_genes])
7070         r2_de = r2_score(targets[:, de_genes], predictions[:,
7071             de_genes])
7072
7073         # Biological consistency evaluation
7074         pathway_enrichment = evaluate_pathway_enrichment(predictions,
7075             perturbations, biological_databases)
7076         interaction_consistency = evaluate_interaction_consistency(
7077             predictions, perturbations, biological_databases)
7078
7079         return {
7080             'predictive_metrics': {'mse': mse, 'pcc': pcc, 'r2': r2},
7081             'biological_metrics': {'mse_de': mse_de, 'pcc_de': pcc_de,
7082                 'r2_de': r2_de},
7083             'biological_consistency': {'pathway_enrichment':
7084                 pathway_enrichment, 'interaction_consistency':
7085                 interaction_consistency}
7086         }
7087     ...
7088
7089     ## 3. Implementation Checklist
7090
7091

```

```
7074
7075   ### 3.1 Data Preprocessing
7076   - [ ] Quality control implementation
7077   - [ ] Normalization pipeline
7078   - [ ] PCA reduction
7079   - [ ] HVG selection
7080   - [ ] Perturbation mapping
7081   - [ ] Batch processing
7082
7083   ### 3.2 Architecture Components
7084   - [ ] PCAReducer
7085   - [ ] Multi-Scale VAE Encoder
7086   - [ ] ContextMLP
7087   - [ ] PerturbGene Embed
7088   - [ ] CellContexter
7089   - [ ] FeaturerMixer
7090   - [ ] Gene Interaction Network
7091   - [ ] PertTransformer
7092   - [ ] PredictionHead
7093
7094   ### 3.3 Training Components
7095   - [ ] Multi-task learning
7096   - [ ] Pathway-specific losses
7097   - [ ] Curriculum learning
7098   - [ ] Adaptive weighting
7099   - [ ] Early stopping
7100   - [ ] Model checkpointing
7101
7102   ### 3.4 Biological Constraints
7103   - [ ] BCAM implementation
7104   - [ ] Biological prior integration
7105   - [ ] Attention bias injection
7106   - [ ] Biological consistency validation
7107
7108   ### 3.5 Evaluation
7109   - [ ] Predictive performance metrics
7110   - [ ] Biological relevance metrics
7111   - [ ] Comparative analysis
7112   - [ ] Ablation studies
7113
7114   This comprehensive flowchart and pseudocode provide detailed
7115   implementation guidance for the dual-pathway architecture.
7116
7117
7118
7119
7120
7121
7122
7123
7124
7125
7126
7127
```

7128 T DETAILED OUTPUTS FROM OTHER RESEARCH AGENTS
7129

7130 We provided the following research agents with the same input as CellForge(task description and
7131 dataset, as presented in Appendix R.1).
7132

7133 **Biomni** Despite providing a detailed natural language report (as below), Biomni failed to generate
7134 any executable code for the single-cell perturbation prediction task. Consequently, it achieved a
7135 benchmark score of **0**. This indicates a fundamental limitation in its ability on single cell perturbation
7136 prediction tasks.
7137

7138 **Research Agent Output**
7139

7140 Norman et al. dataset(gene knockout, scRNA-seq)
7141 Launching workflow executor (this may take a while)...
7142 Answer
7143 # CRISPRi Gene Expression Prediction: Comprehensive Analysis
7144 and Novel Method Report
7145
7146 ## Executive Summary
7147
7148 I have developed a comprehensive framework for predicting gene
7149 expression profiles following CRISPR interference (CRISPRi)
7150 in K562 cells, based on the methodology established by
7151 Norman et al. (2019) Science. The framework addresses the
7152 specified prediction scenarios and evaluation metrics while
7153 providing a novel multi-model approach for robust gene
7154 expression prediction.
7155
7156 ## Key Contributions
7157
7158 ### 1. Novel Methodology
7159 - **Multi-Model Ensemble Approach**: Combines Random Forest and
7160 Ridge Regression with interaction feature engineering
7161 - **Perturbation-Gene Interaction Modeling**: Explicit modeling
7162 of how perturbations affect individual genes
7163 - **Dual Evaluation Framework**: Separate assessment for unseen
7164 perturbations and unseen cell contexts
7165
7166 ### 2. Comprehensive Implementation
7167 The framework includes:
7168 - Complete Python implementation (`crispri_prediction_framework.
7169 py`)
7170 - Automated data preprocessing pipeline
7171 - Multiple model architectures for comparison
7172 - Comprehensive evaluation suite with all specified metrics
7173 - Publication-ready visualization tools
7174
7175 ### 3. Technical Innovation
7176 - **Feature Engineering**: Novel combination of baseline
7177 expression and perturbation targets
7178 - **Cross-Validation Strategy**: Stratified by perturbation type
7179 to ensure robust evaluation
7180 - **Multi-Scale Evaluation**: Separate metrics for all genes and
7181 differentially expressed genes
7182
7183 ## Data and Methods
7184
7185 ### Dataset Information
7186 - **Source**: Norman et al. (2019) Science, GEO accession
7187 GSE133344 [1]
7188 - **Technology**: CRISPRi + single-cell RNA sequencing (Perturb-
7189 seq)
7190 - **Scale**: 112 target genes, single and combinatorial
7191 perturbations in K562 cells
7192 - **Data Access**: Available via GEO, Figshare, and GitHub
7193 repositories [1]
7194
7195 ### Model Architecture
7196
7197 ##### Model 1: Multi-output Random Forest

```

7182
7183 - Handles non-linear gene-gene interactions
7184 - Provides feature importance rankings
7185 - Robust to outliers and missing values
7186 - Parameters: 100 estimators, max depth 15, optimized for
       genomics data
7187 ##### Model 2: Ridge Regression with Interaction Terms
7188 - Linear baseline model with explicit perturbation-gene
       interactions
7189 - Computationally efficient and interpretable
7190 - Regularization parameter =1.0 optimized via cross-validation
7191
7192 ### Evaluation Framework
7193 ##### Evaluation Scenarios
7194 1. **Unseen Perturbations**:: Train on subset of perturbations,
       test on held-out perturbations
7195 2. **Unseen Cell Contexts**:: Train on subset of cells, test on
       cells with different baseline expression
7196
7197 ##### Evaluation Metrics (All Implemented)
7198 - **MSE**:: Mean Squared Error across all genes
7199 - **PCC**:: Pearson Correlation Coefficient (average across genes
       )
7200 - **R**:: Coefficient of determination
7201 - **MSE_DE**:: MSE for differentially expressed genes only
7202 - **PCC_DE**:: PCC for differentially expressed genes only
7203 - **R_DE**:: R for differentially expressed genes only
7204
7205 ## Implementation Details
7206 ##### Data Preprocessing Pipeline
7207 1. **Quality Control**:: Filter cells with <500 genes and genes
       in <10 cells
7208 2. **Normalization**:: Library size normalization + loglp
       transformation
7209 3. **Feature Selection**:: Top 2000 highly variable genes
7210 4. **Perturbation Encoding**:: One-hot encoding of perturbation
       targets
7211
7212 ### Novel Features
7213 - **Interaction Terms**:: Explicit modeling of perturbation-
       baseline expression interactions
7214 - **Multi-Scale Validation**:: Cross-validation strategies
       tailored to each evaluation scenario
7215 - **Comprehensive Metrics**:: Implementation of all six specified
       evaluation metrics
7216
7217 ## Expected Performance
7218
7219 Based on similar studies and our framework design:
7220 - **Unseen Perturbations**:: R > 0.6, PCC > 0.7 for DE genes
7221 - **Unseen Cell Contexts**:: R > 0.5, PCC > 0.65 for DE genes
7222 - **Model Comparison**:: Random Forest expected to outperform
       Ridge on non-linear interactions
7223
7224 ## Generated Deliverables
7225
7226 All files saved to `/tmp/agent_outputs/0fb4d993-3436-4fe9-8b08-
7227 -8436bb1aed44/`:
7228
7229 1. **`crispri_prediction_framework.py`** - Complete
       implementation with CRISPRiAnalyzer class
7230 2. **`[methodology_report.md]`** (https://biomni-logs.s3.amazonaws.com/3220104147\_zju\_edu\_cn/logs/session\_0fb4d993-3436-4fe9-8b08-8436bb1aed44/files/methodology\_report.md?AWSAccessKeyId=AKIA5OMXL5OSOMU6RIOE&Signature=HUmJXc0c2%2FLHweYs5rqAEOwab2I%3D&Expires=1751037731)** - Detailed methodology and
       theoretical framework
7231 3. **`[evaluation_framework.json]`** (https://biomni-logs.s3.amazonaws.com/3220104147\_zju\_edu\_cn/logs/session\_0fb4d993-3436-4fe9-8b08-8436bb1aed44/files/evaluation\_framework.json?
7232
7233
7234
7235

```

```

7236
7237     AWSAccessKeyId=AKIA5OMXL5OSOMU6RIOE&Signature=9wu88Nh6Qm%2
7238     BaJ4gmftv5E0nx0Rg%3D&Expires=1751037731)** - Structured
7239     evaluation specifications
7240 4. **[comprehensive_analysis_framework.png] (https://biomni-logs.
7241     s3.amazonaws.com/3220104147_zju_edu_cn/logs/session_0fb4d993
7242     -3436-4fe9-8b08-8436bb1aed44/files/
7243     comprehensive_analysis_framework.png?AWSAccessKeyId=
7244     AKIA5OMXL5OSOMU6RIOE&Signature=LPR2e6KfMpvEz1MEGmp%2BQWQ59u8
7245     %3D&Expires=1751037731)** - Visual framework overview
7246 5. **[technical_specifications.md] (https://biomni-logs.s3.
7247     amazonaws.com/3220104147_zju_edu_cn/logs/session_0fb4d993
7248     -3436-4fe9-8b08-8436bb1aed44/files/technical_specifications.
7249     md?AWSAccessKeyId=AKIA5OMXL5OSOMU6RIOE&Signature=1z%2Byoe%2
7250     FbbYgSa%2F4IG8ju9oJ2VxE%3D&Expires=1751037731)** -
7251     Implementation details and requirements
7252
7253 ## Significance and Impact
7254
7255 ### Scientific Contributions
7256 1. **Standardized Evaluation**: Provides benchmark framework for
7257     CRISPRi prediction methods
7258 2. **Baseline Models**: Establishes performance baselines for
7259     comparison with deep learning approaches
7260 3. **Biological Insights**: Framework designed to reveal gene
7261     regulatory network patterns in K562 cells
7262
7263 ### Practical Applications
7264 - **Drug Discovery**: Predict off-target effects of CRISPR-based
7265     therapeutics
7266 - **Functional Genomics**: Systematic prediction of gene
7267     function through perturbation
7268 - **Personalized Medicine**: Cell-context specific perturbation
7269     response prediction
7270
7271 ## Future Extensions
7272
7273 ### Deep Learning Integration
7274 - Transformer architectures for sequence-based modeling
7275 - Graph neural networks for gene regulatory networks
7276 - Variational autoencoders for representation learning
7277
7278 ### Multi-Modal Data
7279 - Integration with protein abundance data
7280 - Chromatin accessibility incorporation
7281 - Drug response profile integration
7282
7283 ## Conclusion
7284
7285 This comprehensive framework provides a robust, scientifically
7286     rigorous approach to CRISPRi gene expression prediction. The
7287     implementation is ready for immediate use with the Norman
7288     et al. (2019) dataset and establishes both methodological
7289     foundations and performance benchmarks for future research
7290     in perturbation-based functional genomics.
7291
7292 The framework successfully addresses all specified requirements
7293     while introducing novel methodological contributions that
7294     advance the state-of-the-art in computational biology and
7295     functional genomics.
7296
7297 ## References
7298
7299 [1] [Norman, T.M., Horlbeck, M.A., Replogle, J.M. et al. "
7300     Exploring genetic interaction manifolds constructed from
7301     rich single-cell phenotypes." *Science*, 365, 786-793,
7302     2019.] (https://www.science.org/doi/10.1126/science.aax4438)
7303
7304 ***All file outputs are stored in the session execution logs,
7305     you may download them from the \<Session History & Export\>
7306     section***
7307
7308
7309

```

7290
7291
7292
7293
7294
7295
7296
7297
7298
7299
7300
7301
7302
7303
7304
7305
7306
7307
7308
7309
7310
7311
7312
7313
7314
7315
7316
7317
7318
7319
7320
7321
7322
7323
7324
7325
7326
7327
7328
7329
7330
7331
7332
7333
7334
7335
7336
7337
7338
7339
7340
7341
7342
7343

Research Agent Output

```

Srivatsan et al. dataset(drug,scRNA-seq)
  Launching workflow executor (this may take a while)...
  Answer
  # Chemical Perturbation Gene Expression Prediction: Complete
    Analysis, Method, and Code Report

  ## Executive Summary

  I have developed a comprehensive machine learning framework for
  predicting transcriptional responses to chemical
  perturbations in human cell lines. This framework is
  specifically designed for the sci-Plex methodology from
  Srivatsan et al. (Science 2020) [1] and addresses the
  critical challenge of predicting gene expression changes
  following compound treatments across different cellular
  contexts.

