# scAgents: A Multi-Agent Framework for Fully Autonomous End-to-End Single-Cell Perturbation Analysis

#### Anonymous Authors<sup>1</sup>

#### Abstract

While AI agents for scientific discovery have proliferated, truly end-to-end autonomous solutions 015 remain scarce, particularly in complex interdisciplinary fields like single-cell genomics. We 018 introduce scAgents, a fully autonomous multi-019 agent framework that transforms raw single-cell 020 data and task descriptions directly into optimized computational solutions. Given only a dataset and research objective as input, scAgents outputs both a novel model architecture and executable code for training and inference without human 025 intervention. When evaluated on the scPerturb datasets and benchmarks, scAgents consistently outperforms task-specific state-of-the-art methods, achieving up to 49% reduction in predic-028 029 tion error compared to scGPT for gene knock-030 outs and Pearson correlation increases of up to 20% in expression predictions versus ChemCPA for drug perturbations. scAgents' ability to succeed where existing foundation models struggle is particularly significant, adapting effectively to 034 035 different data types (scRNA-seq, scATAC-seq, CITE-seq) and various perturbation categories with consistent performance across modalities. Our code and some scAgents-designed models are 039 available at https://anonymous.4open. science/r/scAgents-2025-242E/.

### 1. Introduction

041

043

044

045

046

047

049

050

051

052

053

054

000

002 003

004 005

010

Scientific discovery is undergoing a transformation with the rise of *AI Scientists* autonomous systems designed to conduct research with minimal human intervention. Recent progress in large language models (LLMs) and AI agents has en-

abled impressive capabilities in literature analysis (Hsu et al., 2024), hypothesis generation (Qi et al., 2023; Yang et al., 2023), and experimental implementation (Tian et al., 2024; Huang et al., 2024b), as demonstrated by emerging benchmarks and agentic toolkits (Skarlinski et al., 2024; Chen et al., 2024; Majumder et al., 2024; Huang et al., 2024a). Yet these systems remain siloedcompetent at individual tasks, but incapable of orchestrating complete scientific workflows. Existing AI scientists lack the ability to propose *methods*, refine them through collaborative reasoning, and carry out empirical validation in a unified framework (Lu et al., 2024; Li et al., 2024b). This gap becomes especially critical in domains that demand interdisciplinary reasoning, such as computational biology, where effective modeling requires integrating biological priors, statistical rigor, and machine learning design (Boiko et al., 2023; Roohani et al., 2024b).

A particularly compelling instantiation of this challenge is single-cell perturbation analysis, a frontier task in computational biology that demands predicting how cells respond to diverse interventionssuch as gene knockouts, drug treatments, or cytokine stimulationby modeling state transitions in high-dimensional expression space (Figure 1). This problem exemplifies the limitations of current AI scientists: Despite progress in foundation models (Cui et al., 2024; Theodoris et al., 2023), existing approaches remain heavily fragmented, relying on human experts to bridge gaps in dataset interpretation, model design, and empirical validation (Lotfollahi et al., 2023; Hao et al., 2024). Moreover, these models often fail to generalize to novel biological scenarios, such as unseen cell types or experimental modalities (Levine et al., 2024; Wenteler et al., 2024). The challenge is further exacerbated by the heterogeneity of data types (e.g., scRNA-seq, scATAC-seq, CITE-seq) and perturbation mechanisms, each requiring domain-specific reasoning to extract meaningful insights (Peidli et al., 2024; Bendidi et al., 2024).

We present scAgents, a fully autonomous multi-agent framework that conducts end-to-end single-cell perturbation analysis by integrating three core modules: (1) **Task Analysis**,

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

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

which characterizes the dataset and retrieves relevant prior work via a hybrid breadth-first and depth-first agentic retrieval strategy; (2) Method Design, where role-specific 058 agents collaboratively analyze baselines, identify design 059 limitations, and iteratively refine novel modeling strategies 060 through graph-based discussions; and (3) Experiment Exe-061 cution, which translates high-level architectural proposals 062 into executable code, complete with training routines and 063 validation protocols. At the heart of scAgents is a persis-064 tent knowledge graph that accumulates insights across all 065 phases and ensures a coherent information flow. The agentic 066 retrieval system integrates perturbation-specific biological 067 knowledge, such as regulatory pathway relationships and 068 modality-aware considerations, into model design without 069 requiring human input. The multi-agent discussion module 070 leverages confidence-based scoring and self-critique mechanisms to converge on scientifically grounded and technically feasible solutions. The final system output includes both the model architecture and the complete implementation code 074 for training and inference.

075 We evaluate scAgents on six diverse perturbation datasets 076 from the scPerturb on single cell perturbation benchmark, 077 encompassing gene knockouts, drug treatments, and cy-078 tokine stimulations across modalities such as scRNA-seq, 079 scATAC-seq, and CITE-seq. In all cases, models designed by scAgents significantly outperform prior baselines. For 081 instance, in the drug perturbation task, scAgents improves 082 Pearson correlation by 20% over the next best method, 083 ChemCPA (Hetzel et al., 2022). On the challenging, sparse, and high-dimensional scATAC-seq dataset, it achieves a 085  $\sim$ **16-fold increase** compared to the second-best baseline, a 086 linear regression model. In surface protein prediction (CITE-087 seq), scAgents boosts Pearson correlation by  $\sim$ 177% over 088 the next best baseline method, Random Forest. These results 089 highlight scAgents' strong generalization and substantial 090 performance gains across diverse settings, including the 091 novel application scenarios. 092

### 2. Related Work

095 Agent Systems for Scientific Discovery Researchers 096 have developed specialized AI systems spanning the en-097 tire research workflow: from literature analysis tools like 098 PaperQA2 (Skarlinski et al., 2024) and CHIME (Hsu et al., 099 2024), to hypothesis generation frameworks that range from 100 domain-specific idea creation (Baek et al., 2024; Qi et al., 2023; Yang et al., 2023) to comparative evaluations with expert proposals (Si et al., 2024). These systems increasingly leverage multi-agent architectures (Ghafarollahi & Buehler, 104 2025; Schmidgall & Moor, 2025; Guo et al., 2024) to fa-105 cilitate collaborative scientific reasoning. Implementation 106 capabilities have advanced through scientific coding frame-107 works like SciCode (Tian et al., 2024) and MLAgentBench



093

094



Figure 1: Our framework predicts cellular responses to unseen interventions by modeling state transitions in high-dimensional expression space. Given multi-modal single-cell data (scRNA-seq, scATAC-seq, CITE-seq) and a training set of known perturbations (gene knockouts, drugs, cytokines), scAgents learns to map control states to perturbed states, enabling accurate prediction of gene expression profiles under novel experimental conditions without additional wet-lab experimentation.

(Huang et al., 2024a), while benchmarks evaluate these capabilities across diverse domains (Qiu et al., 2025; Ruan et al., 2024; Chen et al., 2025; Jing et al., 2024). The integration of literature analysis with data-driven approaches has proven particularly effective for hypothesis generation (Liu et al., 2024b; Zhong et al., 2023; Majumder et al., 2024), with several frameworks enhancing research ideation through structured feedback mechanisms (Pu et al., 2024; Garikaparthi et al., 2025) and approaches to improve novelty and diversity (Hu et al., 2024; Radensky et al., 2024; Gao et al., 2025). End-to-end systems now attempt to unify these capabilities, including domain-general approaches like AI Scientist (Lu et al., 2024) and MLR-Copilot (Li et al., 2024b), alongside domain-specific implementations for chemistry (Boiko et al., 2023), genomics (Roohani et al., 2024b), materials science (Ghafarollahi & Buehler, 2024), and medicine (Tang et al., 2024; Naumov et al., 2025). Despite these advances, significant challenges remain in developing truly autonomous scientific systems, particularly regarding experimental rigor (Kon et al., 2025), falsification mechanisms (Liu et al., 2024c), and comprehensive evaluation metrics (Beel et al., 2025; Friel et al., 2025), as highlighted in recent surveys (Reddy & Shojaee, 2025; Eger et al., 2025; Kulkarni et al., 2025; Ren et al., 2025).

**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 (Dixit et al., 2016) or random forest feature selection (Skinnider et al., 2021), treated each gene or cell type in isolation. Deep generative models (Lotfollahi et al., 2019; Hetzel

et al., 2022; Lotfollahi et al., 2023), conceptualize pertur-111 bations as latent space transformations through linear shifts 112 or decompositions that separate biological covariates. In 113 contrast, network-based methods (Qiu et al., 2022; Roohani 114 et al., 2024a; Bai et al., 2024; Kamimoto et al., 2023) explic-115 itly incorporate biological knowledge via gene regulatory 116 networks or cellular relationships. To further address the is-117 sue of cell heterogeneity, distribution alignment approaches 118 such as optimal transport (Bunne et al., 2023; Dong et al., 119 2023) have been applied to machine learning models (Jiang 120 et al., 2024), matching the distribution of control cells with 121 perturbed cells. The emergence of transformer architectures 122 represents the latest paradigm shift. These architectures (Hao et al., 2024; Theodoris et al., 2023; Cui et al., 2024; 124 Levine et al., 2024) leverage pre-training at scale and self-125 attention mechanisms to model complex gene dependencies 126 without explicit biological structure. This theoretical diversity creates a vast design space where selecting optimal 128 architectures, representation strategies, and biological con-129 straints remains highly context-dependent. 130

131 **Notations.** Let  $X \in \mathbb{R}^{n \times d}$  denote the matrix of single-132 cell profiles, where n represents the number of cells and 133 d represents the dimensionality of the measured features 134 (including gene expression counts, chromatin accessibility 135 or surface protein markers, depending on the modality). A 136 dataset  $D = \{(x_i, p_i, y_i)\}_{i=1}^N$  and a task description S are given, where  $x_i \in \mathbb{R}^d$  represents the pre-perturbation profile, 137 138  $p_i \in \mathcal{P}$  denotes the applied perturbation, and  $y_i \in \mathbb{R}^{d'}$ 139 corresponds to the observed post-perturbation profile. To 140 verify the generalizability of our scAgents, we divide D141 into  $D_{\text{train}} = \{(x_i, p_i, y_i)\}_{i=1}^M$  and  $D_{\text{test}} = \{(x_i, p_i, y_i)\}_{i=1}^K$ , 142 where  $p_i \in \mathcal{P}_{\text{test}}$  constitutes held-out perturbations and  $x_i \in$ 143  $\mathcal{X}_{\text{test}}$  represents held-out cell profiles. 144

Problem Formulation. We formalize the challenge of predicting cellular responses to perturbations as learning a 147 mapping function between pre-perturbation states and their 148 corresponding post-perturbation outcomes. Specifically, for 149 each perturbation  $p \in \mathcal{P}$  (e.g., gene knockout, drug treat-150 ment, cytokine stimulus) applied to a subset of cells, we 151 model the induced change in the cellular profile as a function 152  $f_p: \mathbb{R}^d \to \mathbb{R}^{d'}$ , where  $f_p(x)$  predicts the post-perturbation 153 profile of a cell x, potentially in a different modality. 154

145

155 The training objective involves learning a function  $f_{\theta}$ : 156  $\mathbb{R}^d \times \mathcal{P} \to \mathbb{R}^{d'}$  that generalizes effectively to unseen pertur-157 bations and cell states, where  $\phi$  is a trainable parameter(i.e., 158 weights and biases). To capture intrinsic cell-state structure 159 and enable efficient modeling, we incorporate learnable en-160 coders  $g_{\phi} : \mathbb{R}^d \to \mathbb{R}^h$ , where  $\phi$  is a trainable parameter and  $z_i = g_{\phi}(x_i) \in \mathbb{R}^{n \times h}$  represents the latent embedding 161 162 that preserves geometric relationships between control and 163 perturbed states, thus facilitating accurate prediction of post-164



Figure 2: **The scAgents architecture and workflow.** The framework operates through three sequential phases: Task Analysis (dataset characterization and literature retrieval), Hypothesis Generation (collaborative development of novel approaches by role-based, self-refinement agent systems), and Experiment Execution (code generation, training implementation, and results analysis). All phases communicate through a shared knowledge graph that evolves throughout the workflow.



Figure 3: **The scAgents protocol overview.** The protocol framework integrates JSON-RPC with a persistent knowledge graph, combining the strengths of A2A and MCP protocols while adding scientific domain knowledge representation.

perturbation profiles.

**Evaluation.** For evaluation, we assess  $f_{\theta}(x_i, p_i)$  for all  $(x_i, p_i) \in \mathcal{X}_{\text{test}} \times \mathcal{P}_{\text{test}}$  and evaluate the quality of the learned representation  $g_{\phi}(x)$  in terms of its ability to reconstruct gene expression profiles. Our evaluation protocols incorporate metrics such as mean squared error, Pearson's correlation coefficient, and perturbation consistency adapted from (Roohani et al., 2024a; Bendidi et al., 2024) to ensure biological significance. Appendix B provides a detailed explanation of these metrics.

**Evaluation.** For evaluation, we assess  $f_{\theta}(x_i, p_i)$  for all  $(x_i, p_i) \in \mathcal{X}_{\text{test}} \times \mathcal{P}_{\text{test}}$  and evaluate the quality of the learned representation  $g_{\phi}(x)$  in terms of its ability to reconstruct gene expression profiles. Our evaluation protocols incorporate metrics such as mean squared error, Pearson's correlation coefficient, and perturbation consistency adapted from (Roohani et al., 2024a; Bendidi et al., 2024) to ensure biological significance. Appendix B provides a detailed explanation of these metrics.

scAgents orchestrates an end-to-end scientific workflow into a coherent framework (Figure 2): it autonomously (1) analyzes dataset characteristics and identifies optimal modeling strategies, (2) designs neural architectures incorporating bi165 ological knowledge about gene regulatory networks, and (3) generates executable code implementing the complete 167 prediction pipeline, and evaluates prediction accuracy using 168 various metrics. Central to scAgents is a hybrid communi-169 cation protocol combining JSON-RPC data exchange with 170 a persistent knowledge graph, which is further described in 171 Appendix A.3. This shared representation continuously in-172 tegrates outputs from individual agents, enabling complex 173 reasoning chains while maintaining traceability throughout 174 the scientific process. Detailed information of agent com-175 munication protocol are provided in Appendix E. The exact 176 configurations, prompts, and outputs of each agent are listed 177 in Appendix A.2 and F.

## 179 2.1. Task Analysis Module

178

193

203

204

206

208

209

210

211

212

213

214

215

216

217

218

219

The Task Analysis phase begins with comprehensive char acterization of the scientific problem through three components:

184 Data Parser. This component extracts key metadata from 185 single-cell datasets, including perturbation types, gene features, and cell populations. It standardizes information 186 187 across diverse modalities (RNA-seq, ATAC-seq, CITE-seq) 188 and generates summary statistics to establish the data foundation. The parser identifies critical experimental param-189 190 eters such as perturbation methods, organism details, and 191 cellular contexts without human intervention (detailed examples in Appendix G.1).

Agentic Retrieval. Our retrieval system combines a static corpus of 45 specialized articles with dynamic search capa-195 bilities through PubMed and GitHub APIs. Starting with an 196 initial query  $Q^{(0)}$  derived from the task description, the sys-197 tem employs Sentence-BERT (Reimers & Gurevych, 2019) 198 embeddings and performs multi-layer retrieval alternating 199 between breadth-first and depth-first search strategies. The 200 detailed algorithm and mathematical formulation are pre-201 sented in the Appendix D.1. 202

- BFS layer (t odd): Retrieves top-K documents  $\mathcal{N}_t = \text{TopK}(Q^{(t)}, \text{mode} = \text{BFS}).$
- *DFS layer* (*t* even): Follows highest-scoring paths from  $\mathcal{N}_{t-1}$  in depth.