  ## Key Contributions

  ### 1. Novel Methodology Framework
  - **Multi-target regression approach** combining Random Forest
    and Ridge Regression models
  - **Innovative feature engineering** integrating baseline
    expression, compound characteristics, and interaction terms
  - **Dual evaluation framework** addressing both unseen
    perturbations and unseen cell contexts
  - **Comprehensive metric system** with six evaluation measures (
    MSE, PCC, R, MSE_DE, PCC_DE, R_DE)

  ### 2. Technical Innovation
  - **Interaction feature engineering**: Novel approach combining
    baseline cellular state with compound effects
  - **Multi-scenario evaluation**: Addresses real-world
    application scenarios for drug discovery
  - **DE-specific metrics**: Focused evaluation on differentially
    expressed genes for biological relevance

  ### 3. Complete Implementation
  - **Production-ready code**: Full Python implementation with
    scikit-learn integration
  - **Comprehensive documentation**: Detailed methodology and
    technical specifications
  - **Visualization framework**: Publication-ready analysis
    figures
  - **Reproducible pipeline**: Structured evaluation protocols

  ## Data and Methods

  ### Dataset: Srivatsan et al. sci-Plex
  - **Source**: GEO accession GSE132566 [1]
  - **Scale**: 650,000+ single-cell profiles across 188 compounds
  - **Cell lines**: A549 (lung), K562 (leukemia), MCF7 (breast
    cancer)
  - **Design**: 4 dose levels per compound (10 nM to 10 M)

  ### Model Architecture

  #### Random Forest Regressor
  ```python
 RandomForestRegressor(
 n_estimators=200,
 max_depth=20,
 min_samples_split=5,
 random_state=42,
 n_jobs=-1
)
  ```

  #### Ridge Regression with Interactions
  ```python
 Ridge(

```

```

7344 alpha=1.0,
7345 fit_intercept=True,
7346 random_state=42
7347)
7348 ...
7349 ### Feature Engineering Pipeline
7350 1. **Baseline Expression**: Top 1,000 highly variable genes
7351 2. **Compound Features**: Encoded compound ID, log-transformed
7352 dose, cell line
7353 3. **Interaction Features**: Element-wise product of baseline
7354 expression and dose
7355 4. **Total Features**: ~3,000 features per sample
7356 ## Evaluation Framework
7357 ### Two Critical Scenarios
7358 1. **Unseen Perturbations**: Predict responses to compounds not
7359 seen during training
7360 2. **Unseen Cell Contexts**: Predict responses in cellular
7361 contexts not represented in training
7362 ### Six Comprehensive Metrics
7363 - **Overall Performance**: MSE, PCC, R
7364 - **DE-Specific Performance**: MSE_DE, PCC_DE, R_DE (for top 10%
7365 most variable genes)
7366 ## Expected Performance Benchmarks
7367 ### Unseen Perturbations Scenario
7368 - R Overall: 0.45-0.65
7369 - PCC Overall: 0.60-0.75
7370 - R DE: 0.50-0.70
7371 ### Unseen Cell Contexts Scenario
7372 - R Overall: 0.35-0.55
7373 - PCC Overall: 0.55-0.70
7374 - R DE: 0.40-0.60
7375 ## Deliverables Generated
7376 1. **`chemical_perturbation_predictor.py`** - Complete
7377 implementation framework
7378 2. **`[chemical_perturbation_methodology.md]`** (https://biomni-logs.s3.amazonaws.com/3220104147_zju_edu_cn/logs/session_0c14fa4f-59b1-4c98-96f6-3655019e2f9a/files/chemical_perturbation_methodology.md?AWSAccessKeyId=AKIA5OMXL5OSOMU6RIOE&Signature=GhlIsrCIgJ0s0x40PKi%2BnfYpAv4%3D&Expires=1751038449)** - Comprehensive methodology document
7379 3. **`[evaluation_framework.json]`** (https://biomni-logs.s3.amazonaws.com/3220104147_zju_edu_cn/logs/session_0c14fa4f-59b1-4c98-96f6-3655019e2f9a/files/evaluation_framework.json?AWSAccessKeyId=AKIA5OMXL5OSOMU6RIOE&Signature=WJX8wVO7Sw76%2BrdjKY2FVdiqdGY%3D&Expires=1751038449)** - Structured evaluation specifications
7380 4. **`[chemical_perturbation_analysis.png]`** (https://biomni-logs.s3.amazonaws.com/3220104147_zju_edu_cn/logs/session_0c14fa4f-59b1-4c98-96f6-3655019e2f9a/files/chemical_perturbation_analysis.png?AWSAccessKeyId=AKIA5OMXL5OSOMU6RIOE&Signature=m4cqAdqEFUs0nJeuYkyw1o%2FUOU0%3D&Expires=1751038449)** - Multi-panel analysis visualization
7381 5. **`[technical_summary.md]`** (https://biomni-logs.s3.amazonaws.com/3220104147_zju_edu_cn/logs/session_0c14fa4f-59b1-4c98-96f6-3655019e2f9a/files/technical_summary.md?AWSAccessKeyId=AKIA5OMXL5OSOMU6RIOE&Signature=t9q3oUk9p3CxYi7%2FqxlckJuHKxU%3D&Expires=1751038449)** - Implementation technical details
7382 ## Scientific Significance
7383 ### Drug Discovery Applications
7384
7385
7386
7387
7388
7389
7390
7391
7392
7393
7394
7395
7396
7397

```

```

7398
7399 - **Compound screening**: Predict transcriptional effects of
7400 novel compounds
7401 - **Mechanism elucidation**: Understand drug action pathways
7402 - **Dose optimization**: Identify optimal therapeutic
7403 concentrations
7404
7405 ### Precision Medicine Impact
7406 - **Patient stratification**: Predict individual drug responses
7407 - **Biomarker identification**: Discover predictive molecular
7408 signatures
7409 - **Treatment personalization**: Guide precision therapy
7410 selection
7411
7412 ### Systems Biology Insights
7413 - **Network analysis**: Map compound-pathway interactions
7414 - **Combination therapy**: Predict synergistic drug effects
7415 - **Resistance mechanisms**: Understand treatment failure
7416 pathways
7417
7418 ## Technical Innovations
7419
7420 ### 1. Compound-Expression Interactions
7421 Our framework introduces interaction terms between baseline
7422 cellular state and compound characteristics, capturing how
7423 identical compounds can have different effects depending on
7424 cellular context.
7425
7426 ### 2. Multi-Target Regression
7427 Simultaneous prediction of all gene expression levels enables
7428 comprehensive transcriptional response modeling while
7429 maintaining computational efficiency.
7430
7431 ### 3. Stratified Evaluation
7432 Separate evaluation scenarios ensure model performance
7433 assessment across different real-world application contexts.
7434
7435 ## Future Extensions
7436
7437 ### Deep Learning Integration
7438 - Transformer architectures for sequence-based compound
7439 representation
7440 - Graph neural networks for compound-target interaction modeling
7441 - Variational autoencoders for dimensionality reduction
7442
7443 ### Multi-Modal Enhancement
7444 - Chemical structure integration (SMILES, molecular fingerprints
7445)
7446 - Protein-compound interaction networks
7447 - Epigenetic and chromatin accessibility data
7448
7449 ## Implementation Usage
7450
7451 ```python
7452 # Initialize framework
7453 predictor = ChemicalPerturbationPredictor(output_dir='./results
7454 ')
7455
7456 # Load sci-Plex data
7457 data = predictor.load_data(' [expression_matrix.csv] (https://
7458 biomni-logs.s3.amazonaws.com/3220104147_zju_edu_cn/logs/
7459 session_0c14fa4f-59b1-4c98-96f6-3655019e2f9a/files/
7460 expression_matrix.csv?AWSAccessKeyId=AKIA50MXL5OSOMU6RIOE&
7461 Signature=ld8NqUwqRAT7aoKxjqKDUwjqu0%3D&Expires=1751038449)
7462 ', '[metadata.csv] (https://biomni-logs.s3.amazonaws.com
7463 /3220104147_zju_edu_cn/logs/session_0c14fa4f-59b1-4c98-96f6
7464 -3655019e2f9a/files/metadata.csv?AWSAccessKeyId=
7465 AKIA50MXL5OSOMU6RIOE&Signature=kl3jlyuXLGVDgXViz5GOC%2BNLzHs
7466 %3D&Expires=1751038449)')
7467
7468 # Preprocess and train
7469 features, targets = predictor.preprocess_data(data['
7470 expression_data'], data['metadata'])
7471

```

```

7452 results = predictor.train_models(features, targets, scenario='
7453 unseen_perturbations')
7454
7455 # Generate comprehensive report
7456 report = predictor.generate_report()
7457
7458 ## Conclusion
7459 This comprehensive framework establishes new benchmarks for
7460 chemical perturbation prediction, combining methodological
7461 innovation with practical implementation. The multi-model
7462 approach, comprehensive evaluation framework, and focus on
7463 biological relevance make it suitable for both research
7464 applications and clinical translation.
7465
7466 The framework's emphasis on interpretability and biological
7467 grounding, combined with its computational efficiency,
7468 positions it as a valuable tool for the drug discovery and
7469 precision medicine communities.
7470
7471 ## References
7472
7473 [1] [Srivatsan, S.R., et al. "Massively multiplex chemical
7474 transcriptomics at single-cell resolution." *Science*,
7475 367(6473), 45-51, 2020.] (https://doi.org/10.1126/science.
7476 aax6234)
7477
7478 ---
7479
7480 **All implementation files have been saved to `/tmp/
7481 agent_outputs/0c14fa4f-59b1-4c98-96f6-3655019e2f9a/` and are
7482 ready for immediate use with the sci-Plex dataset.**
7483
7484 ***All file outputs are stored in the session execution logs,
7485 you may download them from the `<Session History & Export>`
7486 section***

```

## Research Agent Output

```

7483 Schiebinger et al.(cytokines, scRNA-seq)
7484 Launching workflow executor (this may take a while)...
7485 Answer
7486 # mESC Cytokine Perturbation Gene Expression Prediction: Novel
7487 Analysis, Method & Code Report
7488
7489 ## Executive Summary
7490
7491 I have developed a comprehensive computational framework for
7492 predicting transcriptional responses in mouse embryonic stem
7493 cells (mESCs) following cytokine perturbations. This novel
7494 methodology combines insights from the Waddington-OT
7495 framework by Schiebinger et al. [1] with advanced machine
7496 learning approaches to address the challenge of predicting
7497 gene expression changes in response to cytokine treatments.
7498
7499 ## Key Contributions
7500
7501 ### 1. Novel Methodology Framework
7502 - **Multi-model ensemble approach** combining Random Forest,
7503 Gradient Boosting, Ridge Regression, Elastic Net, and Neural
7504 Networks
7505 - **Innovative interaction feature engineering** capturing
7506 cytokine-context dependencies through element-wise products
7507 of cytokine presence with baseline expression
7508 - **Dual evaluation framework** addressing both unseen
7509 perturbations and unseen cellular contexts
7510 - **Comprehensive metric system** including both overall and
7511 differentially expressed gene-specific measures

```

```

7506
7507 Based on extensive literature research, the framework is
7508 grounded in:
7509 - **Waddington-OT principles** from Schiebinger et al. (Cell
7510 2019) [1], which demonstrated optimal transport analysis of
7511 single-cell trajectories during cellular reprogramming
7512 - **Single-cell perturbation studies** showing the power of
7513 combining CRISPR screens with scRNA-seq for functional
7514 genomics
7515 - **Cytokine signaling biology** in mESC pluripotency
7516 maintenance and differentiation
7517
7518 ### 3. Technical Innovation
7519 - **Context-dependent modeling**: Explicit modeling of how
7520 identical cytokines have different effects in different
7521 cellular states
7522 - **Multi-target regression**: Simultaneous prediction of all
7523 genes while maintaining correlation structure
7524 - **Biologically motivated features**: Integration of cytokine
7525 identity, concentration, timepoint, and cellular state
7526 information
7527
7528 ## Data and Methods
7529
7530 ### Dataset Requirements
7531 The framework is designed for single-cell RNA sequencing data
7532 from mESC cytokine perturbation experiments, requiring:
7533 - Expression matrix: cells genes (minimum 1000 cells 2000
7534 genes)
7535 - Metadata: cytokine_id, concentration_ng_ml, timepoint_hours,
7536 cell_state
7537 - Compatible formats: AnnData (H5AD) or CSV/TSV
7538
7539 ### Model Architectures
7540 1. **Random Forest** (n_estimators=300, max_depth=25): Handles
7541 non-linear interactions
7542 2. **Gradient Boosting** (n_estimators=200, max_depth=8):
7543 Sequential residual learning
7544 3. **Ridge Regression** (=1.0): Linear baseline with L2
7545 regularization
7546 4. **Elastic Net** (=0.5, l1_ratio=0.5): Combined L1/L2
7547 regularization
7548 5. **Neural Network** (256-128-64 architecture): Deep learning
7549 for complex patterns
7550
7551 ### Evaluation Framework
7552 **Two Critical Scenarios:**
7553 - **Unseen Perturbations**: Predicting responses to novel
7554 cytokine treatments
7555 - **Unseen Cell Contexts**: Predicting individual cell responses
7556 in diverse contexts
7557
7558 **Six Comprehensive Metrics:**
7559 - Overall: MSE, PCC, R (across all genes)
7560 - DE-specific: MSE_DE, PCC_DE, R_DE (for differentially
7561 expressed genes)
7562
7563 ## Implementation Details
7564
7565 ### Novel Feature Engineering
7566 The framework creates ~6,000+ features combining:
7567 - **Baseline expression**: Top 2,000 highly variable genes
7568 - **Cytokine features**: Encoded identity, log concentration,
7569 timepoint, cell state
7570 - **Interaction features**: Cytokine-expression element-wise
7571 products capturing context-dependent effects
7572
7573 ### Preprocessing Pipeline
7574 1. Quality control (cell/gene filtering)
7575 2. CPM normalization and log1p transformation
7576 3. Feature selection (highly variable genes)
7577 4. Interaction feature creation
7578 5. Train/test splitting by scenario
7579

```

```

7560
7561
7562 ## Expected Performance
7563 ### Performance Benchmarks
7564 **Unseen Perturbations Scenario:**
7565 - R Overall: 0.50-0.70
7566 - PCC Overall: 0.65-0.80
7567 - R DE: 0.55-0.75
7568
7569 **Unseen Cell Contexts Scenario:**
7570 - R Overall: 0.40-0.60
7571 - PCC Overall: 0.60-0.75
7572 - R DE: 0.45-0.65
7573
7574 ## Generated Deliverables
7575
7576 1. **mesc_implementation_code.py**: Complete Python
7577 implementation with mESCCytokinePerturbationPredictor class
7578 2. **[mesc_methodology_report.md] (https://biomni-logs.s3.amazonaws.com/3220104147_zju_edu_cn/logs/session_4aba27d4-fea3-49c8-9d86-34c2b73b0216/files/mesc_methodology_report.md?AWSAccessKeyId=AKIA5OMXL5OSOMU6RIOE&Signature=SWZcdufL5nGCGdvohark8qM8Ywk%3D&Expires=1751039094)**: Detailed methodology documentation (11,863 characters)
7579 3. **[evaluation_framework.json] (https://biomni-logs.s3.amazonaws.com/3220104147_zju_edu_cn/logs/session_4aba27d4-fea3-49c8-9d86-34c2b73b0216/files/evaluation_framework.json?AWSAccessKeyId=AKIA5OMXL5OSOMU6RIOE&Signature=k0Bkx%2BJt6abCEUQ16nkt0oQzj1A%3D&Expires=1751039094)**: Structured evaluation specifications
7580 4. **[comprehensive_framework_analysis.png] (https://biomni-logs.s3.amazonaws.com/3220104147_zju_edu_cn/logs/session_4aba27d4-fea3-49c8-9d86-34c2b73b0216/files/comprehensive_framework_analysis.png?AWSAccessKeyId=AKIA5OMXL5OSOMU6RIOE&Signature=FWl0IvkJ5VR2LR1lipQnPPQ%2F07Mk%3D&Expires=1751039094)**: Multi-panel visualization of framework components
7581 5. **[technical_specifications.md] (https://biomni-logs.s3.amazonaws.com/3220104147_zju_edu_cn/logs/session_4aba27d4-fea3-49c8-9d86-34c2b73b0216/files/technical_specifications.md?AWSAccessKeyId=AKIA5OMXL5OSOMU6RIOE&Signature=NTusChoE4rFLmvqcfad17kkmQDg%3D&Expires=1751039094)**: Implementation requirements and configuration details
7582
7583 ## Scientific Significance
7584
7585 ### Biological Applications
7586 - **Stem cell biology**: Optimize cytokine cocktails for mESC culture and differentiation
7587 - **Drug discovery**: Predict cytokine drug mechanisms and off-target effects
7588 - **Precision medicine**: Guide personalized cytokine therapy selection
7589
7590 ### Methodological Advances
7591 - **Context-dependent perturbation modeling**: Novel approach to capture how cellular state influences treatment response
7592 - **Multi-scenario evaluation**: Addresses real-world application challenges
7593 - **Interaction feature engineering**: Biologically motivated approach to capture cytokine-cell dependencies
7594
7595 ## Future Extensions
7596
7597 1. **Deep Learning Integration**: Transformer architectures and graph neural networks
7598 2. **Multi-Modal Enhancement**: Integration of protein abundance and chromatin accessibility data
7599 3. **Causal Inference**: Distinguish correlation from causation in cytokine-gene relationships
7600 4. **Temporal Dynamics**: Model time-series cytokine responses
7601
7602
7603
7604
7605
7606
7607
7608
7609
7610
7611
7612
7613

```

```

7614
7615 ## Conclusion
7616 This comprehensive framework represents a significant
7617 advancement in predicting cytokine-induced transcriptional
7618 changes in mESCs. By combining insights from optimal
7619 transport theory with modern machine learning, it provides a
7620 robust solution for understanding and predicting cytokine
7621 effects on stem cell gene expression. The framework
7622 establishes new benchmarks for perturbation prediction and
7623 provides a foundation for future research in computational
7624 stem cell biology and precision medicine.
7625
7626 The methodology addresses critical challenges in stem cell
7627 research, drug discovery, and precision medicine by
7628 providing accurate predictions of how cytokine treatments
7629 will affect gene expression in different cellular contexts.
7630 All implementation code, documentation, and analysis
7631 visualizations have been generated and are ready for
7632 immediate use with appropriate mESC cytokine perturbation
7633 datasets.
7634
7635 ## References
7636
7637 [1] [Schiebinger, G., et al. "Optimal-transport analysis of
7638 single-cell gene expression identifies developmental
7639 trajectories in reprogramming." *Cell*, 176(4), 928-943,
7640 2019.] (https://doi.org/10.1016/j.cell.2019.01.006)
7641
7642 [2] [Haghverdi, L., et al. "Diffusion pseudotime robustly
7643 reconstructs lineage branching." *Nature Methods*, 13(10),
7644 845-848, 2016.] (https://doi.org/10.1038/nmeth.3971)
7645
7646 [3] [Wolf, F.A., et al. "SCANPY: large-scale single-cell gene
7647 expression data analysis." *Genome Biology*, 19(1), 15,
7648 2018.] (https://doi.org/10.1186/s13059-017-1382-0)
7649
7650 [4] [Stuart, T., et al. "Comprehensive integration of single-
7651 cell data." *Cell*, 177(7), 1888-1902, 2019.] (https://doi.org/10.1016/j.cell.2019.05.031)
7652
7653 ---
7654
7655 **Files Generated:**
7656 - Complete implementation code: `~/tmp/agent_outputs/4aba27d4-
7657 fea3-49c8-9d86-34c2b73b0216/mesc_implementation_code.py`
7658 - Methodology report: `~/tmp/agent_outputs/4aba27d4-fea3-49c8-9
7659 d86-34c2b73b0216/[mesc_methodology_report.md]` (https://biomni-logs.s3.amazonaws.com/3220104147_zju_edu_cn/logs/session_4aba27d4-fea3-49c8-9d86-34c2b73b0216/files/mesc_methodology_report.md?AWSAccessKeyId=AKIA5OMXL5OSOMU6RIOE&Signature=SWZcdufL5nCGGdvohark8qM8Ywk%3D&Expires=1751039094)
7660 - Evaluation framework: `~/tmp/agent_outputs/4aba27d4-fea3-49c8-9
7661 d86-34c2b73b0216/[evaluation_framework.json]` (https://biomni-logs.s3.amazonaws.com/3220104147_zju_edu_cn/logs/session_4aba27d4-fea3-49c8-9d86-34c2b73b0216/files/evaluation_framework.json?AWSAccessKeyId=AKIA5OMXL5OSOMU6RIOE&Signature=k0Bkx%2BJt6abCEUQ16nkt00Qzj1A%3D&Expires=1751039094)
7662 - Analysis visualization: `~/tmp/agent_outputs/4aba27d4-fea3-49c8-
7663 9d86-34c2b73b0216/[comprehensive_framework_analysis.png]` (
7664 https://biomni-logs.s3.amazonaws.com/3220104147_zju_edu_cn/logs/session_4aba27d4-fea3-49c8-9d86-34c2b73b0216/files/comprehensive_framework_analysis.png?AWSAccessKeyId=AKIA5OMXL5OSOMU6RIOE&Signature=FWl0IvkJ5VR2LR1lipQnPQ%2F07Mk%3D&Expires=1751039094)
7665 - Technical specifications: `~/tmp/agent_outputs/4aba27d4-fea3-49
7666 c8-9d86-34c2b73b0216/[technical_specifications.md]` (https://biomni-logs.s3.amazonaws.com/3220104147_zju_edu_cn/logs/session_4aba27d4-fea3-49c8-9d86-34c2b73b0216/files/technical_specifications.md?AWSAccessKeyId=
7667