Document relevance is computed via cosine similarity:  $Score(Q, d) = \frac{e(Q) \cdot e(d)}{\|e(Q)\|\|e(d)\|}$ . The retrieval terminates upon reaching any of three conditions: (1) maximum layers  $L_{max} = 10$ , (2) query overlap exceeding threshold  $\tau = 0.8$ , or (3) document scores below  $\epsilon = 0.5$ . Results are stored in a vector database for subsequent access.

Agent Collaboration. Three specialized agents *Dataset Analyst*, *Problem Investigator*, *Baseline Assessor* process the retrieved information. The Dataset Analyst examines data



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, the agentic system will judge if current state satisfy the stopping criteria. If not, each expert will refine their ideas based on the self-critic's suggestions and other experts' viewpoint. This graph-based self-critic refinement continues until reaching the terminate state. Complete multirounds of discussions are in Appendix D.2.

integrity and characteristics, the Problem Investigator defines research questions and analytical approaches, and the Baseline Assessor establishes reference models and benchmarks. Their outputs are synthesized by a *Refinement Agent* that structures the analysis into a JSON file.

Analysis Report At the core of our framework's performance gains is its ability to autonomously discover optimal model architectures for each specific perturbation task. Similar to how human experts approach new domains, the agents systematically analyze scientific literature to identify candidate architectures appropriate for the given data characteristics. The Baseline Assessor specifically evaluates different model architectures (e.g. Transformers, GANs, VAEs, MLPs, etc.) based on their strengths and limitations for the particular perturbation type and data modality. For each architecture component, the agent generates a comprehensive analysis report detailing advantages, disadvantages, and suitability for specific aspects of the task (e.g., handling sparsity in scATAC-seq data). This literature-grounded comparative analysis forms the foundation for the Method Design module, where expert agents collaboratively refine these insights into optimized custom architectures.

#### 2.2. Method Design Module

**Overview.** This module produces a research plan through collaborative agent discussion, concurrently producing three integrated components: (1) data preprocessing strategies, (2) model architecture design, and (3) concrete model implementation details. Architecture design involves textual descriptions of the neural network components selected for the specific perturbation task and their theoretical justifications, while model design translates these concepts into executable pseudocode. Unlike approaches that merely tune hyperparameters of fixed architectures, scAgents fun-

damentally focuses on discovering optimal architectural
combinations tailored to each dataset's unique biological
characteristics. This architectural discovery process, rather
than hyperparameter optimization, is the primary source of
our framework's performance advantages in perturbation
prediction tasks.

227

228

229

230

231

232

233

234

235

236

**Multi-Expert Critic System.** For each task, the system dynamically selects a subset of domain experts  $E^{(k)}$  (e.g., Data Engineer, Single-Cell Biologist, Deep Learning Expert) based on task requirements, along with a permanent Self-Critic. 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.

237 **Graph-based Discussion.** The system performs T rounds 238 of graph neural message passing, where experts propose 239 architectural solutions, the Self-Critic identifies weaknesses, 240 and all nodes update their perspectives based on collective 241 feedback. After each round, confidence scores are updated 242 according to the evaluation quality by both the Self-Critic 243 and peer experts. The discussion ends when all experts' con-244 fidence scores exceed the threshold  $\tau = 0.8$  with minimal variance  $(\|c_{t^*}^{(i)} - c_{t^*}^{(i)}\| < \epsilon)$ . If not, experts refine their pro-245 posals using historical context and proceed to the next round. 246 247 This process ensures convergence toward scientifically valid 248 and technically feasible model designs with explicit reason-249 ing chains. Further information of Experts selection and 250 discussion construction is in Appendix A.4, detailed algo-251 rithm and mathematical formulation are presented in the 252 Appendix D.2. 253

## 254 **2.3. Experiment Execution Module**

The Experiment Execution module transforms conceptual
designs into operational implementations:

(1) Code Generation. The Code Generator translates method 258 designs into production-ready implementations with com-259 plete dependency management. It produces modular code with built-in error detection mechanisms that automate the 261 correction of common implementation issues. (2) Training Process. Model training proceeds through an automated 263 scheduler implementing best-practice optimization routines 264 and early stopping mechanisms. Hyperparameter optimiza-265 tion runs in parallel with model checkpointing to ensure 266 reproducibility and enable efficient experimentation. (3) 267 Validation & Refinement. Rather than one-shot generation, the Validation Agent implements an iterative refinement 269 cycle. It evaluates performance against established met-270 rics, identifies failure modes, and progressively enhances 271 implementation quality. This evolutionary approach allows 272 the system to address complex edge cases and optimize 273 implementation based on empirical feedback. 274

#### **3.** Experiments

#### 3.1. Evaluation setup

We evaluate the models designed and implemented by scAgents in various types of perturbation from scPerturb (Peidli et al., 2024), including gene knockouts, drug treatments, and cytokine stimulation across multiple modalities (scRNA-seq, scATAC-seq, CITE-seq).

Each dataset represents distinct biological challenges: The Adamson (Adamson et al., 2016) and Norman (Norman et al., 2019) datasets capture CRISPR gene knockouts in different cell lines, providing fundamental test cases for genetic perturbation. The Papalexi (Papalexi et al., 2021) dataset offers both RNA and protein measurements (CITE-seq), enabling assessment of cross-modality prediction. The Liscovitch (Liscovitch-Brauer et al., 2021) dataset presents the distinct challenge of predicting chromatin accessibility changes (scATAC-seq) rather than gene expression. The Srivatsan (Srivatsan et al., 2020) dataset assesses the prediction of cellular responses to chemical compounds, while the Schiebinger (Schiebinger et al., 2019) dataset examines responses to immune signaling molecules (cytokines).

To assess generalization to **unseen perturbations**, we select baselines accordingly: scGPT (Cui et al., 2024) for gene knockouts, ChemCPA (Hetzel et al., 2022) for drug treatments. For modalities lacking established models (scCITEseq, scATAC-seq, cytokine), we employ Random Forest and Linear Regression using one-hot encoded perturbations concatenated with expression profiles as inputs.

#### **3.2. Predictive Performance**

Table 1 evaluates the prediction accuracy of models designed by scAgents across diverse perturbation datasets. We employ multiple complementary metrics: Basic prediction accuracy 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  $(R^2 \uparrow)$ , measuring the proportion of expression variance explained by the model. Biological relevance is assessed through metrics on differentially expressed genes (DE) -those exhibiting significant expression changes after perturbation and thus most biologically meaningful. For each dataset, top 20 DE genes are selected based on ground truth expression changes under perturbation. We then compute the same metrics restricted to these DE gene subsets, denoted as MSE<sub>DE</sub>,  $PCC_{DE}$ , and  $R^2_{DE}$ . Across these settings, models designed by scAgents consistently outperform baselines. For gene knockouts (Adamson dataset), our models achieve 49% reduction in prediction error (MSE = 0.0051 vs. 0.0100for scGPT) with improved variance explained ( $R^2 = 0.9761$ vs. 0, 9649). This improvement extends to differentially

#### Submission and Formatting Instructions for ICML 2025 GenBio Workshop

275 276

329

Table 1: Post-perturbation gene expression prediction results, where scAgents-Model is the prediction model automatically designed and implemented by scAgents. The reported metrics for *scAgents-Models* represent the best performance from three separate models automatically designed and trained by the framework, with standard deviations reflecting the variability across these runs.

Mong	MOD	DCC A	$D^2$ $\Rightarrow$	MCE (DE)	DCC (DE) A	$\mathbf{D}^2$ (DE) 4
MODEL	MSE↓	PCC↑	$R^{2}\uparrow$	MSE (DE)↓	PCC (DE)↑	<i>R</i> <sup>~</sup> (DE) ↑
G	ene Knock Out Pertu	rbation – scRNAseq	Dataset (Adamson e	et al. (Adamson et al	., 2016))	
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
scGPT (Cui et al., 2024)	0.0100	0.9861	0.9649	0.2562	0.9088	0.7911
scAgents-Models	$0.0051 \pm 0.0063$	<b>0.9883</b> ± 0.0459	<b>0.9761</b> ± 0.0803	$0.2013 \pm 0.0444$	$0.9474 \pm 0.0601$	<b>0.8912</b> ± 0.0518
	Gene Knock Out Pert	urbation – scRNAse	q Dataset (Norman e	et al. (Norman et al.,	2019))	
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
scGPT (Cui et al., 2024)	0.0076	0.9823	0.9536	0.5318	0.8630	0.5652
scAgents-Models	$0.0034 \pm 0.0023$	<b>0.9846</b> ± 0.0418	<b>0.9609</b> ± 0.0081	$0.1736 \pm 0.0677$	$0.8109 \pm 0.0133$	$0.5975 \pm 0.0539$
Gene	Knock Out Perturba	ution – scCITEseq (K	NA) Dataset (Papal	exi et al. (Papalexi e	t al., 2021))	
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
scAgents-Models	$0.0417 \pm 0.0051$	$0.6935 \pm 0.1995$	$0.3687 \pm 0.0651$	$0.0535 \pm 0.1566$	$0.6406 \pm 0.1940$	$0.2354 \pm 0.0224$
Gene 1	Knock Out Perturbat	ion – scCITEseq (Pr	otein) Dataset (Papa	ılexi et al. (Papalexi	et al., 2021))	
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
scAgents-Models	$0.0070 \pm 0.0387$	$0.7495 \pm 0.0653$	$\textbf{0.6872} \pm 0.0956$	$0.2921 \pm 0.0045$	$0.7409 \pm 0.0970$	$0.5489 \pm 0.0749$
Gene K	nock Out Perturbatio	on – scATACseq Dat	aset (Liscovitch et a	l. (Liscovitch-Braue)	r et al., 2021))	
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
scAgents-Models	$\textbf{0.0327} \pm 0.0430$	$\textbf{0.0855} \pm 0.0357$	$\textbf{0.0678} \pm 0.0120$	$\textbf{0.0406} \pm 0.0268$	$\textbf{0.6991} \pm 0.3173$	$\textbf{0.6400} \pm 0.2779$
	Drug Perturbatio	n – scRNA-seq Data	set (Srivatsan et al.	(Srivatsan et al., 202	20))	
Random Forest	0.5289	0.0527	0.0986	0.6138	0.0945	0.0817
Linear Regression	0.6703	0.0711	0.2826	0.5625	0.0763	0.0421
ChemCPA (Hetzel et al., 2022)	0.0847	0.7221	0.6930	0.1035	0.8053	0.7412
scAgents-Models	$\textbf{0.0053} \pm 0.0290$	$\textbf{0.8664} \pm 0.1332$	$\textbf{0.7137} \pm 0.0740$	$\textbf{0.0080} \pm 0.0835$	$\textbf{0.9278} \pm 0.1001$	$\textbf{0.7887} \pm 0.0548$
(	Cytokine Perturbation	n – scRNA-seq Datas	set (Schiebinger et al	l. (Schiebinger et al.	, 2019))	
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

expressed genes ( $MSE_{DE} = 0.2013$  vs. 0.2562), demonstrating enhanced biological fidelity.

306 The performance advantages become more pronounced 307 in challenging cross-modality scenarios. For CITE-seq 308 protein measurements, scAgents achieves 177% improve-309 ment in correlation (PCC = 0.7495 vs. 0.2704 for Random Forest). For drug perturbations, our models de-311 liver near-perfect predictions (MSE = 0.0053) and 20% 312 higher correlation (PCC = 0.8664 vs. 0.7221) com-313 pared to ChemCPA. Perhaps most remarkably, scAgents 314 maintains superior performance even on fundamentally dif-315 ferent modalities such as chromatin accessibility (scATAC-316 seq), achieving tremendous improvement in variance ex-317 plained ( $R^2 = 0.0678$  vs. 0.0040) and correlation for key 318 regulatory regions (PCC<sub>DE</sub> = 0.6991 vs. 0.0509).

319 The input datasets, task descriptions, and the LLM interface 320 (Claude 3.7) were held constant for each perturbation type throughout these ablation studies; only the internal modules 322 were varied. Table 2 reveals the critical contributions of 323 our key framework components. The agentic retrieval sys-324 tem improves performance substantially over the baseline 325 (PCC from 0.0087 to 0.5643 on the Adamson dataset), demonstrating the importance of domain knowledge inte-327 gration. Similarly, the graph-based discussion component 328

Table 2: Ablation study on the impact of key framework components. Performance comparison of different scAgents components across gene knockout, drug, and cytokine perturbation datasets.

			0			
MODEL	MSE↓	PCC↑	$R^2 \uparrow$	MSE(DE)↓	PCC (DE) ↑	$R^2$ (DE) $\uparrow$
Gene Knock Out Pertur	bation (Aa	lamson Da	taset ( <mark>Ada</mark>	mson et al., 201	6))	
scAgents(baseline)	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
Drug Perturbation (Srivatsan Dataset (Srivatsan et al., 2020))						
scAgents(baseline)	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	0.985	0.6784	0.7548
+ Agentic Retrieval & Graph-Based Discussion	0.0053	0.9881	0.9665	0.0080	0.9953	0.9802
Cytokine Perturbation	(Schiebin	ger Datase	et (Schiebi	nger et al., 2019	))	
scAgents-Model (baseline)	0.5892	0.0065	0.0021	0.5876	0.0797	0.0999
+ Normal RAG	0.4321	0.1765	0.0243	0.4756	- 0.1987	0.0934
+ Agentic Retrieval	0.3456	0.2034	0.2421	0.3076	0.2068	0.1176
+ Graph-Based Discussion	0.3512	0.2051	0.2765	0.2454	0.2239	0.1123
+ Normal RAG & Graph-Based Discussion	0.0987	0.4875	0.4654	0.1065	0.2534	0.1053
+ Agentic Retrieval & Graph-Based Discussion	0.0428	0.5697	0.5042	0.0144	0.3396	0.1240

with self-critic refinement provides complementary benefits (PCC from 0.0087 to 0.5310). The combination of both components yields synergistic effects far exceeding their individual contributions (PCC reaching 0.9883), highlighting how knowledge-guided collaborative reasoning enables effective scientific discovery. This pattern remains consistent across all perturbation types, suggesting scAgents addresses fundamental challenges rather than exploiting dataset-specific characteristics.



Figure 5: The performance of scAgents' RAG compared to stan-339 dard RAG methods. (1) hal: hallucination detection, (2) rel: con-340 text relevance, (3) utl: context utilization. Results are stratified 341 by perturbation type (Drug, Cytokine, Gene). Detailed evaluation 342 methods are stated in Appendix C.

343 Table 3: Performance comparison on scPerturb datasets and benchmark tasks (all values are in %). Results show scAgents 345 consistently outperforms both scGPT and Geneformer across multiple metrics and perturbation types. Each score represents the average of five independent runs, with higher values indicating 347 better performance. 348

MODEL	top5 lin↑	top1 lin ↑	PERT CONS †	top5 knn↑	top1 knn↑	SPEAR CORR $\uparrow$	STRUCT
	Drug	g Perturbation (	Srivatsan Datas	et (Srivatsan et a	1., 2020))		
scGPT(Cui et al., 2024)	5.2	4.4	11.4	5.6	5.1	18.8	54.3
Geneformer(Cui et al., 2023)	4.4	3.1	0.9	5.1	4.8	17.3	54.
scAgents-Model	7.0	4.2	11.4	6.4	5.3	19.1	54.
	Gene Kno	ck Out Perturbe	ation (Adamson I	Dataset (Adamso	n et al., 2016))		
scGPT (Cui et al., 2024)	2.2	0.8	5.6	26.2	25.5	87.3	96.
Geneformer (Cui et al., 2023)	2.1	0.8	4.3	25.9	24.1	86.6	95.
scAgents-Model	2.4	0.9	6.9	26.6	25.9	89.9	96.
	Cytokine	Perturbation (.	Schiebinger Data	set (Schiebinger	et al., 2019))		
scGPT(Cui et al., 2024)	2.1	4.8	4.6	8.2	5.5	66.9	57.
Geneformer(Cui et al., 2023)	1.4	4.2	4.4	8.3	9.9	68.2	57.
scAgents-Model	2.5	5.3	4.9	8.6	8.8	68.5	59.

## 356 3.3. Component Contributions 3.4. Information Integration

357

374

375

358 Evaluation on RAGBench (Friel et al., 2025) with the Pub-359 MedQA dataset (Jin et al., 2019) demonstrates scAgents' 360 ability to accurately identify and contextualize relevant sci-361 entific literature. In this evaluation: hallucination detection 362 (measured by AUROC ) evaluates the system's ability to 363 avoid generating false information; context relevance (measured by RMSE ) assesses how well the retrieved documents match the query; and context utilization (measured by RMSE 366 ) quantifies how effectively the system incorporates retrieved 367 information into responses. In Figure 5, for gene perturba-368 tions, context utilization showed the greatest improvement. 369 Performance of scAgents remains consistent across all per-370 turbation types, suggesting scAgents addresses fundamental 371 challenges rather than exploiting dataset-specific character-372 istics. 373

#### 3.5. Representation Quality Evaluation on scPerturb Benchmark

376 While scAgents primarily performs gene expression predic-377 tion following perturbations, the quality of learned repre-378 sentations is equally important for biological interpretabil-379 ity. Following evaluation practices established in previ-380 ous works (Cui et al., 2024; Theodoris et al., 2023), we 381 benchmark scAgents against specialized foundation models 382 (scGPT & Geneformer) on representation quality metrics 383 (Table 3). We assess different aspects of latent space or-384



Figure 6: We manually prompt four different DeepResearch variants to generate research plans, which were then evaluated by five independent LLMs across eight dimensions, with scores ranging from 1 to 10. Detailed prompt, outputs and scores are provided in Appendix I.

ganization: (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 scAgents 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 scAgents 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, scAgents achieves a consistency of 0.069 vs. 0.056 for scGPT, representing a 23.2% improvement. This indicates that scAgents 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, scAgents achieves a TOP5\_KNN score of 0.064 vs. 0.056 for scGPT, indicating that a higher proportion of test samples have correctly labeled neighbors in embedding space. (4) The Spearman correlation metric (SPEAR\_CORR) evaluates how accurately the latent embeddings can be mapped back to original gene expression space using a linear transformation. The score of 0.191 for drug perturbations (vs. 0.188 for scGPT) represents a higher rank correlation between predicted and actual expression values after linear decoding. (5) Structural integrity (STRUCT\_INT) measures how well control-perturbation relationships are preserved in the latent space. 0.596 for cytokine perturbations (vs. 0.571 for scGPT) indicates that scAgents better maintains the biological relationship between control and perturbed states for complex signaling cascades. We include additional UMAP visualizations in the Appendix O.

Submission and Formatting Instructions for ICML 2025 GenBio Workshop



Figure 7: Human expert evaluations show strong correlation with agent-generated confidence and LLM judge scores. Detailed scores are provided in Appendix I and J.

410

411

412 413

414

415

416

#### 3.6. Evaluation on Task Analysis Module and Method Design Module

Beyond predictive performance, Figure 6 presents expert 417 evaluations comparing scAgents with DeepResearch vari-418 ants. We employ an evaluation protocol involving five inde-419 420 pendent LLM judges that evaluated research plans across eight dimensions. Each judge assesses outputs in a random-421 ized, blinded manner without knowledge of which system 422 423 generated each plan. The evaluation shows that scAgents demonstrates superior performance in dataset analysis and 424 425 baseline defect identification. In method design, scAgents consistently outperforms alternatives on scientific validity, 426 innovation level, experimental design, and impact potential, 427 while maintaining comparable technical feasibility. These 428 advantages are particularly pronounced for cytokine per-429 turbations, where immune signaling complexity demands 430 sophisticated biological reasoning. 431

432<br/>433Notably, both the LLM-assigned scores for task analysis and<br/>research plans, and the confidence scores generated by our<br/>agent system, show strong correlation with evaluations from<br/>three human domain experts (Pearson r = 0.83, p < 0.01 in<br/>Figure 7). Each expert independently spent approximately<br/>10 hours scoring the same outputs under identical criteria,<br/>blinded to the source system.



Figure 8: A manual post-hoc analysis where we categorized and quantified the architectures designed by scAgents for six datasets.

To better understand how scAgents adapts to different perturbation scenarios, we conducted a post-hoc analysis of the architectures it generated. For each of the six datasets, we manually categorized and quantified the neural network components present in scAgents designs. Figure 8 shows the distribution of models across datasets. This analysis reveals how the framework naturally tailored its designs to different data types - for example, preferring Transformers (36.4%) for cytokine perturbations to capture complex signaling dependencies. These adaptive design choices emerged organically from the system's literature-based knowledge and expert discussions rather than through predefined rules or templates.

## 4. Conclusion

scAgents demonstrates how autonomous multi-agent systems can successfully improve single-cell perturbation analysis through integrated expertise across computational, biological, and statistical domains. Our framework achieves consistent performance improvements across diverse datasets, perturbation types, and modalities. The synergistic combination of knowledge integration (via agentic retrieval) and collaborative reasoning (via graph-based discussion) enables discovery of optimal modeling strategies without human intervention, with emergent architecture adaptation to dataset-specific challenges.

## References

- Adamson, B., Norman, T. M., Jost, M., Cho, M. Y., Nuñez, J. K., Chen, Y., Villalta, J. E., Gilbert, L. A., Horlbeck, M. A., Hein, M. Y., et al. A multiplexed single-cell CRISPR screening platform enables systematic dissection of the unfolded protein response. *Cell*, 167(7):1867–1882, 2016.
- Baek, J., Jauhar, S. K., Cucerzan, S., and Hwang, S. J. ResearchAgent: Iterative research idea generation over scientific literature with large language models. *arXiv* preprint arXiv:2404.07738, 2024.

Bai, D., Ellington, C. N., Mo, S., Song, L., and Xing, E. P.

- 440 AttentionPert: accurately modeling multiplexed genetic
  441 perturbations with multi-scale effects. *Bioinformatics*, 40
  442 (Supplement\_1):i453–i461, 2024.
  - Beel, J., Kan, M.-Y., and Baumgart, M. Evaluating sakana's
    AI scientist for autonomous research: Wishful thinking
    or an emerging reality towards' artificial research intelligence'(ARI)? arXiv preprint arXiv:2502.14297, 2025.
  - Bendidi, I., Whitfield, S., Kenyon-Dean, K., Yedder, H. B.,
    Mesbahi, Y. E., Noutahi, E., and Denton, A. K. Benchmarking transcriptomics foundation models for perturbation analysis: one PCA still rules them all, 11 2024. URL
    http://arxiv.org/abs/2410.13956.
  - Boiko, D. A., MacKnight, R., and Gomes, G. Autonomous chemical research with large language models. *Nature*, 623:760–768, 2023. doi: 10.1038/s41586-023-06792-0.
  - Bunne, C., Stark, S. G., Gut, G., Del Castillo, J. S., Levesque, M., Lehmann, K.-V., Pelkmans, L., Krause, A., and Rätsch, G. Learning single-cell perturbation responses using neural optimal transport. *Nature methods*, 20(11):1759–1768, 2023.
  - Chen, T., Anumasa, S., Lin, B., Shah, V., Goyal, A., and Liu,
    D. Auto-Bench: An automated benchmark for scientific discovery in LLMs. *arXiv preprint arXiv:2502.15224*, 2025.
  - Chen, Z., Chen, S., Ning, Y., Zhang, Q., Wang, B., Yu,
    B., Li, Y., Liao, Z., Wei, C., Lu, Z., et al. ScienceAgentBench: Toward rigorous assessment of language agents for data-driven scientific discovery. *arXiv preprint arXiv:2410.05080*, 2024.
  - Cui, H., Wang, C., Maan, H., Pang, K., Luo, F., Duan, N., and Wang, B. scGPT: toward building a foundation model for single-cell multi-omics using generative ai. *Nature Methods*, 21(8):1470–1480, 08 2024. ISSN 1548-7091.
  - Cui, Z., Liao, Y., Xu, T., and Wang, Y. Geneformer: Learned gene compression using transformer-based context modeling. *arXiv preprint arXiv:2212.08379*, 2023.
- Dixit, A., Parnas, O., Li, B., Chen, J., Fulco, C. P., Jerby-Arnon, L., Marjanovic, N. D., Dionne, D., Burks, T., Raychowdhury, R., et al. Perturb-Seq: dissecting molecular circuits with scalable single-cell RNA profiling of pooled genetic screens. *Cell*, 167(7):1853–1866, 2016.
- Dong, M., Wang, B., Wei, J., de O. Fonseca, A. H., Perry, C. J., Frey, A., Ouerghi, F., Foxman, E. F., Ishizuka, J. J., Dhodapkar, R. M., et al. Causal identification of singlecell experimental perturbation effects with CINEMA-OT. *Nature methods*, 20(11):1769–1779, 2023.

- Eger, S., Cao, Y., D'Souza, J., Geiger, A., Greisinger, C., Gross, S., Hou, Y., Krenn, B., Lauscher, A., Li, Y., et al. Transforming science with large language models: A survey on AI-assisted scientific discovery, experimentation, content generation, and evaluation. *arXiv preprint arXiv:2502.05151*, 2025.
- Friel, R., Belyi, M., and Sanyal, A. RAGBench: Explainable benchmark for retrieval-augmented generation systems, 2025. URL http://arxiv.org/abs/2407. 11005.
- Gao, X., Zhang, Z., Xie, M., Liu, T., and Fu, Y. Graph of AI ideas: Leveraging knowledge graphs and llms for AI research idea generation. *arXiv preprint arXiv:2503.08549*, 2025.
- Garikaparthi, A., Patwardhan, M., Vig, L., and Cohan, A. IRIS: Interactive research ideation system for accelerating scientific discovery. *arXiv preprint arXiv:2504.16728*, 2025.
- Ghafarollahi, A. and Buehler, M. J. AtomAgents: Alloy design and discovery through physics-aware multimodal multi-agent artificial intelligence. *arXiv preprint arXiv:2407.10022*, 2024.
- Ghafarollahi, A. and Buehler, M. J. Sparks: Multi-agent artificial intelligence model discovers protein design principles. *arXiv preprint arXiv:2504.19017*, 2025.
- Gu, K., Shang, R., Jiang, R., Kuang, K., Lin, R.-J., Lyu, D., Mao, Y., Pan, Y., Wu, T., Yu, J., et al. BLADE: Benchmarking language model agents for data-driven science. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 13936–13971, 2024.
- Guo, T., Chen, X., Wang, Y., Chang, R., Pei, S., Chawla, N. V., Wiest, O., and Zhang, X. Large language model based multi-agents: A survey of progress and challenges, 02 2024. URL https://arxiv.org/abs/2402. 01680.
- Hao, M., Gong, J., Zeng, X., Liu, C., Guo, Y., Cheng, X., Wang, T., Ma, J., Zhang, X., and Song, L. Large-scale foundation model on single-cell transcriptomics. *Nature methods*, 21(8):1481–1491, 2024.
- Hetzel, L., Boehm, S., Kilbertus, N., Günnemann, S., Theis, F., et al. Predicting cellular responses to novel drug perturbations at a single-cell resolution. *Advances in Neural Information Processing Systems*, 35:26711–26722, 2022.
- Honovich, O., Scialom, T., Levy, O., and Schick, T. True few-shot learning with prompts – a real-world perspective. *Transactions of the Association for Computational Linguistics*, 10:1032–1047, 2022.

Submission and Formatting Instructions for ICML 2025 GenBio Workshop

495 496 497 498 499	House, F. FutureHouse platform: Super- intelligent AI agents for scientific dis- covery. https://www.futurehouse. org/research-announcements/ launching-futurehouse-platform-ai-agent	<ul> <li>Kulkarni, A., Alotaibi, F., Zeng, X., Wu, L., Zeng, T., Yao,</li> <li>B. M., Liu, M., Zhang, S., Huang, L., and Zhou, D.</li> <li>Scientific hypothesis generation and validation: Methods, datasets, and future directions. <i>arXiv preprint</i> arXiv:2505.04651, 2025</li> </ul>
500 501 502 503 504 505 506 507	<ul> <li>2024. Accessed on 2025-05-06; Nonprofit organization developing AI scientist tools for automated research workflows.</li> <li>Hsu, CC., Bransom, E., Sparks, J., Kuehl, B., Tan, C., Wadden, D., Wang, L. L., and Naik, A. CHIME: LLM-assisted hierarchical organization of scientific studies for literature review support. <i>Findings of ACL 2024</i>, 2024.</li> </ul>	<ul> <li>Levine, D., Rizvi, S. A., Lévy, S., Pallikkavaliyaveetil, N., Zhang, D., Chen, X., Ghadermarzi, S., Wu, R., Zheng, Z., Vrkic, I., et al. Cell2Sentence: teaching large language models the language of biology. <i>BioRxiv</i>, pp. 2023–09, 2024.</li> <li>Li, L., Xu, W., Guo, J., Zhao, R., Li, X., Yuan, Y., Zhang, B., Jiang Y. Xin, Y. Dang R. Zhao, D. Rong, Y. Feng, T.</li> </ul>
508 509 510 511 512	Hu, X., Fu, H., Wang, J., Wang, Y., Li, Z., Xu, R., Lu, Y., Jin, Y., Pan, L., and Lan, Z. Nova: An iterative planning and search approach to enhance novelty and diversity of llm generated ideas. <i>arXiv preprint arXiv:2410.14255</i> , 2024.	<ul> <li>and Bing, L. Chain of ideas: Revolutionizing research via novel idea development with llm agents. <i>arXiv preprint arXiv:2410.13185</i>, 2024a.</li> <li>Li, R., Patel, T., Wang, O., and Du, X. MLR-Copilot: Au-</li> </ul>
513 514 515 516	Huang, Q., Vora, J., Liang, P., and Leskovec, J. MLAgent- Bench: Evaluating language agents on machine learning experimentation. In <i>ICML 2024</i> , 2024a.	tonomous machine learning research based on large lan- guage models agents. <i>arXiv preprint arXiv:2408.14033</i> , 2024b.
517 518 519 520	<ul><li>Huang, Y., Luo, J., Yu, Y., Zhang, Y., Lei, F., Wei, Y., He,</li><li>S., Huang, L., Liu, X., Zhao, J., and Liu, K. DA-Code:</li><li>Agent data science code generation benchmark for large language models. <i>EMNLP 2024</i>, 2024b.</li></ul>	Lin, Z. Beyond principlism: practical strategies for ethical ai use in research practices. <i>AI and Ethics</i> , pp. 1–13, 2024.
521 522 523 524 525 526	Jansen, P., Tafjord, O., Radensky, M., Siangliulue, P., Hope, T., Mishra, B. D., Majumder, B. P., Weld, D. S., and Clark, P. Codescientist: End-to-end semi-automated sci- entific discovery with code-based experimentation. <i>arXiv</i> <i>preprint arXiv:2503.22708</i> , 2025.	Liscovitch-Brauer, N., Montalbano, A., Deng, J., Méndez- Mancilla, A., Wessels, HH., Moss, N. G., Kung, C Y., Sookdeo, A., Guo, X., Geller, E., et al. Profiling the genetic determinants of chromatin accessibility with scalable single-cell crispr screens. <i>Nature biotechnology</i> , 39(10):1270–1277, 2021.
527 528 529 530 531	Jiang, Q., Chen, S., Chen, X., and Jiang, R. scPRAM accurately predicts single-cell gene expression perturbation response based on attention mechanism. <i>Bioinformatics</i> , 40(5):btae265, 2024.	Liu, H., Li, Y., Jian, J., Cheng, Y., Lu, J., Guo, S., Zhu, J., Zhang, M., Zhang, M., and Wang, H. Toward a team of AI-made scientists for scientific discovery from gene expression data. <i>arXiv preprint arXiv:2402.12391</i> , 2024a.
532 533 534	Jin, Q., Dhingra, B., Liu, Z., Cohen, W. W., and Lu, X. PubMedQA: A dataset for biomedical research question answering. <i>arXiv preprint arXiv:1909.06146</i> , 2019.	Liu, H., Zhou, Y., Li, M., Yuan, C., and Tan, C. Litera- ture meets data: A synergistic approach to hypothesis generation. <i>arXiv preprint arXiv:2410.17309</i> , 2024b.
535 536 537 538 539	Jing, L., Huang, Z., Wang, X., Yao, W., Yu, W., Ma, K., Zhang, H., Du, X., and Yu, D. DSBench: How far are data science agents to becoming data science experts? <i>arXiv preprint arXiv:2409.07703</i> , 2024.	Liu, Z., Liu, K., Zhu, Y., Lei, X., Yang, Z., Zhang, Z., Li, P., and Liu, Y. AIGS: Generating science from AI-powered automated falsification. <i>arXiv preprint arXiv:2411.11910</i> , 2024c.
540 541 542 543	Kamimoto, K., Stringa, B., Hoffmann, C. M., Jindal, K., Solnica-Krezel, L., and Morris, S. A. Dissecting cell iden- tity via network inference and in silico gene perturbation. <i>Nature</i> , 614(7949):742–751, 2023.	Lotfollahi, M., Wolf, F. A., and Theis, F. J. scGen predicts single-cell perturbation responses. <i>Nature methods</i> , 16 (8):715–721, 2019.
544 545 546 547 548 549	Kon, P. T. J., Liu, J., Ding, Q., Qiu, Y., Yang, Z., Huang, Y., Srinivasa, J., Lee, M., Chowdhury, M., and Chen, A. Curie: Toward rigorous and automated sci- entific experimentation with AI agents. <i>arXiv preprint</i> <i>arXiv:2502.16069</i> , 2025.	Lotfollahi, M., Klimovskaia Susmelj, A., De Donno, C., Hetzel, L., Ji, Y., Ibarra, I. L., Srivatsan, S. R., Naghipour- far, M., Daza, R. M., Martin, B., et al. Predicting cellular responses to complex perturbations in high-throughput screens. <i>Molecular systems biology</i> , 19(6):e11517, 2023.
	1	0