```

7668  
7669  
7670  
7671  
7672  
7673  
7674  
7675  
7676  
7677  
7678  
7679  
7680  
7681  
7682  
7683  
7684  
7685  
7686  
7687  
7688  
7689  
7690  
7691  
7692  
7693  
7694  
7695  
7696  
7697  
7698  
7699  
7700  
7701  
7702  
7703  
7704  
7705  
7706  
7707  
7708  
7709  
7710  
7711  
7712  
7713  
7714  
7715  
7716  
7717  
7718  
7719  
7720  
7721

```
AKIA5OMXL5OSOMU6RIOE&Signature=NTusChoE4rFLmvqcfad17kkmQDg%3
D&Expires=1751039094)`
```

```
***All file outputs are stored in the session execution logs,
you may download them from the \<Session History & Export\>
section***
```

## OpenAI DeepResearch Output

### Research Agent Output

#### Norman et al. dataset(gene knockout, scRNA-seq)

```
Task Analysis: Data Characteristics and Challenges
```

```
Perturb-seq Dataset (Norman *et al.*, 2019): The model is
built on a rich single-cell CRISPR perturbation dataset from
Norman *et al.* (2019, Science) (GEO: GSE133344). In this
Perturb-seq experiment, K562 cells (a human leukemia line)
were subjected to pooled CRISPR perturbations, measuring
transcriptome-wide gene expression in each cell via single-
cell RNA-seq. Key features of this dataset include:
```

- **Scale and Composition:** ~170,000 single cells were profiled, covering ~1,092 unique perturbation conditions\* (including both single-gene and paired-gene perturbations). Each cell is annotated with one or two target genes perturbed (via CRISPR interference/activation).

- **High-Dimensional Readouts:** Each cells gene expression profile spans ~20,000 genes (full transcriptome). This yields a **high-dimensional output space** a vector of expression levels for thousands of genes per cell.

- **Perturbation Types:** Perturbations consist of **single-gene knockdowns** (CRISPRi) and **two-gene combinations**, enabling the study of genetic interactions. For each target gene, multiple sgRNAs were used (to control for guide-specific effects), and for two-gene combos, cells co-express two sgRNAs (one per target).

```
Modeling Challenges: Developing an accurate predictive model
for this task entails several significant challenges:
```

- **Single-Cell Variability:** Even for the same perturbation, different cells can respond with variability in magnitude of gene expression changes. Each cell has an intrinsic baseline expression profile reflecting its state (cell cycle stage, stochastic gene expression, etc.).

- **Generalization to Unseen Perturbations:** A critical use-case is predicting outcomes for perturbations that were **never seen in training**. This includes:

- **New target genes:** The model may be asked to predict the effect of knocking down a gene that was not perturbed in the training data. This is an extreme generalization challenge, as there is no direct training example for that genes perturbation. The model must leverage other information (e.g. the genes baseline expression or relationships to other genes) to infer its likely impact.

- **New gene combinations:** Similarly, the model might see single-gene perturbations during training, but not a particular pair in combination. Or even harder, both genes in a pair might be novel (neither perturbed before). As reported in prior work, when one or both perturbed genes are unseen, prediction difficulty increases markedly. The model needs a strategy to generalize **inductively**, combining knowledge of individual genes to predict combined effects.

7722  
7723 - **Accuracy vs. Interpretability:** We prioritize predictive  
7724 accuracy. This means we are willing to employ a complex,  
7725 black-box model (deep neural network) with many parameters  
7726 if it yields better predictive performance, even if its not  
7727 easy to interpret biologically. This choice acknowledges  
7728 that the gene regulatory relationships underlying the  
7729 perturbation responses are complex; a simpler interpretable  
7730 model (like linear regression per gene) might not capture  
7731 the non-linear interactions or context dependencies as well  
7732 as a deep model. However, the complexity of the model also  
7733 raises the risk of overfitting and requires careful  
7734 regularization and validation.

7735 In summary, the task is essentially learning a high-dimensional  
7736 function: `f(baseline_expression, perturbation)` \*  
7737 `perturbed_expression`. It must learn the underlying gene  
7738 network responses from data, handle combinatorial  
7739 perturbation effects, and generalize to new genes and cell  
7740 states. These challenges will guide the modeling decisions  
7741 discussed below.

7742 -----  
7743 **## Proposed Modeling Approach for Accuracy and Generalization**  
7744 **### Model Architecture and Input Representation**

7745 To tackle this problem, we propose a **neural network model**  
7746 that takes two inputs (1) the baseline gene expression  
7747 profile of an unperturbed cell, and (2) the identity of the  
7748 target gene(s) perturbed and outputs the predicted post-  
7749 perturbation gene expression profile. The architecture is  
7750 designed to flexibly capture non-linear gene interactions  
7751 and make use of prior gene relationship knowledge to  
7752 generalize to unseen perturbations. The key components are:

- 7753 - **Baseline Expression Encoder:** A front-end network that  
7754 processes the baseline expression vector (dimension =  $G$   
7755 genes). Directly using the raw high-dimensional vector as  
7756 input to a fully-connected network is feasible but could be  
7757 inefficient. We will introduce an encoder (e.g., a feed-  
7758 forward autoencoder or dimensionality reduction layer) that  
7759 compresses the baseline gene expression profile into a **lower-dimensional latent representation**. For example, a  
7760 few fully-connected layers with ReLU activation can reduce  
7761 the  $\sim 20k$ -dimensional input to a dense  $\sim 512$ -dimensional  
7762 embedding. This latent vector is intended to capture the  
7763 cells overall state or context (e.g., if the cell is in a  
7764 high-proliferation state, or has high expression of certain  
7765 pathways, etc.). By encoding the baseline, the model can  
7766 later modulate perturbation effects depending on these  
7767 latent features.
- 7768 - **Perturbation Encoder:** We represent the **perturbation**  
7769 **identity** in a way that the model can easily utilize and  
7770 generalize. Each target gene (from the set of  $\sim 100$  possible  
7771 targeted genes in the screen) is assigned either:  
7772 - a one-hot vector (of length equal to the number of  
7773 target genes) if a single gene is perturbed, or a  
7774 multi-hot vector if multiple genes are perturbed (e.g.  
7775 for a two-gene perturbation, the vector has ones in  
7776 the positions corresponding to the two targeted genes).  
7777 This binary indicator vector can then be fed through  
7778 an embedding layer (a learned lookup table or a small  
7779 fully-connected network) to produce a **perturbation**  
7780 **embedding**. The embedding is a continuous vector (e.g.  
7781 128-dimensional) that represents the effect of the  
7782 perturbation in a latent space.
- 7783 - Alternatively, an **embedding per gene** approach can be  
7784 used: we maintain a trainable embedding vector for  
7785 each gene in the target set, and for a combination

7776 perturbation, we combine the embeddings of the  
7777 individual genes (e.g. by summation or an attention  
7778 mechanism). Using a learned embedding for each gene  
7779 gives the model a chance to encode each genes  
7780 characteristic perturbation impact. For a multi-gene  
7781 perturbation, a simple summation assumes independence  
7782 of effects, while a more sophisticated combination (  
7783 see below) can capture interactions.

- 7784 - **Combination Module:** The baseline context and perturbation  
7785 effect must be integrated. We concatenate the baseline  
7786 latent vector and the perturbation embedding vector into a  
7787 single combined latent representation. This combined vector  
7788 (of length ~640 in our example, if baseline latent is 512  
7789 and perturbation embedding 128) now contains information  
7790 about where the cell started and what perturbation was  
7791 applied. This is passed through further layers (a **fusion  
7792 network**) to compute the output. For instance, a multilayer  
7793 perceptron (MLP) with one or two hidden layers (e.g., 512  
7794 neurons, ReLU activation) can mix these features. This stage  
7795 allows for non-linear interactions between cell state and  
7796 perturbation e.g., the effect of perturbing gene X might  
7797 depend on the level of gene Y in the baseline state, which a  
7798 multiplicative interaction in the MLP can learn.
- 7799 - **Output Layer (Prediction Head):** The final layer of the  
7800 network produces a vector of length \*G\* (the number of genes)  
7801 , which is the predicted post-perturbation expression for  
7802 each gene. To ensure the model easily handles the fact that  
7803 many genes dont change, we design the output to predict a **change  
7804 (delta)** from baseline for each gene rather than an  
7805 absolute expression.
- 7806 - **Non-linear Interaction Modeling:** While a simple  
7807 concatenation of embeddings treats multi-gene perturbations  
7808 roughly as an additive combination of single effects, we can  
7809 enhance the model to capture **genegenegene interaction effects**  
7810 \*\*. One idea is to use an **attention mechanism or gating**  
7811 in the perturbation encoder: for example, if two genes A and  
7812 B are perturbed, instead of just summing their embeddings,  
7813 we pass them through an attention network that can learn a  
7814 pairwise interaction term. Another approach is to include  
7815 pairwise products of gene embeddings in the combined feature  
7816 (allowing the network to learn a unique contribution for  
7817 the pair \*A&B\* beyond A + B). Given that Norman *et al.*  
7818 tested primarily pairwise perturbations, we can explicitly  
7819 include a learned parameter or small network for each pair  
7820 of genes in the training set to capture any deviation from  
7821 additivity. However, to generalize to unseen pairs, a better  
7822 strategy is to learn a **function** for combining  
7823 embeddings (like attention) rather than a fixed lookup for  
7824 each pair.
- 7825 - **Incorporating Prior Knowledge (for Generalization):** To  
7826 improve inductive generalization to unseen genes, we can  
7827 draw inspiration from GEARS and similar methods. We could  
7828 initialize or regularize the gene perturbation embeddings  
7829 using external knowledge:
  - 7830 - Use a **gene co-expression network** (computed from the  
7831 baseline single-cell data or external data) as a graph,  
7832 and pass gene embeddings through a Graph Neural  
7833 Network (GNN) layer. This encourages genes that have  
7834 similar roles or expression patterns to have  
7835 embeddings that produce similar effects. Thus, if an  
7836 unseen gene has a similar co-expression profile to a  
7837 seen gene, the model might infer similar perturbation  
7838 outcomes.
  - 7839 - Use **pathway or GO (Gene Ontology) information** to  
7840 place genes in a relational graph (as GEARS did with a  
7841 GO-derived graph for perturbation embeddings). Two

7830  
7831 genes in the same pathway might be expected to produce  
7832 related downstream effects; by training on the known  
7833 genes, the model can generalize to a new gene by its  
7834 connections in the GO graph. Technically, this can be  
7835 done by adding a loss term that encourages the learned  
7836 embedding to correspond to the genes position in the  
7837 memory module, or by a GNN that propagates influence  
7838 from neighbors during training.

- These additions make the model more complex but aim to imbue it with **biological inductive bias**: (i) genes with similar baseline functions yield similar perturbation responses, and (ii) genes in related pathways affect overlapping sets of genes when perturbed. We will prioritize implementing a simpler version (like using co-expression PCA or clusters to initialize embeddings), and note that full graph-based learning could further improve generalization if needed.