- Lu, C., Lu, C., Lange, R. T., Foerster, J., Clune, J., and Ha,
  D. The AI scientist: Towards fully automated open-ended
  scientific discovery, 09 2024. URL http://arxiv.
  org/abs/2408.06292.
- Majumder, B. P., Surana, H., Agarwal, D., Mishra, B. D.,
  Meena, A., Prakhar, A., Vora, T., Khot, T., Sabharwal,
  A., and Clark, P. DiscoveryBench: Towards data-driven
  discovery with large language models. *arXiv preprint arXiv:2407.01725*, 2024.
- Moerman, T., Aibar Santos, S., Bravo González-Blas, C.,
  Simm, J., Moreau, Y., Aerts, J., and Aerts, S. GRNBoost2
  and Arboreto: efficient and scalable inference of gene
  regulatory networks. *Bioinformatics*, 35(12):2159–2161,
  2019.
- Naumov, V., Zagirova, D., Lin, S., Xie, Y., Gou, W., Urban,
  A., Tikhonova, N., Alawi, K., Durymanov, M., Galkin, F.,
  et al. DORA AI scientist: Multi-agent virtual research
  team for scientific exploration discovery and automated
  report generation. *bioRxiv*, 2025.
- Norman, T. M., Horlbeck, M. A., Replogle, J. M., Ge, A. Y.,
  Xu, A., Jost, M., Gilbert, L. A., and Weissman, J. S.
  Exploring genetic interaction manifolds constructed from
  rich single-cell phenotypes. *Science*, 365(6455):786–793,
  2019.
- 578 OpenAI. Introducing deep research. https://openai. 579 com/index/deep-research/, 2025. Accessed: 580 2025-05-08. 581
- Papalexi, E., Mimitou, E. P., Butler, A. W., Foster, S., Bracken, B., Mauck III, W. M., Wessels, H.-H., Hao, Y., Yeung, B. Z., Smibert, P., et al. Characterizing the molecular regulation of inhibitory immune checkpoints with multimodal single-cell screens. *Nature genetics*, 53 (3):322–331, 2021.
- Peidli, S., Green, T. D., Shen, C., Gross, T., Min, J., Garda,
  S., Yuan, B., Schumacher, L. J., Taylor-King, J. P., Marks,
  D. S., et al. scPerturb: harmonized single-cell perturbation data. *Nature Methods*, 21(3):531–540, 2024.
- Pu, K., Feng, K., Grossman, T., Hope, T., Mishra, B. D., Latzke, M., Bragg, J., Chang, J. C., and Siangliulue, P. IdeaSynth: Iterative research idea development through evolving and composing idea facets with literaturegrounded feedback. *arXiv preprint arXiv:2410.04025*, 2024.
- Qi, B., Zhang, K., Li, H., Tian, K., Zeng, S., Chen, Z.-R., Hu, J.-F., and Zhou, B. Large language models are zero shot hypothesis proposers. *Instruction Workshop @ NeurIPS 2023*, 2023.

- Qiu, X., Zhang, Y., Martin-Rufino, J. D., Weng, C., Hosseinzadeh, S., Yang, D., Pogson, A. N., Hein, M. Y., Min, K. H. J., Wang, L., et al. Mapping transcriptomic vector fields of single cells. *Cell*, 185(4):690–711, 2022.
- Qiu, Y., Zhang, H., Xu, Z., Li, M., Song, D., Wang, Z., and Zhang, K. AI Idea Bench 2025: AI research idea generation benchmark. *arXiv preprint arXiv:2504.14191*, 2025.
- Radensky, M., Shahid, S., Fok, R., Siangliulue, P., Hope, T., and Weld, D. S. Scideator: Human-Ilm scientific idea generation grounded in research-paper facet recombination. *arXiv preprint arXiv:2409.14634*, 2024.
- Reddy, C. K. and Shojaee, P. Towards scientific discovery with generative AI: Progress, opportunities, and challenges. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 28601–28609, 2025.
- Reimers, N. and Gurevych, I. Sentence-BERT: Sentence embeddings using siamese BERT-networks. arXiv preprint arXiv:1908.10084, 2019.
- Ren, S., Jian, P., Ren, Z., Leng, C., Xie, C., and Zhang, J. Towards scientific intelligence: A survey of llm-based scientific agents. arXiv preprint arXiv:2503.24047, 2025.
- Roohani, Y., Huang, K., and Leskovec, J. Predicting transcriptional outcomes of novel multigene perturbations with GEARS. *Nature Biotechnology*, 42(6):927–935, 2024a.
- Roohani, Y. H., Vora, J., Huang, Q., Liang, P., and Leskovec, J. BioDiscoveryAgent: An ai agent for designing genetic perturbation experiments. In *ICLR 2024 Workshop on Machine Learning for Genomics Explorations*, 2024b.
- Ruan, K., Wang, X., Hong, J., Wang, P., Liu, Y., and Sun, H. LiveIdeaBench: Evaluating llms' scientific creativity and idea generation with minimal context. *arXiv preprint arXiv:2412.17596*, 2024.
- Schiebinger, G., Shu, J., Tabaka, M., Cleary, B., Subramanian, V., Solomon, A., Gould, J., Liu, S., Lin, S., Berube, P., et al. Optimal-transport analysis of single-cell gene expression identifies developmental trajectories in reprogramming. *Cell*, 176(4):928–943, 2019.
- Schmidgall, S. and Moor, M. AgentRxiv: Towards collaborative autonomous research. arXiv preprint arXiv:2503.18102, 2025.
- Schmidgall, S., Su, Y., Wang, Z., Sun, X., Wu, J., Yu, X., Liu, J., Liu, Z., and Barsoum, E. Agent laboratory: Using llm agents as research assistants. *arXiv preprint arXiv:2501.04227*, 2025.

- Si, C., Yang, D., and Hashimoto, T. Can LLMs generate
  novel research ideas? a large-scale human study with
  100+ NLP researchers. *arXiv preprint arXiv:2409.04109*,
  2024.
- Skarlinski, M. D., Cox, S., Laurent, J. M., Braza, J. D.,
  Hinks, M., Hammerling, M. J., Ponnapati, M., Rodriques,
  S. G., and White, A. D. Language agents achieve superhuman synthesis of scientific knowledge. *arXiv preprint arXiv:2409.13740*, 2024.
- Skinnider, M. A., Squair, J. W., Kathe, C., Anderson, M. A.,
  Gautier, M., Matson, K. J., Milano, M., Hutson, T. H.,
  Barraud, Q., Phillips, A. A., et al. Cell type prioritization
  in single-cell data. *Nature biotechnology*, 39(1):30–34,
  2021.
- Srivatsan, S. R., McFaline-Figueroa, J. L., Ramani, V., Saunders, L., Cao, J., Packer, J., Pliner, H. A., Jackson, D. L.,
  Daza, R. M., Christiansen, L., et al. Massively multiplex chemical transcriptomics at single-cell resolution. *Science*, 367(6473):45–51, 2020.
- Starace, G., Jaffe, O., Sherburn, D., Aung, J., Shern, C. J., Maksin, L., Dias, R., Mays, E., Kinsella, B., Thompson, W., Heidecke, J., Glaese, A., and Patwardhan, T. PaperBench: Evaluating AI's ability to replicate AI research. arXiv preprint arXiv:2504.01848, 2025. URL https://arxiv.org/abs/2504.01848.
- Su, H., Chen, R., Tang, S., Zheng, X., Li, J., Yin, Z., Ouyang,
  W., and Dong, N. Two heads are better than one: A multiagent system has the potential to improve scientific idea
  generation. *arXiv preprint arXiv:2410.09403*, 2024. URL
  https://arxiv.org/abs/2410.09403.
- Tang, X., Zou, A., Zhang, Z., Li, Z., Zhao, Y., Zhang, X.,
  Cohan, A., and Gerstein, M. MedAgents: Large language
  models as collaborators for zero-shot medical reasoning. *Findings of ACL 2024*, 2024.
- Theodoris, C. V., Xiao, L., Chopra, A., Chaffin, M. D.,
  Al Sayed, Z. R., Hill, M. C., Mantelos, H., Brydon, E. M.,
  Zeng, Z., Liu, X. S., and Ellinor, P. T. Transfer learning
  enables predictions in network biology. *Nature*, 618:
  616–624, 05 2023.
- Tian, M., Gao, L., Zhang, S., Chen, X., Fan, C., Guo, X.,
  Haas, R., Ji, P., Krongchon, K., Li, Y., et al. SciCode:
  A research coding benchmark curated by scientists. *Advances in Neural Information Processing Systems*, 37:
  30624–30650, 2024.
- Wenteler, A., Occhetta, M., Branson, N., Huebner, M.,
  Curean, V., Dee, W., Connell, W., Hawkins-Hooker, A.,
  Chung, P., Ektefaie, Y., et al. PertEval-scFM: Benchmarking single-cell foundation models for perturbation effect
  prediction. *bioRxiv*, pp. 2024–10, 2024.

- Yan, X., Feng, S., Yuan, J., Xia, R., Wang, B., Zhang, B., and Bai, L. SurveyForge: On the outline heuristics, memory-driven generation, and multi-dimensional evaluation for automated survey writing. *arXiv preprint arXiv:2503.04629*, 2025.
- Yang, F., Wang, W., Wang, F., Fang, Y., Tang, D., Huang, J., Lu, H., and Yao, J. scBERT as a large-scale pretrained deep language model for cell type annotation of singlecell rna-seq data. *Nature Machine Intelligence*, 4(10): 852–866, 2022.
- Yang, T., Hu, X., Li, X., Tan, M., Zhang, J., Wen, Z., Chang, E., Dai, A. M., Li, Q. V., Gonzalez, J. E., Cardie, C., and Wei, J. AnyBench: Language models evaluate anything. *arXiv preprint arXiv:2312.13771*, 2023.
- Zhong, R., Zhang, P., Li, S., Ahn, J., Klein, D., and Steinhardt, J. Goal driven discovery of distributional differences via language descriptions. In *NeurIPS 2023*, 2023.

## 660 Limitation

661
 662
 663
 663
 664
 665
 665
 665
 666
 666
 667
 668
 668
 669
 669
 669
 660
 661
 661
 661
 662
 663
 664
 665
 665
 665
 665
 666
 667
 668
 668
 669
 669
 669
 669
 669
 669
 660
 661
 661
 661
 661
 662
 663
 664
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665
 665

Computational and Economic Costs. Our framework requires substantial computational resources and API calls to large
 language models. As shown in Table 15, the average cost per experiment ranges from \$0.38 to \$18.90 depending on the LLM
 backend, which may limit accessibility for resource-constrained research groups. Additionally, the multi-agent discussion
 process can require up to 400,000 output tokens per task, resulting in significant latency.

Failure Modes. Our failure analysis (Section N) reveals that 41% of errors stem from computation execution issues,
 particularly tensor dimension mismatches. While we implemented mitigation strategies, these systematic failures highlight
 the challenges of fully autonomous code generation for complex biological data analysis.

Domain Specificity. Currently, scAgents is specialized for single-cell perturbation prediction tasks. Extending the framework
 to other biological domains (e.g., spatial transcriptomics, proteomics, or multi-omics integration) would require substantial
 architectural modifications and domain-specific knowledge engineering.

**Dependence on Task Specification Quality.** The framework's performance is sensitive to the clarity and completeness of input task descriptions. Ambiguous or incomplete specifications can lead to suboptimal model designs, requiring users to have sufficient domain expertise to formulate well-defined research questions.

Limited Biological Novelty Detection. While scAgents excels at combining existing knowledge, it may struggle to propose
 truly novel biological mechanisms beyond the patterns present in its training data and retrieved literature. This could limit
 its utility for discovering fundamentally new biological phenomena.

## 715 Broader Impact

The development of scAgents has several important implications for the scientific community and society:

718 Democratization of Scientific Research. By automating complex analytical workflows, scAgents lowers the technical 719 barriers for conducting sophisticated single-cell analyses. This could enable smaller research groups and institutions with 720 limited computational expertise to engage in cutting-edge genomics research, potentially accelerating scientific discovery 721 globally.

Acceleration of Therapeutic Development. The ability to rapidly predict cellular responses to perturbations has direct applications in drug discovery and personalized medicine. Our framework could significantly reduce the time and cost of identifying therapeutic targets and understanding drug mechanisms of action.

Scientific Reproducibility. The end-to-end automation and explicit documentation of all analytical decisions enhance
 reproducibility in computational biology. The generated code and research plans provide complete audit trails, addressing a
 critical challenge in modern scientific research.

Ethical Considerations. The automation of scientific discovery raises important questions about attribution, responsibility, and validation of AI-generated hypotheses. While scAgents includes human-interpretable outputs and confidence scores, researchers must carefully validate all predictions before drawing biological conclusions or making therapeutic decisions.

Finite Computational Requirements of running multiple LLM agents have a non-negligible carbon
 footprint. Future work should explore more efficient architectures and consider the environmental cost-benefit trade-offs of
 automated scientific discovery.

Research Workforce Evolution. As AI systems become capable of conducting increasingly sophisticated analyses, the role of computational biologists may shift from implementing methods to critically evaluating AI-generated hypotheses and designs. This transition requires careful consideration of training and career development in the scientific workforce.

Open Science Contribution. By releasing our code and model architectures, we aim to foster community-driven improvements and applications. However, we acknowledge the potential for misuse and encourage responsible deployment with appropriate biological validation and ethical oversight.

# 770 A. Experimental Details

# A.1. Datasets Introduction

Our study leverages six publicly available single-cell perturbation datasets from the scPerturb (Peidli et al., 2024) 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. (Adamson et al., 2016) (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. (Norman et al., 2019) (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.

784 Liscovitch et al. (Liscovitch-Brauer et al., 2021) (ATAC-seq): Employing CRISPRsciATAC, a single-cell combinatorial 785 indexing assay, to delineate the genetic determinants of chromatin accessibility in human myelogenous leukemia K562 786 cells. Targeting 105 chromatin-related genes via CRISPR-Cas9, the study generated chromatin accessibility profiles for 787 approximately 30,000 single cells. Key findings include correlations between the loss of specific chromatin remodelers 788 and global changes in chromatin accessibility. Notably, EZH2 depletion was associated with enhanced accessibility in 789 heterochromatic regions linked to embryonic development and with activation of genes in the HOXA and HOXD clusters. 790 This high-throughput approach offers valuable insights into the role of chromatin modifiers in regulating gene expression 791 and their implications in disease states. 792

Papalexi et al. (Papalexi et al., 2021) (CITE-seq): Combining CRISPR-Cas9 perturbations with single-cell RNA and sur face 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. (Srivatsan et al., 2020) (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. (Schiebinger et al., 2019) (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 typesincluding CRISPRi, CRISPRa, CRISPR-Cas9, smallmolecule drugs, and cytokinesacross 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.

## A.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:

814 **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.

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

822

808

809

## 825 A.3. Knowledge Graph Construction

841

849

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

- Collaborative Agents Shared Knowledge Graph in Task Analysis. Instead of operating in isolation, the Dataset Analyst,
  Problem Investigator, and Baseline Assessor interact via the shared knowledge graph and query interface. Each agent
  incrementally updates the knowledge graph with its findings, while continuously polling for updates from other agents.
  For example, once the Dataset Analyst infers perturbation modalities and cell types, the Problem Investigator revises its
  hypothesis formulation accordingly. Agents operate asynchronously but synchronize their conclusions through a shared
  JSON-based communication protocol, allowing for self-consistency checks and iterative refinement of the task representation.
  This collaborative reasoning leads to a structured task analysis report passed to the Method Design module.
- Graph-Based Expert Shared Knowledge Graph in Method Design. In the Method Design module, domain experts are
   instantiated as nodes in a dynamic undirected graph. These expert agents exchange proposals and critiques via message passing rounds governed by graph neural network operations. Throughout the discussion, the Self-Critic agent monitors
   logical coherence and suggests refinements. Each expert agent has read-write access to the shared knowledge graph and can
   retrieve relevant prior knowledge from Agentic Retrieval. Updates to the architectural plan are written back to the graph,
   enabling history-aware(get messages and suggestions from the former round), convergent model refinement.

## A.4. Experts Discussion Construction Details

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

853 Based on the task analysis report, a set of relevant expert agents is selected by matching task attributes against a curated 854 registry of expert types. The selected experts are grouped into five broad categories to ensure comprehensive domain 855 coverage: (i) Data Engineering and Preprocessing. A Data Expert is instantiated to address normalization, quality control, 856 feature selection, and batch correction issues tailored to the input modality. (ii) Model Design and Scalability. The Model 857 Architecture Expert and Deep Learning Expert are responsible for proposing architectures that balance expressiveness, 858 interpretability, and scalability, considering modality-specific modeling needs. (iii) Biological Plausibility. Single Cell 859 Biologists such as the Pathway Analyst, Drug Response Expert, and Omics Modality Expert contribute domain knowledge 860 to align model components with known biological mechanisms, including gene regulatory networks, cytokine signaling, or 861 pharmacodynamics. (iv) Training and Optimization. A Training Expert is responsible for selecting and justifying the 862 learning algorithm, optimization strategy, regularization, and validation scheme suitable for the data structure and model 863 complexity. (v) Self-Critique and Evaluation. A Self-Critic agent is included in every discussion to promote internal 864 scrutiny, consistency checks, and critical reflection over model assumptions and claims. 865

For example, in a gene knockout task, the system may instantiate the Data Expert to inspect whether the scRNA-seq 866 matrix is properly normalized, whether cell and gene identifiers are standardized, and whether preprocessing sufficiently 867 preserves perturbation-related variation. The Model Architecture Expert and Deep Learning Expert are instantiated to 868 co-design a gene-centric model that integrates perturbation-aware attention and captures target gene dependent regulatory 869 effects. The Pathway Analyst is instantiated to evaluate the role of target gene within interferon signaling cascades, while 870 the Omics Modality Expert assesses whether transcriptomic changes resulting from target gene ablation are robustly 871 captured by scRNA-seq alone. The Training Expert selects dropout-regularized contrastive training and a cell-type-aware 872 sampling scheme to stabilize optimization. The Statistics Expert designs a differential expressionbased evaluation framework 873 and quantifies the significance of target gene induced shifts using FDR-corrected effect sizes. Finally, the Self-Critic is 874 instantiated to identify overfitting risks in rare knockout subsets, challenge latent space linearity assumptions, and refine 875 model outputs for interpretability and robustness. 876

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

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

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

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

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

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

### **B. Evaluation Details**

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

#### B.1. Mean Squared Error (MSE)

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

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

#### **B.2.** Pearson Correlation Coefficient (PCC)

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

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

Then

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

## **B.3. Coefficient of Determination** $(R^2)$

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

$$R^{2} = 1 - \frac{\sum_{i=1}^{n} \|Y_{i} - \hat{Y}_{i}\|^{2}}{\sum_{i=1}^{n} \|Y_{i} - \bar{Y}\|^{2}}$$

#### **B.4.** Metrics with Differential Expression (DE)

Differential expression highlights the genes whose changes drive the biological response to a perturbation, 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 expression vectors under perturbation and control, respectively, with  $Y_{p,i}$ ,  $Y_{c,i} \in \mathbb{R}^{d'}$ . For each gene  $g = 1, \ldots, d'$ , compute the mean expression

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

Quantify the change by the logfoldchange

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

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

$$DE = \{ g : \operatorname{rank}_g(|LFC|) \le K \}, \qquad K = 20.$$

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

990 B.5. LatentSpace Linear Separability

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

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

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

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

# **B.6. Perturbation Consistency**

994 995

1005

1008 1009

1012

1016

1018

1028 1029

1040 1041

1003 This metric assesses the consistency with which a model represents perturbations between different samples and batches. 1004 Let  $\mathcal{P}$  be the set of all gene perturbations. For each  $p \in \mathcal{P}$ , suppose we have  $n_p$  embedding vectors

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

1006 1007 Define the *mean cosinesimilarity score* 

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

Let  $\{S_{q_k}\}_{k=1}^K$  be the corresponding scores for K unexpressed gene controls  $q_k$ . The empirical point of perturbation p is

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

<sup>1014</sup> Finally, the overall *consistency rate* is

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

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

#### **B.7. Latent Space Direct Organization**

This metric evaluates the degree to which perturbation clusters are locally organized in the latent 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 embeddings for the query and reference sets, with the corresponding labels  $y_i$  and  $y_j$ . Set

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

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

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

where C denotes the set of all perturbation labels.

#### <sup>33</sup> B.8. Linear Interpretability of Latent Space

1034 1035 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 two metrics: *Spearman* 1036 *correlation* and *structural integrity*.

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

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

<sup>1043</sup> where rank( $\cdot$ ) returns the within sample rank vector.

conditions within each biological batch. For $b = 1,, B$ batches with $n_b$ samples each, let
$\widetilde{Y}^{(b)}_{ ext{pred}} = Y^{(b)}_{ ext{pred}} - Y^{(b)}_{ ext{pred,ctrl}},  \widetilde{Y}^{(b)}_{ ext{act}} = Y^{(b)}_{ ext{act}} - Y^{(b)}_{ ext{act,ctrl}}.$
Then
$\mathbf{D} = \frac{1}{2} \sum_{i=1}^{B} \frac{1}{  \widetilde{\mathbf{Y}}(b) } = \widetilde{\mathbf{Y}}(b)   = \mathbf{D} = \frac{2}{2} \sum_{i=1}^{B} \frac{1}{  \widetilde{\mathbf{Y}}(b)  }$
$D = \overline{B} \sum_{l=1}^{l} \overline{n_b} \  Y_{\text{pred}} - Y_{\text{act}} \ _F,  D_{\max} \approx \overline{B} \sum_{l=1}^{l} \overline{n_b} \  Y_{\text{act}} \ _F,$
b=1 $b=1$
D
$SI = 1 - \frac{D}{D}$ ,
$\mathcal{D}_{\max}$
with higher SI indicating better preservation of control perturbation structure.