- **Residual Connections and Regularization:** We will include skip-connections wherever helpful (for example, the baseline input could be fed not only into the encoder but also concatenated directly to a later layer, or the output head could directly see the raw baseline as well, ensuring the model can easily learn identity for unchanged genes). Regularization techniques like dropout in the MLP layers, L2 weight decay, or even an auxiliary loss to reconstruct the baseline (autoencoder style) can be employed to prevent overfitting and encourage the model to learn meaningful latent features rather than memorizing training examples.

In essence, the architecture is a **conditioned deep neural network**: it conditions on the cells initial state and the perturbation, and produces an output state. This is somewhat analogous to an encoderdecoder model where the encoder is the baseline expression and the condition or control signal is the perturbation identity. Because accuracy is paramount, we allow a fairly large model with enough capacity to capture complex gene regulatory responses.

### ### Training Strategy for Accuracy and Generalization

With the architecture in place, we next focus on **training methodology**, as it greatly affects model generalization and performance:

- **Training Data Construction:** We will pair each perturbed cells data with a baseline profile as input. Since in the actual experiment we typically do not have a *pre-perturbation* measurement of the same cell, we have to simulate a baseline. We can use the expression profiles of control cells (non-targeting sgRNA) as proxies for baseline input. For each perturbed cell in the training set, we can randomly sample a control cells expression as the baseline input. This effectively assumes that any control cell is an example of an unperturbed state the perturbed cell *could* have come from. Over many samples, the model will learn to map from an average baseline state to the perturbed outcome. We can further refine this pairing by matching on cell state: e.g., cluster the baseline cells by expression and pick a baseline from the same cluster as the perturbed cells profile (minus the perturbation effect) to provide a closer starting point. However, random pairing with a large pool of controls adds variability that can help the model not to overfit a one-to-one mapping.
- **Unseen cell context generalization:** By exposing the model to many different baseline samples paired with a given perturbation outcome (through random pairing), we train it to handle diverse baseline inputs for the same perturbation. This should improve the models

7884 robustness to any particular baseline context and  
 7885 enable generalization to new baseline profiles.  
 7886 Essentially, the model sees that the same perturbation  
 7887 can apply to various starting expression patterns.

7888

- 7889 - **Loss Functions:** The primary loss will be **Mean Squared**  
 7890 **Error (MSE)** between the predicted and actual post-  
 7891 perturbation expression vectors. To ensure we adequately  
 7892 learn the important changes, we can modify the loss in two  
 7893 ways:

- 7894 - Compute a weighted MSE that gives higher weight to genes  
 7895 that are truly differentially expressed in that  
 7896 training example. For instance, if we know gene  $j$   
 7897 changed significantly in the real perturbed cell (  
 7898 compared to baseline or compared to controls), we can  
 7899 upweight the error on gene  $j$  for that sample. This  
 9000 forces the model to focus on fitting the genes that  
 9001 move, rather than being dominated by the many near-  
 9002 zero changes.
- 9003 - Alternatively, we can train in two phases: first  
 9004 optimize MSE on the full profile to get general trends,  
 9005 then fine-tune the model on just the top-k DE genes  
 9006 for each perturbation (or using a loss like  
 9007 contrastive that emphasizes getting the direction of  
 9008 change correct).

9009 In practice, a simpler approach is to stick with standard  
 9010 MSE on all genes but monitor the top-k gene  
 9011 performance as a separate metric, ensuring the model  
 9012 doesn't ignore those signals. If we see the model  
 9013 predicting trivial (no-change) solutions, we will  
 9014 adjust the loss weighting.

- 9015 - **Optimizer and Regularization:** We will use **Adam optimizer**  
 9016 **(adaptive learning rate)** which is well-suited for  
 9017 training deep networks on possibly noisy data. A relatively  
 9018 small learning rate (e.g.  $1e-3$  to start) will be used and  
 9019 well monitor validation loss for convergence. Early stopping  
 9020 on the validation MSE/PCC can prevent overfitting. Dropout  
 9021 layers (e.g. dropout rate 0.2-0.5) can be inserted in the  
 9022 MLP to regularize. Weight decay (L2) will help keep  
 9023 embedding weights from growing excessively (important if  
 9024 some genes cause large effects, we don't want the weights to  
 9025 explode to fit one-hot signals).
- 9026 - **Generalization Enhancements:** To specifically improve  
 9027 performance on **unseen genes**:

- 9028 - We will pre-train the **perturbation encoder embeddings**  
 9029 **:** one idea is to train a simple autoencoder or  
 9030 matrix factorization on the **observed perturbation**  
 9031 **effects** in training. For example, compute the average  
 9032 expression change each perturbation causes on each  
 9033 gene (from data), yielding a perturbation-gene effect  
 9034 matrix. We can factorize this matrix to get a low-  
 9035 dimensional embedding for genes and for perturbations.  
 9036 If a gene is unseen, we can still derive a gene  
 9037 embedding from its baseline expression pattern using  
 9038 the gene embedding from baseline encoder. This is a  
 9039 bit complex, but essentially we want the model to know  
 9040 something about gene function even if not perturbed.  
 9041 Incorporating gene ontology or known regulator  
 9042 information for unseen genes could also inform an  
 9043 initial embedding (e.g., if an unseen gene is a  
 9044 transcription factor similar to another seen TF,  
 9045 assign it a similar perturbation embedding).

9046 By combining a carefully-designed architecture with these  
 9047 training strategies, our goal is a model that not only fits  
 9048 the training perturbations well but **generalizes to new**  
 9049

```

7938 perturbations and contexts**. In particular, success would
7939 be a model that, given an arbitrary K562 baseline profile
7940 and a novel target gene, can predict a biologically
7941 plausible outcome capturing which pathways might be up or
7942 down-regulated by that genes perturbation, even without
7943 having seen it before. The inclusion of baseline input
7944 enables context-specific predictions (so the same gene
7945 perturbation could have different effects in different
7946 baseline states if the model learns those interactions), and
7947 the use of gene embeddings with potential graph-based
7948 inductive bias enables inductive generalization across genes.
7949
7948 Next, we present the implementation of this model in PyTorch,
7949 structured for GPU acceleration.
7950
7951 -----
7952 ## PyTorch Implementation (Model & Training Code)
7953
7954 Below is the complete PyTorch code for the proposed model and an
7955 example training loop. This code assumes that the dataset
7956 has been preprocessed into tensors: `X_train` (baseline
7957 expressions), `Y_train` (perturbed expressions), and `
7958 P_train` (perturbation indicator vectors). The model is
7959 designed to run on a GPU for efficiency.
7960
7961 ```python
7962 import torch
7963 import torch.nn as nn
7964 import torch.optim as optim
7965
7966 # Define device for GPU acceleration
7967 device = torch.device('cuda' if torch.cuda.is_available() else '
7968 cpu')
7969 print("Using device:", device)
7970
7971 # Example shapes (for illustration; replace with actual dataset
7972 dimensions)
7973 num_genes = X_train.shape[1] # dimensionality of gene expression
7974 (e.g. ~20000)
7975 num_target_genes = P_train.shape[1] # dimensionality of
7976 perturbation vector (e.g. ~100)
7977
7978 # 1. Define the model architecture
7979 class PerturbationPredictor(nn.Module):
7980 def __init__(self, input_dim_expr, input_dim_pert,
7981 hidden_dim_expr=512, hidden_dim_pert=128,
7982 hidden_dim_comb=512):
7983 """
7984 input_dim_expr: number of genes (baseline expression
7985 vector length)
7986 input_dim_pert: number of possible targeted genes (length
7987 of perturbation indicator vector)
7988 """
7989 super(PerturbationPredictor, self).__init__()
7990 # Baseline expression encoder network
7991 self.expr_encoder = nn.Sequential(
7992 nn.Linear(input_dim_expr, hidden_dim_expr),
7993 nn.ReLU(),
7994 nn.Linear(hidden_dim_expr, hidden_dim_expr),
7995 nn.ReLU()
7996)
7997 # Perturbation encoder network
7998 self.pert_encoder = nn.Sequential(
7999 nn.Linear(input_dim_pert, hidden_dim_pert),
8000 nn.ReLU(),
8001 nn.Linear(hidden_dim_pert, hidden_dim_pert),
8002 nn.ReLU()
8003)
8004 # Fusion network to combine baseline and perturbation
8005 representations

```

```

7992 self.comb_fc1 = nn.Linear(hidden_dim_expr +
7993 hidden_dim_pert, hidden_dim_comb)
7994 self.comb_act = nn.ReLU()
7995 self.comb_fc2 = nn.Linear(hidden_dim_comb, input_dim_expr)
7996 # outputs delta for each gene
7997 def forward(self, baseline_expr, perturb_vec):
7998 # baseline_expr: Tensor of shape (batch_size,
7999 input_dim_expr)
8000 # perturb_vec: Tensor of shape (batch_size, input_dim_pert
8001)
8002 # Encode baseline and perturbation
8003 z_expr = self.expr_encoder(baseline_expr) # shape: (batch,
8004 hidden_dim_expr)
8005 z_pert = self.pert_encoder(perturb_vec) # shape: (batch,
8006 hidden_dim_pert)
8007 # Combine encoded representations
8008 z_comb = torch.cat([z_expr, z_pert], dim=1) # shape: (
8009 batch, hidden_dim_expr+hidden_dim_pert)
8010 h = self.comb_act(self.comb_fc1(z_comb)) # hidden layer
8011 delta = self.comb_fc2(h) # predicted change in expression
8012 # Add the predicted change (delta) to the baseline to get
8013 final prediction
8014 pred_expr = baseline_expr + delta
8015 return pred_expr
8016
8017 # Instantiate the model and move to GPU
8018 model = PerturbationPredictor(num_genes, num_target_genes).to(
8019 device)
8020 print(model)
8021
8022 # 2. Set up loss function and optimizer
8023 criterion = nn.MSELoss() # Mean Squared Error loss
8024 optimizer = optim.Adam(model.parameters(), lr=1e-3)
8025
8026 # 3. Training loop (example)
8027 num_epochs = 20
8028 batch_size = 256
8029
8030 # If data is large, use DataLoader for batching (assuming
8031 X_train, P_train, Y_train are TensorDatasets)
8032 train_dataset = torch.utils.data.TensorDataset(X_train, P_train,
8033 Y_train)
8034 train_loader = torch.utils.data.DataLoader(train_dataset,
8035 batch_size=batch_size, shuffle=True)
8036
8037 for epoch in range(1, num_epochs+1):
8038 model.train()
8039 total_loss = 0.0
8040 for X_batch, P_batch, Y_batch in train_loader:
8041 # Move data to device
8042 X_batch = X_batch.to(device)
8043 P_batch = P_batch.to(device)
8044 Y_batch = Y_batch.to(device)
8045 # Forward pass
8046 pred = model(X_batch, P_batch)
8047 loss = criterion(pred, Y_batch)
8048 # Backpropagation
8049 optimizer.zero_grad()
8050 loss.backward()
8051 optimizer.step()
8052 total_loss += loss.item() * X_batch.size(0)
8053 avg_loss = total_loss / len(train_dataset)
8054 if epoch % 5 == 0 or epoch == 1:
8055 print(f"Epoch {epoch:02d}: Training MSE = {avg_loss:.4f}")
8056
8057 # 4. Evaluation on test set
8058 model.eval()
8059 with torch.no_grad():
8060 X_test = X_test.to(device)
8061 P_test = P_test.to(device)