<b>C. RAGBench Evaluation Details</b> To evaluate the performance of scAgents' Agentic Retrieval system in Task Analysis Module, we employ RAGBench(Frie et al., 2025). We first align our systems outputs to the format expected by RAGBench. Each output record must include: • dc: unique sample identifier: • response: the generated answer. Refer to constants.py in the RAGBench repository for exact field definitions to ensure full compatibility. Then we numiference on our system's outputs with evaluation models. Trulens and dataset PubMedQA(In et al., 2019). Detailed formulations of the metrics used in this RAGBench benchmark are as follow: Halucination Detection (Hal) In Retrieval-Augmented Generation (RAG) systems, <i>hallucination</i> 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 votation such unsupported information. Reliable RAG outputs demand faith/liness 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:	-	Submission and Formating instructions for Forth 2025 Genero Workshop
To evaluate the performance of scAgents' Agentic Retrieval system in Task Analysis Module, we employ RAGBech(Frie et al., 2025). 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 Tratens and dataset PubMedQA(Jin et al., 2019). Detailed formulations of the metrics used in this RAGBench benchmark are as follow: <b>Hallucination Detection (Hal)</b> In Retrieval-Augmented Generation (RAG) systems, <i>hallucination</i> refers to the generation of content not grounded in the retrieved contextin 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: trues <sub>halluc</sub> = ¬trues <sub>halluwence</sub> , preds <sub>halluc</sub> [mask], preds <sub>halluc</sub> [mask]), we define hallucination labels by trues <sub>halluc</sub> = ¬trues <sub>halluwence</sub> , preds <sub>halluc</sub> [mask]), where ROC_AUC denotes the standard implementation (sklearn.metrics.roc_auc_score). <b>Context relevance assesses</b> 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 consures that the generator receives useful evidence. We measure relevance via Root Mean Squared Eirror (RMSE) between true and predicted r		C. RAGBench Evaluation Details
We first align our systems outputs to the format expected by RAGBench. Each output record must include: • i.d: unique sample identifier; • documents: ist 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(Jin et al., 2019). Detailed formulations of the metrics used in this RAGBench benchmark are as follow: <b>Hallucination Detection (Hal)</b> In Retrieval-Augmented Generation (RAG) systems, <i>hallucination</i> refers to the generation of content not grounded in the retrieved contextin 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:	] e	To evaluate the performance of scAgents' Agentic Retrieval system in Task Analysis Module, we employ RAGBech(Friel et al., 2025).
• i.d: unique sample identifier; • document s: 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(Jin et al., 2019). Detailed formulations of the metrics used in this RAGBench benchmark are as follow: <b>Hallucination Detection (Hal)</b> In Retrieval-Angmented Generation (RAG) systems, <i>hallucination</i> refers to the generation of content not grounded in the retrieved context in other words, the model makes up facts. Hallucination detection quantifies the systems protein 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 generation and 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 genomence. Given:	١	We first align our systems outputs to the format expected by RAGBench. Each output record must include:
• 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(Jin et al., 2019). Detailed formulations of the metrics used in this RAGBench benchmark are as follow: <b>Hallucination Detection (Hal)</b> In Retrieval-Augmented Generation (RAG) systems, <i>hallucination</i> refers to the generation of content not grounded in the retrieved contextin 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: rues <sub>adherence</sub> { True, False}, preds <sub>adherence</sub> $\in [0, 1]$ , we define hallucination labels by rues <sub>halluc</sub> = -trues <sub>adherence</sub> , preds <sub>halluc</sub> = 1 - preds <sub>adherence</sub> . Let mask = ¬isnan(preds <sub>halluc</sub> ). Then AUROC = ROC_AUC (trues <sub>halluc</sub> [mask], preds <sub>halluc</sub> [mask]), where ROC_AUC denotes the standard implementation (sklearn.metrics.roc_auc_score). <b>Context Relevance (RB)</b> 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: $RMSE = \sqrt{\frac{1}{n}} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$ , where $y_i$ is the gold relevance score, $\hat{y}_i$ the predicted score, and $n$ the numb		• id: unique sample identifier;
• 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(Jin et al., 2019). Detailed formulations of the metrics used in this RAGBench benchmark are as follow: <b>Hallucination Detection (Hal)</b> In Retrieval-Augmented Generation (RAG) systems, <i>hallucination</i> refers to the generation of content not grounded in the retrieved contextin other words, the model makes up facts. Hallucination detection measures whether the models outputs contain such nusupported 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: ruces <sub>adherence</sub> $\in \{\text{True, False}\}$ , $\operatorname{preds}_{adherence} \in [0, 1]$ , we define hallucination labels by trues <sub>hallace</sub> $\neg$ trues <sub>adherence</sub> , $\operatorname{preds}_{hallac}$ ], $\operatorname{preds}_{hallace}[mask]$ ), where ROC_AUC denotes the standard implementation (sklearn.metrics.roc_auc_score). <b>Context Relevance (Rel)</b> 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. $RMSE = \sqrt{\frac{1}{n}} \sum_{i=1}^{n} (\mu_i - \hat{y}_i)^2$ , where $\mu_i$ is the gold relevance score, $\hat{y}_i$ the predicted score, and $n$ the number of examples. We ignore any NaN predictions		• documents: list of retrieved context documents;
• 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(Jin et al., 2019). Detailed formulations of the metrics used in this RAGBench benchmark are as follow: <b>Hallucination Detection (Hal)</b> In Retrieval-Augmented Generation (RAG) systems, <i>hallucination</i> 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: rutes <sub>adherence</sub> $\in {True, False}$ , $preds_{adherence} \in [0, 1]$ , we define hallucination labels by rutes <sub>halluc</sub> $= \neg trues_{adherence}$ , $preds_{halluc} = 1 - preds_{adherence}$ . Let mask $= \neg isnan(preds_{halluc})$ . Then AUROC $= ROC_AUC(trues_{halluc}[mask], preds_{halluc}[mask])$ , where ROC_AUC denotes the standard implementation (sklearn.metrics.roc_auc_score). <b>Context Relevance (Rel)</b> 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 prorequisite for accurate generation. Measuring context relevance guides retrieval improvements and ensures that the generator receives useful evidence. $RMSE = \sqrt{\frac{1}{n}} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$ , 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 and ensures that the gold relevance score, $\hat{y}_i$ the predicted score, and $n$ the number of examples. We ignor		• question: the query text;
Refer to constants.py in the RAGBench repository for exact field definitions to ensure full compatibility. Then we numericate on our system's outputs with evaluation models Trulens and dataset PubMedQA(lin et al., 2019). Detailed formulations of the metrics used in this RAGBench benchmark are as follow: <b>Hallucination Detection (Hal)</b> In Retrieval-Augmented Generation (RAG) systems, <i>hallucination</i> refers to the generation of content not grounded in the retrieved contextin 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: trues <sub>halluce</sub> = ¬trues <sub>adherence</sub> , preds <sub>halluc</sub> = 1 - preds <sub>adherence</sub> . Let mask = ¬isnan(preds <sub>halluc</sub> ). Then AUROC = ROC_AUC(trues <sub>halluc</sub> [mask], preds <sub>halluc</sub> [mask]), where ROC_AUC denotes the standard implementation (sklearn.metrics.roc_auc.score). <b>Context Relevance (Rel)</b> 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: $RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2},$ where $y_i$ is the gold relevance score, $\hat{y}_i$ the predicted score, and $n$ the number of examples. We ignore any NaN prediction		• response: the generated answer.
Detailed formulations of the metrics used in this RAGBench benchmark are as follow: Hallucination Detection (Hal) In Retrieval-Augmented Generation (RAG) systems, <i>hallucination</i> 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: trues <sub>adherence</sub> $\in [True, False]$ , preds <sub>adherence</sub> $\in [0, 1]$ , we define hallucination labels by trues <sub>halluc</sub> $= \neg trues_{adherence}$ , preds <sub>halluc</sub> $= 1 - preds_{adherence}$ . Let mask $= \neg isnan(preds_{halluc})$ . Then AUROC $= ROC\_AUC(trues_{halluc}[mask], preds_{halluc}[mask])$ , 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: $RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2},$ 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 the prediction for accurate score, and $n$ the number of examples.	H i	Refer to constants.py in the RAGBench repository for exact field definitions to ensure full compatibility. Then we run nference on our system's outputs with evaluation models Trulens and dataset PubMedQA(Jin et al., 2019).
Hallucination Detection (Hal) In Retrieval-Augmented Generation (RAG) systems, <i>hallucination</i> 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: rues <sub>adherence</sub> $\in \{\text{True, False}\}, \text{ preds}_{adherence} \in [0, 1],$ we define hallucination labels by rues <sub>halluc</sub> = $\neg \text{trues}_{adherence}, \text{ preds}_{halluc} = 1 - \text{preds}_{adherence}.$ Let mask = $\neg \text{isnan}(\text{preds}_{halluc}).$ Then AUROC = ROC_AUC (trues <sub>halluc</sub> [mask], preds <sub>halluc</sub> [mask]), where ROC_AUC denotes the standard implementation (sklearn.metrics.roc_auc_score). <b>Context Relevance (Rel)</b> 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 correct answer. High relevance via Root Mean Squared Error (RMSE) between true and predicted relevance scores: $RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2},$ 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	I	Detailed formulations of the metrics used in this RAGBench benchmark are as follow:
In Retrieval-Augmented Generation (RAG) systems, <i>hallucination</i> 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: $trues_{adherence} \in \{True, False\},  preds_{adherence} \in [0, 1],$ we define hallucination labels by $trues_{halluc} = \neg trues_{adherence},  preds_{halluc} = 1 - preds_{adherence}.$ Let $mask = \neg isnan(preds_{halluc}).$ Then $AUROC = ROC_AUC (trues_{halluc}[mask], preds_{halluc}[mask]),$ where ROC_AUC denotes the standard implementation (sklearn.metrics.roc_auc_score). <b>Context Relevance (Rel)</b> 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 consures useful evidence. We measure relevance via Root Mean Squared Error (RMSE) between true and predicted relevance scores: $RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2},$ 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	ł	Hallucination Detection (Hal)
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: trues <sub>adherence</sub> $\in$ {True, False}, preds <sub>adherence</sub> $\in$ [0, 1], we define hallucination labels by trues <sub>halluc</sub> = $\neg$ trues <sub>adherence</sub> , preds <sub>halluc</sub> = 1 - preds <sub>adherence</sub> . Let mask = $\neg$ isnan(preds <sub>halluc</sub> ). Then AUROC = ROC_AUC(trues <sub>halluc</sub> [mask], preds <sub>halluc</sub> [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: $RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2}$ , 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	I r c	In Retrieval-Augmented Generation (RAG) systems, <i>hallucination</i> refers to the generation of content not grounded in the retrieved contextin other words, the model makes up facts. Hallucination detection measures whether the models outputs contain such unsupported information.
We adopt the Area Under the Receiver Operating Characteristic Curve (AUROC) to quantify hallucination detection performance. Given: $trues_{adherence} \in \{True, False\},  preds_{adherence} \in [0, 1],$ we define hallucination labels by $trues_{halluc} = \neg trues_{adherence},  preds_{halluc} = 1 - preds_{adherence}.$ Let $mask = \neg isnan(preds_{halluc}).$ Then $AUROC = ROC\_AUC(trues_{halluc}[mask], preds_{halluc}[mask]),$ where ROC\_AUC denotes the standard implementation (sklearn.metrics.roc\_auc\_score). <b>Context Relevance (Rel)</b> 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: $RMSE = \sqrt{\frac{1}{n}\sum_{i=1}^{n}(y_i - \hat{y}_i)^2},$ 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	I I	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.
$trues_{adherence} \in \{True, False\},  preds_{adherence} \in [0, 1],$ we define hallucination labels by $trues_{halluc} = \neg trues_{adherence},  preds_{halluc} = 1 - preds_{adherence}.$ Let $mask = \neg isnan(preds_{halluc}).$ Then $AUROC = ROC\_AUC(trues_{halluc}[mask], preds_{halluc}[mask]),$ 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: $RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2},$ 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 to result in the sum of the standard in the number of examples.	۱ ۲	We adopt the Area Under the Receiver Operating Characteristic Curve (AUROC) to quantify hallucination detection performance. Given:
we define hallucination labels by $trues_{halluc} = \neg trues_{adherence},  preds_{halluc} = 1 - preds_{adherence}.$ Let $mask = \neg isnan(preds_{halluc}).$ Then $AUROC = ROC\_AUC(trues_{halluc}[mask], preds_{halluc}[mask]),$ where ROC\_AUC denotes the standard implementation (sklearn,metrics,roc_auc_score). <b>Context Relevance (Rel)</b> 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: $RMSE = \sqrt{\frac{1}{n}\sum_{i=1}^{n}(y_i - \hat{y}_i)^2},$ 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		$trues_{adherence} \in \{\texttt{True}, \texttt{False}\},  preds_{adherence} \in [0,1],$
$trues_{halluc} = \neg trues_{adherence},  preds_{halluc} = 1 - preds_{adherence}.$ Let $mask = \neg isnan(preds_{halluc}).$ Then $AUROC = ROC\_AUC(trues_{halluc}[mask], preds_{halluc}[mask]),$ 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: $RMSE = \sqrt{\frac{1}{n}\sum_{i=1}^{n}(y_i - \hat{y}_i)^2},$ 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 the store true and predictions of the score true and predictions of the score true true and predictions of the score true true true true true and predicted relevance scores: $RMSE = \sqrt{\frac{1}{n}\sum_{i=1}^{n}(y_i - \hat{y}_i)^2},$	V	we define hallucination labels by
Let $mask = \neg isnan(preds_{halluc}).$ Then $AUROC = ROC\_AUC(trues_{halluc}[mask], preds_{halluc}[mask]),$ where ROC\_AUC denotes the standard implementation (sklearn.metrics.roc\_auc\_score). <b>Context Relevance (Rel)</b> 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: $RMSE = \sqrt{\frac{1}{n}\sum_{i=1}^{n}(y_i - \hat{y}_i)^2},$ 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		$trues_{halluc} = \neg trues_{adherence},  preds_{halluc} = 1 - preds_{adherence}.$
Then $AUROC = ROC\_AUC(trues_{halluc}[mask], preds_{halluc}[mask]),$ where ROC\_AUC denotes the standard implementation (sklearn.metrics.roc_auc_score). <b>Context Relevance (Rel)</b> 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: $RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2},$ 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	Ι	$mask = \neg isnan(preds_{halluc}).$
where ROC_AUC denotes the standard implementation (sklearn.metrics.roc_auc_score). <b>Context Relevance (Rel)</b> 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: $RMSE = \sqrt{\frac{1}{n}\sum_{i=1}^{n}(y_i - \hat{y}_i)^2},$ 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	]	Then $AUROC = ROC\_AUC(trues_{halluc}[mask], preds_{halluc}[mask]),$
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: $RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2},$ 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	V	where ROC_AUC denotes the standard implementation (sklearn.metrics.roc_auc_score).
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: $RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2},$ 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	(	Context Relevance (Rel)
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: $RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2},$ 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	(	Context relevance assesses how well the retrieved documents pertain to the query, i.e. whether the context can support a correct answer.
We measure relevance via Root Mean Squared Error (RMSE) between true and predicted relevance scores: $RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2},$ 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	H e	High relevance is a prerequisite for accurate generation. Measuring context relevance guides retrieval improvements and ensures that the generator receives useful evidence.
$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2},$ 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	V	We measure relevance via Root Mean Squared Error (RMSE) between true and predicted relevance scores:
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		$ ext{RMSE} = \sqrt{rac{1}{n}\sum_{i=1}^n ig(y_i - \hat{y}_iig)^2},$
by masking.	v t	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.

## 1155 Context Utilization (Utl)

 $\frac{1156}{1157}$  Context utilization evaluates the extent to which the model leverages the retrieved context when generating its responses.

1158 Even with relevant context, a model may underuse it. This metric reveals the generators ability to integrate context 1159 information into its output.

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

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

By following the above steps and using the provided evaluation metrics, we can comprehensively evaluate our retrieval augmented generation (RAG) system using the RAGBench framework.

1210	D. Detailed Algorithm Specifications
1211 1212	D.1. Agentic Retrieval System
1213 1214 1215 1216	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.
1210	D.1.1. QUERY CONSTRUCTION AND INITIALIZATION
1210	Given a task description $T$ and dataset metadata $D$ , we first construct an initial query representation: Algorithm 1 Ouery Construction
1220 1221	1. input <i>T</i> , <i>D</i>
1222	2. set keywords $\leftarrow$ ExtractKeyTerms $(T) \cup$ ExtractKeyTerms $(D)$
1224 1225	3. set embedding $\leftarrow$ SentenceBERT(keywords)
1226 1227	4. set $Q^{(0)} \leftarrow \text{NormalizeVector}(\text{embedding})$
1228 1229	5. output $Q^{(0)}$
1230	
1231 1232 1233 1234	The function ExtractKeyTerms performs domain-specific extraction of biological entities (genes, cell types, perturbation methods) and technical terms (model architectures, evaluation metrics) using named entity recognition enhanced with domain-specific dictionaries.
1235 1236	D.1.2. Alternating Search Strategy
1237 1238	Our multi-layer retrieval process alternates between breadth-first and depth-first search modes to balance exploration and
1239	Algorithm 2 Alternating BFS-DFS Retrieval
1240 1241 1242	1. input $Q^{(0)}, L_{\max},  au, \epsilon$
1242	2. set $t \leftarrow 0$
1244 1245	3. set $\mathcal{N}_0 \leftarrow \emptyset$
1246 1247	4. set $\mathcal{D} \leftarrow \emptyset$ {Document collection}
1248	5. while $t < L_{\max}$ do
1249 1250 1251 1252 1253	(a) if $t \mod 2 = 1$ then {BFS layer (odd $t$ )} i. set $\mathcal{N}_t \leftarrow \text{TopK}(Q^{(t)}, \text{mode} = \text{BFS})$ (b) else {DFS layer (even $t$ )} i. set $\mathcal{N}_t \leftarrow \text{FollowCitations}(\mathcal{N}_{t-1})$
1254 1255 1256	(c) set $\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{N}_t$ (d) set $Q^{(t+1)} \leftarrow \text{UpdateQuery}(Q^{(t)}, \mathcal{N}_t)$ (e) if $\text{Overlap}(Q^{(t+1)}, Q^{(t)}) > \tau$ then break
1257 1258	(f) if $\max_{d \in \mathcal{N}_t} \operatorname{Score}(Q^{(t)}, d) < \epsilon$ then break
1259 1260 1261	(g) set $\iota \leftarrow \iota + 1$ 6. output $\mathcal{D}$

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

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

 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:

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

(1)

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

**Overlap Computation.** Query overlap is calculated as:

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

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

#### **D.2. Graph-based Multi-Expert Discussion**

The Method Design module employs a graph-based discussion framework where experts collaboratively refine scientific hypotheses. Here we formalize the complete algorithm:

- Algorithm 3 Graph-based Expert Discussion
- 1. **input** TaskAnalysis,  $\tau$ ,  $\epsilon$ ,  $T_{max}$
- 2. set  $E^{(k)} \leftarrow$  SelectExperts(TaskAnalysis)
- 3. set  $S \leftarrow$  InitializeSelfCritic()
- 4. set  $G^{(k)} \leftarrow (S, E^{(k)})$  {Initialize collaboration graph}
- 5. for i = 1 to k do
  - (a) set  $c_0^{(i)} \leftarrow 0$  {Initialize confidence scores}
  - (b) set  $m_0^{(i)} \leftarrow \text{InitialProposal}(E^{(i)}, \text{TaskAnalysis})$
- 6. set  $t \leftarrow 0$

```
7. while t < T_{\text{max}} do
                  (a) for i = 1 to k do
                         i. set m_t \leftarrow \text{Integrate}(\{m_t^{(j)}\}_{i=1}^k)
                        ii. set c_t^{(i)} \leftarrow UpdateConfidence(c_{t-1}^{(i)}, m_t, S)
                 (b) if \forall i : c_t^{(i)} \ge \tau \land \|c_t^{(i)} - c_{t-1}^{(i)}\| < \epsilon then break
                  (c) for i = 1 to k do
                 i. set m_{t+1}^{(i)} \leftarrow \text{RefineIdea}(E^{(i)}, m_t)
(d) set t \leftarrow t+1
1314
1315
            8. set ResearchPlan \leftarrow FinalizeProposal(m_t)
1316
1317
            9. output ResearchPlan
```

**Expert Selection.** The expert selection procedure dynamically assembles a team of domain specialists based on task requirements: 1323  $P(E^{(i)}|\text{TaskAnalysis}) \propto \exp(\beta \cdot \text{Relevance}(E^{(i)}, \text{TaskAnalysis}))$ (4)1324 where  $\beta$  is a temperature parameter controlling selection diversity. 1326 1327 Confidence Update Rule. The confidence score update incorporates feedback from both the Self-Critic and other experts: 1328 1329  $c_t^{(i)} = \lambda_1 \cdot c_{t-1}^{(i)} + \lambda_2 \cdot \text{SelfCriticScore}(m_t^{(i)}, S) + \lambda_3 \cdot \frac{1}{k-1} \sum_{i \neq i} \text{PeerScore}(m_t^{(i)}, E^{(j)})$ 1330 (5) 1331 1333 where  $\lambda_1 + \lambda_2 + \lambda_3 = 1$  weights the relative importance of each component. 1334 **Message Integration.** Expert proposals are integrated through a weighted combination:  $m_t = \sum_{i=1}^k w_t^{(i)} \cdot m_t^{(i)}$ 1338 (6)1339 1340 1341 where weights  $w_t^{(i)}$  are derived from normalized confidence scores: 1342 1343  $w_t^{(i)} = \frac{\exp(c_t^{(i)})}{\sum_{j=1}^k \exp(c_t^{(j)})}$ 1344 (7)1345 1346 1347 This soft-voting mechanism ensures that higher-confidence perspectives have greater influence while still preserving diversity 1348 of thought. 1349 1350 **D.3.** Validation and Refinement Process 1351 The Validation Agent employs an iterative refinement process that systematically improves implementation quality: 1352 1353 Algorithm 4 Iterative Implementation Refinement 1355 1. input ModelDesign, Dataset,  $R_{\max}$ 1356 1357 2. set  $Code_0 \leftarrow InitialImplementation(ModelDesign)$ 1358 1359 3. set Performance<sub>0</sub>  $\leftarrow$  Evaluate(Code<sub>0</sub>, Dataset) 1360 4. for r = 1 to  $R_{\text{max}}$  do 1361 1362 (a) set  $\operatorname{Errors}_r \leftarrow \operatorname{IdentifyIssues}(\operatorname{Code}_{r-1}, \operatorname{Performance}_{r-1})$ (b) set  $Code_r \leftarrow RefineImplementation(Code_{r-1}, Errors_r)$ 1364 (c) set Performance<sub>r</sub>  $\leftarrow$  Evaluate(Code<sub>r</sub>, Dataset) 1365 (d) if Performance<sub>r</sub> – Performance<sub>r-1</sub>  $< \delta$  then break 1366 5. output  $Code_r$ 1369 1370 **Error Analysis.** The error identification procedure categorizes implementation issues into distinct types: 1372 • Logical errors: Incorrect algorithm implementation 1374

- 1375 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 N.

# **D.4. Hyperparameter Configuration**

1387 Our framework employs the following hyperparameter settings, determined through empirical validation on held-out 1388 scientific tasks:

Module	Parameter	Value
	$L_{\max}$	10
Agentic Retrieval	au	0.8
-	$\epsilon$	0.5
	$T_{\rm max}$	6
	au	0.8
Expert Discussion	$\epsilon$	1
	$(\lambda_1,\lambda_2,\lambda_3)$	(0.3, 0.4, 0.3)
Mali dati an	$R_{\rm max}$	5
vandation	δ	0.01

Table 4: Hyperparameter Configuration

 $\begin{array}{l} \begin{array}{l} 1404\\ 1405\\ 1406\\ 1407\end{array} \\ \begin{array}{l} \text{These parameters balance convergence speed with solution quality across different perturbation types and dataset character$ istics. We observed that the Expert Discussion module particularly benefits from a higher weight on Self-Critic evaluation $($\lambda_2$), which promotes more rigorous scientific validation. } \end{array}$ 

Table 5: Performance and Cost Impact of Hyperparameter Settings

Module	Setting	Metric	Score	Change	Time Cost (s)	Token $\times$
Agentic Retrieval	$\begin{aligned} \epsilon &= 0.3 \\ \epsilon &= 0.5 \end{aligned}$	Relevance Score Relevance Score	74.5 75.1	(0.6)	14.1 28.3	1.0× 1.6×
Expert Discussion	$\lambda_2 = 0.2$ $\lambda_2 = 0.4$ $\lambda_2 = 0.6$	Validity Score Validity Score Validity Score	69.1 <b>73.4</b> 72.2	(4.3)	25.2 27.0 28.5	$3.5 \times$ $4.2 \times$ $4.6 \times$
	$\lambda_2 = 0.4$	Consistency Score	76.1	(5.9)	17.0	4.2×
Validation	$\begin{split} \delta &= 0.02 \\ \delta &= 0.01 \end{split}$	Acceptance Rate Acceptance Rate	82.1 81.7		3.2 2.9	1.0× 0.7×

## 1430 E. Agent Communication Protocol Details

#### 1431 1432 **E.1. Protocol Design and Comparison**

The SCAGENTS protocol represents an advancement in agent communication architectures designed specifically for scientific
 discovery. Figure 3 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 knowledge graph. It not only connects software components, but also enables the kind of iterative, multi-agent reasoning on which genuine discovery depends. The scAgents' protocol method allows agents to coordinate autonomously when tasked with comprehensive scientific research.

Table 6 provides a detailed comparison of SCAGENTS 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 6:	Comparison	of Agent	Communication	Protocols
----------	------------	----------	---------------	-----------

Protocol	Context	Interop.	Msg. Struct.	Use Cases
MCP (Anthropic)	1	×	JSON-RPC only	Tool Use & Data Access
Agent2Agent (Google)	×	1	JSON-RPC event	Cross-agent Collaboration
ACP (BeeAI/IBM)	1	1	RESTful	Local Orchestration
scAgents	1	1	JSON-RPC event + Knowledge Graph	End-to-end Scientific Discover

# 1455 E.2. Protocol Implementation Details1456

1465

1467

1469

1470 1471

1472

1477

1450 The SCAGENTS protocol implementation consists of two primary components:

IA58
 IA59
 IA59
 IA59
 IA60
 asynchronous event handling. Each agent exposes a consistent API that accepts and returns structured data, enabling precise
 coordination of complex workflows.

Knowledge Graph Integration Layer Beyond simple message passing, SCAGENTS maintains a persistent knowledge
 graph that captures:

- Research entities (datasets, methods, metrics, results)
- Relationships between entities (causal, hierarchical, temporal)
  - Provenance information (confidence scores, reasoning chains, citations)
  - Domain-specific knowledge (regulatory pathway information, gene-gene interactions)

1473 This dual-layer approach provides several advantages compared to prior protocols: 1474

- 1475
   1476
   1. Context-awareness: Agents maintain awareness of the overall research state through the knowledge graph, enabling them to make more informed decisions.
- 1478 2. Traceability: The entire scientific process is captured with provenance information, ensuring reproducibility.
- 1479
  1480
  1481
  3. Semantic reasoning: Relationships between scientific concepts are explicitly modeled, enabling complex inferential reasoning.
- 1482
   1483
   1484
   4. Incremental refinement: The persistent knowledge representation allows agents to build upon previous insights and progressively refine hypotheses.

1485 In scientific research contexts, these capabilities are essential for managing the complexity of cross-disciplinary knowledge 1486 integration required for tasks like single-cell perturbation analysis.

#### <sup>1488</sup> 1489 **F. Prompt Templates**

1487

## 1490 F.1. Task Description input

1491	
1492	Task Desciption Input
1493	
1/0/	Your task is to develop a predictive model that accurately estimates gene
1405	expression profiles of individual K562 cells following CRISPR interference (
1495	CRISPRI), Using the dataset from Norman et al. (2019, Science).
1490	Task Definition:
1497	- Input: Baseline gene expression profile of an unperturbed K562 cell and the
1498	Identity of the target gene(s) for perturbation
1499	Sucpue. Heatered gene expression profile after percurbation
1500	Evaluation Sconarios.
1501	1. Unseen Perturbations: Predict effects of gene perturbations not present during
1502	training
1503	2. Unseen Cell Contexts: Predict responses in cells with gene expression profiles
1504	not observed during training
1505	Evaluation Metrics:
1506	- Mean Squared Error (MSE): Measures the average squared difference between
1507	predicted and observed gene expression.
1508	predicted and observed profiles.
1500	- R\$^2\$ (Coefficient of Determination): Represents the proportion of variance in
1510	the observed gene expression that can be explained by the predicted values.
1510	- MSE for Differentially Expressed (DE) Genes (MSE_DE): Same as MSE but computed
1510	- PCC for Differentially Expressed (DE) Genes (PCC DE): Same as PCC but computed
1512	specifically for genes identified as differentially expressed.
1513	- R\$^2\$ for Differentially Expressed (DE) Genes (R\$^2\$_DE): Same as R\$^2\$ but
1514	computed specifically for genes identified as differentially expressed.
1515	
1516	
1517	
1518	
1519	
1520	
1521	
1522	
1523	
1524	
1525	
1526	
1527	
1528	
1520	
1520	
1521	
1520	
1532	
1533	
1534	
155-	
1535	
1535 1536	
1535 1536 1537	
1535 1535 1536 1537 1538	

1	540	
1	541	
1	542	
1	543	
1	544	
1	545	
1	545	
1	540	
1	547	
1	548	
1	549	
1	550	
1	551	
1	552	
1	553	
1	554	
1	555	
1	556	
1	550	
1	550	
1	550	
1	559	
1	560	
1	561	
1	562	
1	563	
1	564	
1	565	
1	566	
1	567	
1	568	
1	560	
1	570	
1	570	
1	5/1	
1	572	
1	573	
1	574	
1	575	
1	576	
1	577	
1	578	
1	579	
1	580	
1	581	
1	501	
1	202	
	502	
1	583	
1	583 584	
1	583 584 585	
1 1 1	583 584 585 586	
1 1 1 1	583 584 585 586 587	
1 1 1 1 1	583 584 585 586 586 587 588	
1 1 1 1 1 1	583 584 585 586 586 587 588 589	
1 1 1 1 1 1 1 1	583 584 585 586 587 588 589 590	
1 1 1 1 1 1 1 1 1	583 584 585 586 586 587 588 589 590 591	
1 1 1 1 1 1 1 1 1 1	583 584 585 586 587 588 589 590 591 592	
1 1 1 1 1 1 1 1 1 1 1	583 584 585 586 587 588 589 590 591 592	

#### 0 F.2. Task Analysis Collaboration Agents Settings

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

#### **Task Analysis Collaboration Agent Settings**

#### Data Analyst

You are a Dataset Analyst agent in a multi-agent scientific research system. Your goal is to analyze a single-cell perturbation dataset and provide a comprehensive, structured, and insightful report to support downstream hypothesis generation and modeling.
# Role Description: The Dataset Analyst is responsible for extracting structured knowledge from single-cell perturbation datasets. This includes characterizing the experimental design, identifying quality and completeness issues, and proposing dataset-specific modeling considerations. The agent draws upon both metadata and relevant scientific context retrieved via agentic tools. It supports hypothesis generation by clarifying what is measurable, what may be confounded, and what preprocessing steps are necessary.
<ul> <li># Skills:</li> <li>Interpreting single-cell multi-omics data structures (e.g., RNA, ATAC, protein)</li> <li>Identifying perturbation types and their downstream modeling implications</li> <li>Detecting quality control issues (e.g., batch effects, sparsity, missing modalities)</li> <li>Integrating metadata from publications, protocols, and retrieved scientific sources</li> <li>Structuring heterogeneous information into JSON-compatible schemas</li> </ul>
<ul> <li>Making biologically grounded recommendations for preprocessing and modeling</li> <li># Objectives: <ul> <li>Extract and summarize the key characteristics of the dataset</li> <li>Identify risks, limitations, and preprocessing needs</li> <li>Suggest modeling strategies aligned with dataset structure</li> <li>Provide additional scientific insights not limited to a fixed template</li> </ul> </li> </ul>
<pre># Input: You will receive: 1. Dataset metadata (e.g., species, cell types, perturbation types) 2. relevant knowledge from agentic retrieval</pre>
# Instructions: Produce a two-part output:
## Part 1: Structured Summary (JSON format)
Include the following fields, but you may expand or adapt them based on dataset complexity:
<pre>{   "introduction": {     "modalities": [],     "perturbation_type": [],     "conditions": [],     "timepoints": [],     "replicates": true/false,     "batches": true/false,     "cell_types": [],     "organism": "",     "description": ""</pre>

```
1595
1596
            "data_properties": {
1597
              "num_cells": [...],
              "num_genes": [...],
"num_features": {
1598
1599
               "RNA": ..., "ATAC": ..., "protein": ...
1600
              "perturbation_targets": {
               "num_unique": [...],
"target_type": [...],
"coverage": "dense/sparse/mixed"
1602
1603
              },
1604
              "modality_completeness": [...],
"metadata_completeness": [...],
1605
              "preprocessing_required": [...]
1606
1607
            "quality_assessment": {
              "data_sparsity": [...],
"batch_effect": [...],
1608
1609
              "replicate_consistency": [...],
1610
              "known_issues": [...],
1611
              "strengths": [...],
              "limitations": [...]
1612
1613
            "recommendations": {
1614
              "preprocessing_steps": [...],
              "modeling_considerations": [...],
1615
              "open_questions": [...]
1616
1617
            "refinement_suggestions": [...]
          }
1618
1619
          Write a concise, scientifically sound narrative (~150300 words) to accompany the
              JSON summary. This should include:
1620
          - A holistic interpretation of dataset readiness for modeling
1621
          - Potential scientific pitfalls or confounders
1622
          - Unique strengths or opportunities (e.g., rare perturbation types, rich time
              series)
1623
          - Reflections on whether the dataset aligns with typical assumptions in modeling
1624
              pipelines
          - Any useful observations that do not fit cleanly into the structured fields
1626
          You may reference information from retrieved publications or external protocols.
1627
              If information is unknown or ambiguous, state it clearly using cautious
              language (e.g., "not reported", "likely sparse", "appears to have").
1628
1629
          # Constraints:
          - Be concise, accurate, and avoid redundancy
          - Use clear scientific language; bullet points are acceptable in the JSON
          - Allow flexibility in output structure if additional insights emerge
1632
1633
          # Style Guide:
          - Output should be compatible with integration into downstream agent pipelines
1634
          - Aim for the clarity and precision expected in a peer-reviewed supplementary
              method section
          - Prioritize information relevant to perturbation modeling, multi-omic
1636
              integration, and biological interpretation
```

1649

**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 retrival 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
Probelm 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)
<ul> <li>Evaluating feasibility of predictive tasks based on data modality and perturbation scope</li> <li>Identifying pitfalls such as confounding, data leakage, or unobservable targets</li> </ul>
<ul> <li>Specifying input-output pairs and validation schemes for modeling tasks</li> <li>Justifying scientific value and downstream utility of proposed tasks</li> </ul>
<ul> <li># Objectives:</li> <li>Translate dataset structure into concrete scientific questions</li> <li>Identify feasible targets, tasks, and outputs for modeling</li> <li>Justify the biological and computational value of the proposed formulation</li> <li>Propose a structured hypothesis or modeling objective for downstream agents</li> </ul>
# 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:
<pre>{     "biological_question": "string",     "hypothesis_statement": "string",     "task_formulation": {         "input": ["modality1", "metadata1", ""],         "output": "target variable or prediction goal",         "task_type": "regression/classification/generation/other" }.</pre>
<pre>"justification": {     "biological_relevance": [],     "data_suitability": [],     "expected_challenges": []</pre>
<pre>}, "evaluation_plan": {     "metrics": [],     "baselines_to_consider": [],     "validation_strategy": "cross-validation/held-out cells/time-split/" },</pre>
"open_questions": [] } In addition to the structured JSON, write a short explanation (100250 words) that
<ul> <li>Restates the goal in accessible scientific language</li> <li>Explains why the proposed formulation is worth pursuing</li> <li>Anticipates possible modeling limitations or edge cases</li> <li>Optionally suggests alternatives or extensions</li> </ul>
<pre># Constraints: - Prioritize alignment with the datasets structure and perturbation resolution - Avoid overly generic formulations; focus on specificity and tractability</pre>