```

```

8046 Y_test = Y_test.to(device)
8047 pred_test = model(X_test, P_test)
8048 mse_full = criterion(pred_test, Y_test).item()
8049 # Compute Pearson correlation coefficient (PCC) for each test
8050 sample
8051 pred_np = pred_test.cpu().numpy()
8052 Y_np = Y_test.cpu().numpy()
8053 baseline_np = X_test.cpu().numpy()
8054 pcc_list = []
8055 topk = 20
8056 topk_pcc_list = []
8057 for i in range(Y_np.shape[0]):
8058 true_expr = Y_np[i]
8059 pred_expr = pred_np[i]
8060 # PCC for all genes
8061 cov = np.cov(true_expr, pred_expr, bias=True)
8062 # cov matrix 2x2: cov[0,1] is covariance between true and
8063 pred
8064 pcc = cov[0,1] / (np.std(true_expr) * np.std(pred_expr) + 1e
8065 -8)
8066 pcc_list.append(pcc)
8067 # PCC for top-k differentially expressed genes
8068 # Identify top-k genes by absolute change in true expression
8069 vs baseline
8070 base_expr = baseline_np[i]
8071 diff = np.abs(true_expr - base_expr)
8072 topk_idx = np.argsort(diff)[-topk:]
8073 if topk > 0:
8074 true_top = true_expr[topk_idx]
8075 pred_top = pred_expr[topk_idx]
8076 cov_top = np.cov(true_top, pred_top, bias=True)
8077 topk_pcc = cov_top[0,1] / (np.std(true_top) * np.std(
8078 pred_top) + 1e-8)
8079 topk_pcc_list.append(topk_pcc)
8080 # Calculate mean metrics
8081 mean_pcc = float(np.mean(pcc_list))
8082 mean_topk_pcc = float(np.mean(topk_pcc_list))
8083 print(f"Test MSE (all genes): {mse_full:.4f}")
8084 print(f"Test mean PCC (all genes): {mean_pcc:.3f}")
8085 print(f"Test mean PCC (top-{topk} DE genes): {mean_topk_pcc:.3f
8086 ...}")
8087 ...
8088
8089 **Explanation of the Code:**
8090
8091 - We define a `PerturbationPredictor` model class that
8092 implements the architecture described. The baseline
8093 expression encoder (`expr_encoder`) and perturbation encoder
8094 (`pert_encoder`) are simple feed-forward networks. These
8095 could be extended or replaced with more complex sub-networks
8096 (e.g., adding dropout, or using a graph convolution in `
8097 pert_encoder` if incorporating a gene network). The combined
8098 representation is fed through two linear layers (`comb_fc1`
8099 and `comb_fc2`) with a ReLU in between. The output of `
8100 comb_fc2` is a vector of length equal to number of genes,
8101 representing the predicted change in expression. We then add
8102 this to the baseline (`pred_expr = baseline + delta`) to
8103 get the final prediction.
8104
8105 - We move the model to GPU (`.to(device)`). The training loop
8106 iterates over mini-batches of data (using a DataLoader for
8107 efficiency). For each batch, we do a forward pass to compute
8108 predictions and then use `nn.MSELoss` to compute the loss
8109 against the true perturbed expression. We then backpropagate
8110 (`loss.backward()`) and update weights with Adam optimizer.
8111 We print the training MSE every few epochs for monitoring.
8112
8113 - In the evaluation section, after training, we compute metrics
8114 on `X_test, P_test, Y_test`. We calculate the overall MSE on
8115 test (as `mse_full`). Then for each test sample, we compute
8116 the Pearson correlation coefficient (PCC) between predicted
8117 and true expression. We do this manually using numpy:

```

8100  
 8101 computing the covariance and standard deviations (note: one  
 8102 could also use ``scipy.stats.pearsonr`` or `torchmetrics` for  
 8103 this). We also compute **top-20 PCC** by finding the 20  
 8104 genes with largest absolute change in the true data (  
 8105 compared to baseline) and computing PCC on that subset.  
 8106 Finally, we print the average PCC across all test cells, and  
 the average top-20 PCC.

8107 **Note:** In practice, you would likely use vectorized  
 8108 operations or library functions to compute PCC for all cells  
 8109 at once. Also, you might evaluate the model separately on  
 8110 the different generalization scenarios (e.g., one test set  
 8111 for unseen genes, one for seen genes, etc.) and compare  
 8112 performance. The code above provides a template that can be  
 8113 adapted for those purposes.

8114 By executing this model training pipeline, we expect to obtain a  
 8115 model that can accurately predict gene expression changes  
 8116 in K562 cells after CRISPRi perturbations. Its performance  
 8117 will be reflected in low MSE (especially on perturbed genes)  
 8118 and high correlation with actual profiles, even for  
 8119 perturbations or cell states not seen in training.  
 8120 Ultimately, this demonstrates a powerful approach to **in**  
 8121 **silico** predict cellular response to genetic perturbations,  
 8122 leveraging modern deep learning (PyTorch with GPU) for both  
 8123 precision and generalization. The combination of baseline  
 8124 context and perturbation identity as inputs allows the model  
 8125 to capture conditional dependencies (how a perturbations  
 8126 effect can vary by cell state), addressing one of the key  
 8127 challenges in predictive biology.

## 8125 Perplexity DeepResearch

### 8126 Research Agent Output

8127  
 8128 Norman et al. dataset(gene knockout, scRNA-seq)

8129  
 8130 Task Analysis and Technical Challenges

8131 I. Biological Context of CRISPRi Perturbations

8132 The Norman et al. dataset captures transcriptional consequences  
 8133 of CRISPRi-mediated gene repression through single-cell  
 8134 profiling of:

- 8135 - 105 individual gene knockouts
- 8136 - 131 pairwise gene interactions<sup>[1]</sup><sup>[4]</sup>
- 8137 - ~90,000 cells with baseline (unperturbed) and perturbed states<sup>[4]</sup>

8138  
 8139 Key biological considerations:

- 8140 1. **Non-linear genetic interactions**: Pairwise perturbations  
 8141 exhibit emergent effects not predictable from individual  
 8142 gene knockouts<sup>[2]</sup>
- 8143 2. **Cell-state heterogeneity**: Native transcriptional  
 8144 variability in K562 populations influences perturbation  
 responses<sup>[2]</sup><sup>[4]</sup>
- 8145 3. **Technical noise characteristics**: Dropout patterns and UMI  
 8146 count distributions differ between perturbed/control cells<sup>[4]</sup><sup>[6]</sup>

8147  
 8148 II. Machine Learning Challenges

- 8149 1. **High-dimensional output space**: Predicting 20,000+ gene  
 8150 expressions simultaneously
- 8151 2. **Sparse perturbation coverage**: Only 0.6% of possible gene  
 pairs tested in original study<sup>[1]</sup>
- 8152 3. **Context-dependent effects**: Perturbation responses vary  
 8153 with baseline cellular states<sup>[2]</sup><sup>[5]</sup>

```

8154
8155 4. Combinatorial explosion: Quadratic growth of possible
8156 gene interactions with target count
8157
8158 Evaluation Protocol Design
8159
8160 Train/Test Splits:
8161 - Unseen perturbations: Hold out 20% of single-gene and 30%
8162 of pairwise perturbations
8163 - Unseen contexts: Stratify by baseline PCA clusters, hold
8164 out entire cell-state subpopulations
8165 -----
8166 Methodological Framework: GenePerturbNet
8167
8168 I. Architectural Overview
8169
8170 Core Components:
8171
8172 1. Baseline Encoder: Variational autoencoder mapping
8173 unperturbed expression x to latent z .
8174
8175 2. Perturbation Transformer: Attention network processing
8176 target genes g into perturbation vector δ .
8177
8178 3. Interaction Module: Graph neural network modeling gene-
8179 gene interaction effects.
8180
8181 4. Response Decoder: generating perturbed expression profile
8182
8183 II. Training Strategy
8184
8185 Loss Function:
8186
8187
$$\mathcal{L} = \text{ELBO}(x, y) + \lambda_1 \text{MSE}(y, \hat{y}) + \lambda_2 \mathcal{L}_{\text{contrast}}$$

8188
8189 Where:
8190
8191 - Contrastive Loss:
8192
$$\mathcal{L}_{\text{contrast}} = -\log \frac{\exp(s(y, \hat{y}))}{\sum_{y'} \exp(s(y', \hat{y}))}$$

8193
8194 - $s()$: Cosine similarity
8195
8196 - τ : Temperature parameter
8197
8198 IV. Curriculum Learning Schedule:
8199
8200 1. Phase 1: Pretrain on single-gene perturbations
8201
8202 2. Phase 2: Fine-tune with pairwise interactions
8203
8204 3. Phase 3: Joint optimization with contrastive loss
8205
8206 Implementation and Benchmarking
8207
8208 I. Code Implementation
8209
8210 ```python
8211 import scanpy as sc
8212 import scvi
8213 import torch
8214 from torch import nn
8215 from scvi import REGISTRY_KEYS
8216 from scvi.module.base import BaseModuleClass
8217
8218 class GenePerturbNet(scvi.model.SCVI):
8219 def __init__(self, adata, n_hidden=256, n_latent=32,
8220 perturbation_embedding=64, **kwargs):
8221 super().__init__(adata, **kwargs)
8222 self.pert_embed = nn.Embedding(adata.n_vars,
8223 perturbation_embedding)
8224 self.transformer = nn.TransformerEncoderLayer(

```

```

8208 d_model=perturbation_embedding, nhead=4
8209)
8210 self.gnn = GeneInteractionGNN(perturbation_embedding)
8211
8212 self.decoder = nn.Sequential(
8213 nn.Linear(n_latent + perturbation_embedding, n_hidden),
8214 nn.ReLU(),
8215 nn.Linear(n_hidden, adata.n_vars)
8216)
8217
8218 def forward(self, x, perturbations, *args):
8219 qz_m, qz_v, z = self.encoder(x)
8220
8221 pert_emb = self.pert_embed(perturbations)
8222 delta = self.transformer(pert_emb)
8223
8224 z_pert = z + delta.mean(dim=1)
8225
8226 px_scale, px_r, px_rate, px_dropout = self.decoder(z_pert)
8227
8228 return px_rate, qz_m, qz_v, z_pert
8229
8230 class GeneInteractionGNN(nn.Module):
8231 def __init__(self, embed_dim):
8232 super().__init__()
8233 self.conv1 = nn.Conv1d(embed_dim, embed_dim*2, kernel_size
8234 =3)
8235 self.attn = nn.MultiheadAttention(embed_dim, num_heads=4)
8236
8237 def forward(self, gene_pairs):
8238 x = self.conv1(gene_pairs)
8239 x, _ = self.attn(x, x, x)
8240 return x.mean(dim=1)
8241
8242
8243
8244
8245

```

## II. Training Protocol

```

8235 ``python
8236 adata = sc.read("norman_2019.h5ad")
8237 scvi.model.setup_anndata(adata, layer="counts",
8238 categorical_covariate_keys=["perturbation"])
8239
8240 model = GenePerturbNet(adata, n_latent=32,
8241 perturbation_embedding=64)
8242 model.train(max_epochs=100, batch_size=256,
8243 plan_kwargs={"contrastive_weight":0.5})
8244
8245 control_cells = adata[adata.obs["perturbation"] == "control"]
8246 perturbed_pred = model.predict(control_cells, targets=["geneA",
8247 "geneB"])
8248
8249
8250

```

## III. Conclusion and Future Directions

8247 This framework demonstrates that deep learning architectures  
8248 explicitly modeling genetic interactions can significantly  
8249 improve perturbation prediction accuracy. Key  
8250 recommendations:

- 8251 1. **Data Expansion**: Incorporate epigenetic data to better  
8252 model CRISPRi efficiency
- 8253 2. **Dynamic Modeling**: Add temporal component to capture  
8254 perturbation kinetics
- 8255 3. **Transfer Learning**: Adapt model to other cell types using  
8256 few-shot learning

8257 The provided codebase offers a foundation for extending  
8258 perturbation prediction to novel experimental designs and  
8259 therapeutic applications.

8260  
8261

8262  
8263  
8264  
8265  
8266  
8267  
8268  
8269  
8270  
8271  
8272  
8273  
8274  
8275  
8276  
8277  
8278  
8279  
8280  
8281  
8282  
8283  
8284  
8285  
8286  
8287  
8288  
8289  
8290  
8291  
8292  
8293  
8294  
8295  
8296  
8297  
8298  
8299  
8300  
8301  
8302  
8303  
8304  
8305  
8306  
8307  
8308  
8309  
8310  
8311  
8312  
8313  
8314  
8315

## Research Agent Output

Norman et al. dataset(gene knockout, scRNA-seq)

Task Analysis

# Dataset Description

The foundation of this predictive modeling task lies in the comprehensive dataset generated by Norman et al. and published in Science in 2019. This study aimed to explore genetic interaction manifolds constructed from rich single-cell phenotypes using a Perturb-seq approach.<sup>15</sup> The experiments were conducted on human K562 leukemia cells, a commonly used model system in hematological research. The study employed CRISPRi to target a total of 236 genetic perturbations, encompassing 105 single genes and 131 pairs of genes. The scale of the dataset is substantial, comprising approximately 90,000 single-cell RNA-seq profiles, which includes both control cells and cells subjected to various single and combinatorial CRISPRi perturbations.

The dataset's focus on both single and paired gene perturbations is crucial for training a model capable of predicting combinatorial effects, which are often non-additive and complex. Single gene perturbations provide baseline information about the effect of individual gene silencing. Paired perturbations allow the model to learn about genetic interactions, where the combined effect of perturbing two genes is different from the sum of their individual effects. This is essential for capturing the non-linear nature of biological systems.

# Problem Definition

The primary objective is to develop a predictive model that takes as input the baseline gene expression profile of an unperturbed K562 cell and the identity of the target gene(s) for CRISPRi perturbation. The desired output of the model is the predicted gene expression profile of the same cell after the specified perturbation has been applied.

The model's performance will be evaluated under two key scenarios to assess its generalizability:

1. **\*\*Unseen Perturbations:\*\*** The model should be able to accurately predict the effects of CRISPRi targeting genes or gene pairs that were not included in the training data. This scenario tests the model's ability to extrapolate its learned knowledge to novel genetic manipulations.
2. **\*\*Unseen Cell Contexts:\*\*** The model should be capable of predicting the response to a perturbation in cells with baseline gene expression profiles that were not observed during the training phase. This evaluates the model's robustness to the inherent heterogeneity within the K562 cell population.

The evaluation scenarios highlight the need for a model that can generalize beyond the specific perturbations and cell states seen during training, which is a significant challenge in biological systems due to their complexity and inherent variability. A model trained only on observed perturbations might not accurately predict the effects of novel perturbations or responses in slightly different cellular contexts. Evaluating on unseen data is crucial for assessing the model's robustness and real-world applicability.

# Challenges and Considerations

8316  
8317 Developing an accurate predictive model for gene expression  
8318 profiles following CRISPRi perturbation presents several  
8319 significant challenges stemming from the inherent complexity  
8320 of biological systems and the nature of single-cell data.

8321 The gene expression data obtained from scRNA-seq experiments is  
8322 inherently high-dimensional, with measurements for thousands  
8323 of genes in each individual cell. This high dimensionality  
8324 necessitates the use of sophisticated modeling techniques  
8325 capable of handling a large number of input features.  
8326 Furthermore, the underlying biological networks that govern  
8327 gene expression are incredibly complex, involving intricate  
8328 regulatory relationships between genes, transcription  
8329 factors, and other cellular components.

8330 The relationship between gene perturbations and the resulting  
8331 changes in gene expression is often non-linear.<sup>1</sup> Genetic  
8332 interactions, such as synergistic or epistatic effects, mean  
8333 that the combined effect of perturbing multiple genes  
8334 cannot simply be predicted by summing the individual effects  
8335 of each perturbation.<sup>1</sup> The model must be able to capture  
8336 these complex, non-additive relationships to make accurate  
8337 predictions, especially for the 131 gene pairs targeted in  
8338 the Norman et al. (2019) dataset.

8339 While CRISPRi is generally considered to have fewer off-target  
8340 effects compared to CRISPRko, the possibility of unintended  
8341 gene repression at secondary genomic loci still exists. The  
8342 specificity of CRISPR-Cas9 systems relies heavily on the  
8343 design of the single guide RNA (sgRNA), and mismatches  
8344 between the sgRNA and off-target DNA sequences can lead to  
8345 unintended gene silencing. Although the Norman et al. (2019)  
8346 study likely employed well-designed sgRNAs, the potential  
8347 for some level of off-target activity should be considered.

8348 Single-cell RNA-seq data is also characterized by inherent noise  
8349 and technical variability.<sup>11</sup> Factors such as cell lysis  
8350 efficiency, reverse transcription efficiency, and sequencing  
8351 depth can introduce variability in the measured gene  
8352 expression levels. The predictive model needs to be robust  
8353 to this noise and focus on capturing the true biological  
8354 signals induced by the CRISPRi perturbations.

8355 Finally, clonal expansion, a phenomenon where certain perturbed  
8356 cells proliferate more rapidly than others, can occur during  
8357 single-cell CRISPR screens. This can lead to an  
8358 overrepresentation of certain perturbations in the  
8359 sequencing data, potentially biasing the results and  
8360 affecting the training of the predictive model.  
8361 Computational methods exist to identify and potentially  
8362 adjust for clonal cell expansion.

8363 Accurately modeling the intricate web of gene regulatory  
8364 interactions and the non-linear effects of perturbations  
8365 requires sophisticated computational approaches that can  
8366 capture these complexities beyond simple additive models.  
8367 Biological systems are highly interconnected, and the effect  
8368 of perturbing one gene can propagate through the network,  
8369 affecting the expression of many other genes in non-  
predictable ways using linear models. Therefore, models  
capable of learning complex, non-linear relationships are  
needed.

8363 -----  
8364  
8365 New Method Plan

8366 To address the challenges outlined above and develop an accurate  
8367 predictive model for gene expression profiles following  
8368 CRISPRi, a deep learning-based approach is proposed.  
8369 Specifically, a Graph Neural Network (GNN) architecture,

8370  
8371  
8372  
8373  
8374  
8375  
8376  
8377  
8378  
8379  
8380  
8381  
8382  
8383  
8384  
8385  
8386  
8387  
8388  
8389  
8390  
8391  
8392  
8393  
8394  
8395  
8396  
8397  
8398  
8399  
8400  
8401  
8402  
8403  
8404  
8405  
8406  
8407  
8408  
8409  
8410  
8411  
8412  
8413  
8414  
8415  
8416  
8417  
8418  
8419  
8420  
8421  
8422  
8423

inspired by the success of models like GEARS, appears to be a promising candidate.

**\*\*Proposed Model Architecture\*\***

The proposed model will leverage a GNN to incorporate prior knowledge about gene-gene relationships and model the perturbation effects as changes within this network. GNNs are well-suited for learning representations of nodes in a graph by aggregating information from their neighbors, allowing the model to capture the dependencies and interactions between genes.

The input to the model will consist of two components: (1) the baseline gene expression profile of an unperturbed cell, represented as a vector of gene expression counts, and (2) the identity of the target gene(s) for CRISPRi. For single gene perturbations, the target gene will be directly specified. For paired gene perturbations, the identities of both target genes will be provided.

The model architecture will comprise the following key components:

1. **\*\*Gene Embedding Layer:\*\*** Each gene in the dataset will be assigned a low-dimensional embedding vector. These embeddings will capture intrinsic properties of the genes and will be learned during the training process. Prior biological knowledge, such as gene co-expression networks or functional annotations from databases like Gene Ontology (GO) 36, can be used to initialize these embeddings or to inform the GNN architecture.
2. **\*\*Perturbation Embedding Layer:\*\*** The identity of the perturbed gene(s) will also be encoded into an embedding vector. For single perturbations, a dedicated embedding will be learned for each targeted gene. For paired perturbations, the embeddings of the two target genes can be combined (e.g., through summation or concatenation) to represent the combined perturbation.
3. **\*\*Graph Neural Network (GNN):\*\*** A gene regulatory network (GRN) will be constructed, where genes are represented as nodes and edges represent regulatory relationships between them. This GRN can be derived from publicly available databases or inferred from the unperturbed single-cell expression data. The gene embeddings will serve as initial node features in this graph. The GNN will then propagate information across the network, allowing each gene's representation to be informed by its neighbors and their interactions. The perturbation embedding will be incorporated into the GNN, potentially by modifying the node features of the perturbed gene(s) or by influencing the message passing process.
4. **\*\*Cell State Encoding Layer:\*\*** The baseline gene expression profile of the unperturbed cell will be passed through a separate neural network (e.g., a multi-layer perceptron) to learn a low-dimensional representation of the cell's initial transcriptional state. This encoding will capture the cell's context and will be used to condition the prediction of the perturbed state.
5. **\*\*Prediction Layer:\*\*** The output of the GNN (representing the perturbed gene embeddings) and the cell state encoding will be combined (e.g., through concatenation followed by another neural network) to predict the gene expression profile after the perturbation. The output will be a vector of the same dimensionality as the input gene expression profile, representing the predicted expression levels for each gene in the cell.

8424  
8425 The rationale behind choosing this architecture is that it  
8426 allows for the integration of prior biological knowledge  
8427 about gene-gene interactions through the GRN. This can help  
8428 the model to better understand the potential downstream  
8429 effects of a perturbation. Furthermore, the use of  
8430 embeddings allows the model to learn meaningful  
8431 representations of genes and perturbations, potentially  
8432 enabling better generalization to unseen perturbations.

8433 **\*\*Feature Engineering and Data Preprocessing\*\***

8434 The Norman et al. (2019) dataset will require careful  
8435 preprocessing before being used to train the model. The  
8436 steps involved will include:

- 8437 1. **\*\*Data Loading and Normalization:\*\*** The processed gene  
8438 expression matrices will be loaded using appropriate  
8439 libraries like Scanpy or AnnData. The gene expression counts  
8440 will be normalized to account for differences in sequencing  
8441 depth between cells. Log transformation (e.g., using a  
8442 natural logarithm after adding a pseudocount) will be  
8443 applied to stabilize the variance of gene expression levels.
- 8444 2. **\*\*Perturbation Information Encoding:\*\*** The perturbation  
8445 information, specifying the targeted gene(s) for each cell,  
8446 will be extracted from the dataset's metadata. For single  
8447 gene perturbations, the gene name will be used. For paired  
8448 gene perturbations, both gene names will be used. These gene  
8449 names will then be mapped to their corresponding indices or  
8450 identifiers in the gene expression matrix. The perturbation  
8451 information will be encoded as input to the model,  
8452 potentially using one-hot encoding initially, where a binary  
8453 vector indicates which genes are targeted. Alternatively,  
8454 learned embeddings for each gene could be used to represent  
8455 the perturbation.
- 8456 3. **\*\*Control Sample Handling:\*\*** Cells labeled as control (  
8457 unperturbed) will be identified and used to establish the  
8458 baseline gene expression profiles. These control profiles  
8459 will be crucial for training the model to predict the  
8460 changes in expression induced by the perturbations.
- 8461 4. **\*\*Feature Selection:\*\*** Given the high dimensionality of the  
8462 gene expression data, feature selection techniques may be  
8463 employed to focus on the most relevant genes. One common  
8464 approach is to identify highly variable genes (HVGs) across  
8465 the cell population and use only these genes as input to the  
8466 model.<sup>16</sup> This can reduce the dimensionality of the input,  
8467 potentially improving model training and performance.
- 8468 5. **\*\*GRN Construction (if applicable):\*\*** If a GNN is used, a  
8469 gene regulatory network will need to be constructed. This  
8470 could involve using publicly available databases of known  
8471 gene interactions or inferring a network from the  
8472 unperturbed single-cell expression data using methods like  
8473 co-expression analysis or network inference algorithms.

8474 **\*\*Training Strategy\*\***

8475 The training of the predictive model will involve the following  
8476 steps:

- 8477 1. **\*\*Data Splitting:\*\*** The dataset will be split into training,  
8478 validation, and test sets. To address the evaluation  
8479 scenarios of unseen perturbations, the split will be  
8480 performed at the level of perturbations. This means that all  
8481 cells corresponding to certain perturbations (both single  
8482 and paired) will be held out in the test set and will not be  
8483 seen by the model during training. A separate validation  
8484 set, also containing held-out perturbations, will be used  
8485 for hyperparameter tuning and model selection.<sup>16</sup> Carefully  
8486 designing the data splitting strategy is crucial to ensure

8478 that the model is truly evaluated on unseen perturbations  
8479 and cell contexts, avoiding information leakage from the  
8480 training set.

- 8481 2. **\*\*Loss Function:\*\*** The model will be trained to minimize the  
8482 difference between the predicted gene expression profiles  
8483 and the observed gene expression profiles. The Mean Squared  
8484 Error (MSE) will be used as the primary loss function, as it  
8485 directly measures the average squared difference between  
8486 the predicted and observed values.
- 8487 3. **\*\*Optimizer:\*\*** An appropriate optimization algorithm, such as  
8488 Adam, will be used to update the model's parameters during  
8489 training. A learning rate schedule, which gradually reduces  
8490 the learning rate over time, may be employed to improve  
8491 convergence and prevent overfitting.
- 8492 4. **\*\*Hyperparameter Tuning:\*\*** The model architecture and  
8493 training process will have several hyperparameters (e.g.,  
8494 the dimensionality of the embeddings, the number of layers  
8495 in the neural networks, the learning rate). These  
8496 hyperparameters will be tuned using the validation set.  
8497 Techniques like grid search or random search can be used to  
8498 explore different combinations of hyperparameters and select  
8499 the configuration that yields the best performance on the  
8500 validation set.
- 8501 5. **\*\*Overfitting Prevention:\*\*** Techniques such as dropout,  
8502 weight regularization (e.g., L1 or L2 regularization), and  
8503 early stopping (monitoring the performance on the validation  
8504 set and stopping training when it starts to degrade) will  
8505 be used to prevent the model from overfitting to the  
8506 training data and improve its generalization ability.

8507 **\*\*Strategies for Handling Unseen Perturbations and Cell Contexts**  
8508 **\*\***

8509 The proposed GNN-based architecture offers several advantages  
8510 for handling the challenges of unseen perturbations and cell  
8511 contexts:

- 8512 - **\*\*Generalization to Unseen Genes (for GNN):\*\*** If a GNN is used  
8513 and the GRN includes genes that are not targeted in the  
8514 training set, the model may still be able to make informed  
8515 predictions about the effects of perturbing these unseen  
8516 genes by leveraging their relationships with other genes in  
8517 the network.<sup>36</sup> The model can learn general principles of how  
8518 perturbations propagate through the network, allowing it to  
8519 extrapolate to new nodes (genes). Predicting the effects of  
8520 completely novel perturbations (genes not seen during  
8521 training) is a significant challenge.
- 8522 - **\*\*Cell State Conditioning:\*\*** By explicitly encoding the  
8523 baseline gene expression profile of the unperturbed cell,  
8524 the model can condition its prediction on the specific  
8525 context of that cell. This allows the model to capture some  
8526 of the inherent heterogeneity within the cell population and  
8527 potentially make more accurate predictions for cells with  
8528 unseen baseline profiles.
- 8529 - **\*\*Learned Embeddings:\*\*** The use of learned embeddings for  
8530 genes and perturbations can help the model to capture  
8531 semantic relationships between different genes and  
8532 perturbations. If the embedding space is learned effectively,  
8533 the model may be able to generalize to unseen perturbations  
8534 that are functionally similar to those seen during training,  
8535 even if the specific genes were not encountered before.

8536 While more advanced techniques like meta-learning or domain  
8537 adaptation could potentially further improve the model's  
8538 ability to handle unseen perturbations and cell contexts,  
8539 the proposed GNN architecture with cell state conditioning  
8540

```

8532
8533 and learned embeddings provides a strong foundation for
8534 addressing these challenges.
8535
8536 -----
8537 Generate Prediction Model Code
8538
8539 The prediction model will be implemented using the Python
8540 programming language and several key libraries commonly used
8541 in machine learning and single-cell data analysis.
8542
8543 **Implementation Details**
8544 - **Programming Language:** Python
8545 - **Key Libraries:**
8546 - **PyTorch** or **TensorFlow/Keras:** For implementing the
8547 neural network architecture, including the GNN and other
8548 layers.
8549 - **Scanpy** or **AnnData:** For efficient handling and
8550 preprocessing of the single-cell RNA-seq data.
8551 - **NumPy:** For numerical computations and array
8552 manipulations.
8553 - **SciPy:** For scientific computing, including statistical
8554 functions.
8555 - **scikit-learn:** For evaluation metrics (MSE, PCC) and
8556 potential utility functions.
8557
8558 **Code Structure and Key Functions**
8559 The codebase will be organized into several modules or classes
8560 to ensure modularity and maintainability:
8561
8562 1. **`DataLoader` Class:**
8563 - **Purpose:** Responsible for loading the Norman et al.
8564 (2019) dataset from the specified files (e.g., MTX, TSV).
8565 - **Functionality:**
8566 - Reads the gene expression matrix and metadata.
8567 - Performs normalization and log transformation.
8568 - Encodes perturbation information.
8569 - Splits the data into training, validation, and test sets
8570 based on perturbations.
8571 - Provides data loaders for efficient batching during
8572 training.
8573
8574 2. **`Model` Class:**
8575 - **Purpose:** Defines the neural network architecture,
8576 including the gene embedding layer, perturbation
8577 embedding layer, GNN (if chosen), cell state encoding
8578 layer, and prediction layer.
8579 - **Functionality:**
8580 - Initializes the model parameters.
8581 - Implements the forward pass of the network, taking
8582 baseline expression and perturbation information as
8583 input and outputting the predicted perturbed
8584 expression.
8585
8586 3. **`Trainer` Class:**
8587 - **Purpose:** Handles the training process of the model.
8588 - **Functionality:**
8589 - Takes the model, data loaders, loss function, and
8590 optimizer as input.
8591 - Implements the training loop, including forward pass,
8592 loss calculation, backpropagation, and parameter
8593 updates.
8594 - Monitors performance on the validation set and
8595 implements early stopping if needed.
8596 - Saves the trained model parameters.
8597
8598 4. **`Evaluator` Class:**
8599 - **Purpose:** Evaluates the performance of the trained model
8600 on the test set.
8601 - **Functionality:**
8602 - Loads the trained model parameters.
8603 - Iterates through the test data.
8604 - Makes predictions for each test sample.
8605 - Calculates the evaluation metrics (MSE, PCC, Top-k PCC)
8606 by comparing predictions with the observed data.

```

```

8586
8587 - Reports the evaluation results.
8588 5. Predictor` Class:
8589 - Purpose: Allows for making predictions on new, unseen
8590 data (e.g., a baseline expression profile and a
8591 perturbation).
8592 - Functionality:
8593 - Loads the trained model parameters.
8594 - Takes a baseline expression profile and perturbation
8595 information as input.
8596 - Preprocesses the input data as needed.
8597 - Uses the trained model to generate the predicted gene
8598 expression profile after the perturbation.
8599
8600 Instructions for Usage
8601
8602 1. Install Libraries: Install the necessary Python libraries
8603 using pip:
8604 bash
8605 pip install torch torchvision torchaudio scapy numpy scipy
8606 scikit-learn
8607 (or tensorflow and keras if using TensorFlow).
8608
8609 2. Download Data: Download the processed data files for
8610 GSE133344 from GEO or Figshare and place them in a
8611 designated directory.
8612
8613 3. Prepare Input Data: Ensure that the input data for
8614 prediction (baseline gene expression profile and
8615 perturbation identity) is in the expected format (e.g., a
8616 NumPy array or a Pandas DataFrame).
8617
8618 4. Run the Code: Execute the Python script containing the
8619 model implementation, training, and evaluation steps. The
8620 script should:
8621 - Initialize the DataLoader with the path to the data files.
8622 - Split the data into training, validation, and test sets.
8623 - Initialize the Model with the appropriate architecture
8624 and hyperparameters.
8625 - Initialize the Trainer with the model, data loaders, loss
8626 function, and optimizer.
8627 - Train the model using the Trainer.
8628 - Initialize the Evaluator with the trained model and test
8629 data loader.
8630 - Evaluate the model and print the results.
8631 - (Optionally) Use the Predictor to make predictions on new
8632 data.
8633
8634 Evaluation and Results
8635
8636 A comprehensive evaluation of the developed predictive model is
8637 crucial to assess its performance and generalizability.
8638
8639 Experimental Setup
8640
8641 The Norman et al. (2019) dataset will be split into training,
8642 validation, and test sets, ensuring that perturbations are
8643 held out for both validation and testing to evaluate the
8644 model's ability to predict the effects of unseen genetic
8645 manipulations. The split will likely involve a significant
8646 portion of the data for training, with smaller portions
8647 allocated for validation (for hyperparameter tuning and
8648 model selection) and testing (for final performance
8649 evaluation).16 The specific hyperparameters of the chosen
8650 model architecture (e.g., number of layers, dimensionality
8651 of embeddings, learning rate) will be determined through
8652 experimentation and tuning on the validation set.
8653
8654 To make a prediction for a given perturbation, the model will
8655 take as input the baseline gene expression profile of an
8656

```

8640  
8641  
8642  
8643  
8644  
8645  
8646  
8647  
8648  
8649  
8650  
8651  
8652  
8653  
8654  
8655  
8656  
8657  
8658  
8659  
8660  
8661  
8662  
8663  
8664  
8665  
8666  
8667  
8668  
8669  
8670  
8671  
8672  
8673  
8674  
8675  
8676  
8677  
8678  
8679  
8680  
8681  
8682  
8683  
8684  
8685  
8686  
8687  
8688  
8689  
8690  
8691  
8692  
8693

unperturbed cell (which could be an average profile of control cells or a specific control cell's profile) and the identity of the target gene(s).

**\*\*Performance Metrics\*\***

The model's performance on the test set will be quantified using the three evaluation metrics defined earlier: Mean Squared Error (MSE), Pearson Correlation Coefficient (PCC), and Top-k PCC. These metrics will be calculated by comparing the model's predicted gene expression profiles with the actual observed profiles in the test set for the held-out perturbations. The results will be reported separately for unseen single-gene perturbations and unseen paired-gene perturbations to assess the model's ability to handle both types of genetic manipulations. It may also be informative to report the performance on different subsets of genes, such as the highly variable genes, as these are often the most biologically relevant. Visualizations, such as scatter plots of predicted vs. observed gene expression for representative perturbations, can provide further insights into the model's predictive capabilities.