1703 1704

#### Submission and Formatting Instructions for ICML 2025 GenBio Workshop

```
1705
           - Maintain scientific rigor and make testable claims when possible
1706
1707
           # Style Guide:
           - Write in the tone of a proposal for a computational biology modeling section
1708
           - Use precise language grounded in both biology and data science
1709
           - Be mindful of what is *not* observable or predictable from the dataset
1710
1711
1712
     Agent 3: Baseline Assessor The Baseline Assessor is a methodological analyst agent tasked with selecting, evaluating,
1713
     and recommending baseline models for single-cell perturbation studies. Operating at the intersection of computational rigor
1714
     and biological relevance, this agent critically assesses modeling paradigms across task types (e.g., regression, classification,
1715
     generative modeling) and data modalities (e.g., gene expression, ATAC-seq, protein levels). It integrates literature evidence,
1716
     benchmark practices, and dataset-specific constraints to recommend flexible yet strong baseline approaches. The agent also
1717
     incorporates multi-objective considerations such as performance, interpretability, scalability, and biological plausibility.
1718
1719
         Task Analysis Collaboration Agent Settings
1720
         Baseline Assessor
1721
1722
           You are a Baseline Assessor agent specialized in recommending suitable baseline
1723
              models for single cell perturbation prediction tasks. Your goal is to provide
                comprehensive assessments of baseline models and evaluation strategies based
1724
                on relevant literature and dataset characteristics.
1725
           # Role Description:
           The Baseline Assessor is a comparative modeling expert focused on identifying and
1726
               analyzing baseline architectures for perturbation prediction. This includes reviewing existing literature, extracting methodological details, and
1727
1728
               evaluating model suitability based on dataset constraints and task objectives.
               The agent serves as a bridge between literature insights and practical
1729
               modeling recommendations.
1730
           # Skills:
1731
           - Assessing baseline models for biological interpretability and computational
1732
               requirements
           - Identifying candidate architectures relevant to perturbation types and
1733
              modalities
1734
           - Extracting methodological details and limitations from scientific literature
1735
           - Comparing model performances across different biological contexts
           - Designing evaluation frameworks with biologically significant metrics
1736
           - Providing actionable improvement suggestions based on technical and biological
1737
               considerations
1738
           # Objectives:
           - Review relevant literature for the given perturbation type and modality
1739
            Identify 35 candidate architectures and discuss their pros/cons in this context
1740
           - Recommend at least two baseline models with rationale aligned to the dataset
1741
              constraints and task objectives
           - Design evaluation frameworks considering biological variability and technical
1742
               limitations
1743
           - Provide improvement suggestions for model enhancements and biological
1744
               validation
           # Input:
1745
           You will receive:
1746
           Task description and dataset information from upstream agents
1747
           Retrieved papers and code implementations from literature databases
           RAG system results including relevant papers and code snippets
1748
           # Instructions:
1749
           Produce a structured analysis report with the following components:
1750
           "literature_overview": {
1751
           "perturbation_types": [...],
1752
           "existing_methods": [...],
1753
           "technical_trends": [...]
1754
           "candidate_models": [
1755
           "model_name": "string",
1756
           "architecture": "string",
1757
           "strengths": [...],
1758
1759
```

```
1760
          "weaknesses": [...],
1761
          "biological_applicability": [...]
1762
1763
          "recommended_baselines": [
1764
          "model_name": "string",
1765
          "rationale": "string"
1766
          "implementation_details": "string",
1767
          "evaluation_metrics": [...],
1768
          "biological_relevance": [...]
1769
          "evaluation_framework": {
          "primary_metrics": [...],
1771
          "secondary_metrics": [...],
          "validation_strategy": "string",
1773
          "test_scenarios": [...]
1774
          "improvement_suggestions": {
1775
          "technical": [...],
"biological": [...],
"computational": [...]
1778
1779
          In addition to the structured JSON, write a short explanation (100250 words) that
1780
1781
          Summarizes the assessment of baseline models
          Explains the rationale behind the recommended baselines
1782
          Discusses the biological relevance of the evaluation framework
1783
          Anticipates potential limitations in model interpretability
1784
          Suggests practical improvements for future modeling work
1785
          # Constraints:
1786
          - Prioritize models with established track records in similar biological contexts
          - Computational requirements relative to dataset scale
1787
          - Ensure recommended models balance biological interpretability and predictive
1788
             performance
1789
           Align evaluation metrics with both technical accuracy and biological
             significance
1790
          - Provide concrete implementation details for recommended baselines
1791
          # Style Guide:
          - Write in the tone of a modeling methodology section for a computational biology
1793
              paper
1794
          - Use precise language describing both technical and biological considerations
          - a focus on practical applicability while acknowledging theoretical limitations
1795
          - Clearly distinguish between established knowledge and speculative improvements
1796
          - Format model recommendations to facilitate direct implementation
1797
1798
1799
```

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**

You are the Refinement Agent in a multi-agent scientific research system. Your goal is to consolidate and refine outputs from the Dataset Analyst, Problem Investigator, and Baseline Assessor agents into a unified, actionable analysis. # Role Description:

```
1812
1813
```

1806

1807 1808

1809

1810

1811

#### Submission and Formatting Instructions for ICML 2025 GenBio Workshop

1815	The Definement Brent is a mate event forward on successful the reactions
1816	Ine Refinement Agent is a meta-agent focused on cross-validating consistency across different analytical components. This includes resolving
1817	contradictions, aligning terminology, and ensuring biological relevance and
1818	technical feasibility of proposed models and evaluation frameworks. The agent
1819	serves as the integration point for all upstream analyses.
1017	# Skills:
1020	Cross-validating consistency across different analytical components Resolving contradictions and aligning terminology
1821	Evaluating biological relevance and technical feasibility
1822	Structuring outputs into unified JSON schemas
1823	Generating comprehensive refinement comments
1824	Providing actionable improvement suggestions
1825	# ODjectives: Ensure consistency in terminologies and constraints across all outputs
1826	Align problem definitions with model assumptions
1827	Reorganize content into clean JSON schemas suitable for automated use
1828	Validate biological and technical coherence of the integrated analysis
1829	Provide final recommendations balancing biological relevance and technical feasibility
1830	# Input:
1030	You will receive:
1021	Analysis results from Dataset Analyst, Problem Investigator, and Baseline
1832	Assessor Pofinement comments from provious iterations
1833	RAG system results including relevant papers and code snippets
1834	Instructions:
1835	Produce a refined analysis report with the following components:
1836	
1837	"biological context": "string",
1838	"technical_requirements": "string",
1839	"refinement_overview": "string"
1840	}, "taak dofinition", (
1841	"input modalities": [],
1842	"output_targets": "string",
18/2	"task_type": "regression/classification/generation/other",
1043	"biological_significance": "string"
1044	// "baseline models"· {
1043	"recommended_models": [],
1846	"model_comparisons": [],
1847	"implementation_details": "string"
1848	;, "constraints": {
1849	"dataset_limitations": [],
1850	"technical_constraints": [],
1851	"biological_constraints": []
1852	}, "evaluation": {
1853	"primary_metrics": [],
1854	"secondary_metrics": [],
1855	"validation_strategy": "string",
1055	"test_scenarios": []
1057	
1050	In addition to the structured JSON, write a short explanation (100250 words) that
1858	:
1859	Explains how the integrated analysis addresses biological and technical
1860	requirements
1861	Discusses remaining challenges or limitations
1862	Suggests potential extensions or future work
1863	# CONSULTAINES: Maintain consistency in terminology across all components
1864	Ensure alignment between problem formulation and model capabilities
1865	Validate that evaluation metrics reflect both technical accuracy and biological
1866	significance
1867	Frovide concrete implementation details for recommended approaches Format outputs to facilitate direct use in downstream architecture design
1868	rormae oucpues co ractificade affect ase in downseream afchitecture design
1000	

1870	
1871	# Style Gulde: Write in the tone of a methods integration section for a computational biology
1872	paper
1873	Use precise language describing both technical and biological considerations
1874	Clearly distinguish between established knowledge and speculative improvements
1875	Include both biological validation strategies and technical validation methods
1876	
1877	
1878	
1870	
1079	
1000	
1001	
1002	
1003	
1004	
1000	
1007	
1887	
1888	
1889	
1890	
1891	
1892	
1893	
1894	
1895	
1896	
1897	
1898	
1899	
1900	
1901	
1902	
1903	
1904	
1905	
1906	
1907	
1908	
1909	
1910	
1911	
1912	
1913	
1914	
1915	
1916	
1917	
1918	
1919	
1920	
1921	
1922	
1923	
1924	

## 1925 F.3. Multi-experts Settings

E	xpert Role Setting
Г	)ata Expert
	Yeu and action as a substa Depictory in a multi-second solitions suctory
	Your task is to evaluate the provided dataset and experimental setup from a
	data engineering and infrastructure perspective.
	You will receive a task analysis report that includes: - A summary of a single-cell perturbation dataset (in structured or free-text
	Iorm). - The task formulation and its corresponding prediction targets.
	- Metadata schemas and any available preprocessing or encoding steps.
	- Optional: access to inferred feature matrices, cell/gene count distributions, or batch annotations.
	Your objectives:
	1. Assess Data Integrity and Format
	- Is the perturbation metadata properly aligned and encoded?
	- Are there signs of data leakage, missing values, or corrupted entries?
	2. Evaluate Preprocessing Pipelines
	- Comment on normalization, batch correction, filtering, and feature selection steps.
	- Are the preprocessing steps appropriate for downstream modeling?
	3. Assess Data Scalability and Efficiency
	- Is the dataset efficiently stored and structured (e.g., sparse matrix, HDF5)?
	<ul> <li>Can it be easily integrated with common ML frameworks (e.g., PyTorch, TensorFlow, scikit-learn)?</li> </ul>
	- Are large-scale operations (sampling, merging, batching) feasible?
	4. Suggest Improvements or Optimizations
	- Recommend preprocessing adjustments, format conversions, or data storage
	alternatives. - Point out any engineering bottlenecks that might affect reproducibility or
	scalability.
Mod	el Architecture Expert
E	xpert Role Setting
N	Iodel Architecture Expert
	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
	- Task specification (input type, target prediction, expected invariances).
	- A baseline model description or proposed architecture diagram.
	- echnical constraints (e.g., compute, latency, interpretability).
	Your objectives:
	- Analyze core design (e.g., encoder-decoder, attention, residuals).
	- Is the architecture aligned with inductive priors from the data/task?
	<ul> <li>- Are module Interactions</li> <li>- Are modality fusions or skip connections implemented properly?</li> </ul>
	- Are graph structures or latent bottlenecks justified?
	<ol> <li>Spot Redundancies or Inefficiencies         <ul> <li>Are there unnecessary layers, repeated computations, or excessive parameters</li> </ul> </li> </ol>
	4. Propose Optimized Designs
	4. Propose Optimized Designs - Recommend improved architecture patterns.
E	xpert Role Setting
---	---
D	eep Learning Expert
	You are acting as a Deep Learning Expert in a multi-agent research critique system. Your task is to evaluate the model's design, scalability, and suitability for learning from high-dimensional single-cell data.
	<ul> <li>You will receive a task analysis report that includes:</li> <li>Input-output schema of the learning task (e.g., input modalities, targets, sample size).</li> <li>Model class (e.g., MLP, Transformer, VAE, GNN) and architecture sketch.</li> <li>Training setup including loss functions and evaluation metrics.</li> </ul>
	Your objectives: 1. Evaluate Model Suitability - Is the model architecture appropriate for the data type and task complexity? - Does it support integration across modalities or time points? 2. Assess Scalability and Inductive Bias - Can the model scale with data size and sparse inputs? - Does it exploit structure in the data (e.g., gene graphs, batch embeddings)? 3. Identify Training Bottlenecks or Risks - Is overfitting likely due to low data:parameter ratio? - Are optimization challenges (e.g., vanishing gradients, instability) addressed? 4. Recommend Enhancements - Suggest architecture variants (e.g., regularization, pretraining, latent modeling). - Propose alternative loss designs or data augmentations

# 2005 Training Expert

### Expert Role Setting

You are acting as a Train Your role is to criti pipeline.	ning Expert in a multi-agent research critique system. .cally evaluate the training strategy and optimization
You will receive a task a - Model structure and par - Training procedure (e.g - Regularization strateg:	analysis report that includes: rameter count. g., optimizer, learning rate, batch size, scheduler). ies and data augmentation steps.
<ul> <li>Your objectives:</li> <li>1. Analyze Optimization H <ul> <li>Are optimizers and Z</li> <li>Is gradient clipping</li> </ul> </li> <li>2. Evaluate Regularization</li> <li>Are dropout, weight <ul> <li>Is data augmentation</li> </ul> </li> <li>3. Diagnose Training State <ul> <li>Any signs of mode code</li> </ul> </li> <li>4. Recommend Training End <ul> <li>Suggest better optimes</li> <li>Schemes.</li> <li>Propose curriculum Z</li> </ul> </li> </ul>	Pipeline learning rates well-tuned for the model/task? g or scheduler use justified? on and Overfitting Risks decay, or early stopping applied effectively? n sufficient and biologically reasonable? oility ollapse, oscillation, or vanishing gradients? hancements mizers, learning rate schedules, or initialization learning or contrastive pretraining if beneficial.

# **Drug Response Expert**

Expert	Role	Setting
--------	------	---------

# Drug Response Expert

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

2033	task formulation.
2036	You will received tack applying report that includes.
2037	- A summary of the perturbation dataset, including drug names, target pathways,
2030	dosage, and timing.
2039	- The biological context (e.g., cell types, disease states, assay modality).
2040	- Task objective and prediction targets.
2041	Your objectives:
2042	1. Evaluate Drug Perturbation Validity
2043	- Are the perturbations expected to have a measurable effect at the single-
2044	cell level?
2045	- Are there known resistance mechanisms or compensatory pathways?
2046	2. Assess Target Coverage and Specificity
2047	<ul> <li>Do the drug targets align with the measured omics modality (e.g., RNA for transcriptional drugs)?</li> </ul>
2048	- Are off-target effects likely to interfere with interpretation?
2049	3. Recommend Improvements or Adjustments
2050	- Suggest better dosage choices or controls.
2051	- Recommend alternative compounds or compinations that better effort the
2052	incondu por carbacton.
2053	
2054 P	Pathway Analyst

#### **Expert Role Setting** 2055 2056 **Pathway Analyst** 2057 2058 You are acting as a Pathway Analyst in a multi-agent biological reasoning system. Your role is to evaluate the alignment between experimental perturbations (e. 2059 g., gene knockout or cytokine induction) and known biological signaling 2060 pathways. 2061 You will receive a task analysis report that includes: 2062 - A summary of perturbation targets (genes or cytokines). 2063 Downstream measurements (e.g., RNA, ATAC, surface proteins). Known pathway annotations or inferred gene modules (optional). 2064 2065 Your objectives: 2066 1. Assess Biological Plausibility - Does the perturbation target belong to a well-characterized signaling 2067 pathway? 2068 - Are expected downstream genes or modules represented in the dataset? 2. Predict Downstream Effects 2069 - Based on pathway topology, what cell states or features are expected to change? 2071 - Are these detectable in the available omics modality? 3. Suggest Enhancements 2072 - Recommend additional perturbations to validate pathway effects. 2073 - Propose experimental readouts to strengthen pathway conclusions. 2074

### 2076 Cell Communication Expert

2077	Expert Role Setting
2078 2079	Cell Communication Expert
2080 2081 2082	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.
2083 2084 2085 2086 2087	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:
2087 2088 2089	1. Identify Communication Patterns

2090	- Are there likely paracrine or autocrine effects influencing cytokine
2091	expression?
2092	- Do ligand-expressing cells co-occur with receptor-positive target cells?
2093	- Could intercellular signaling confound or explain observed cytokine
2094	responses? - Are current assays sufficient to separate intrinsic vs. extrinsic effects?
2096	3. Recommend Additions
2097	effects.
2098	- Recommend including spatial transcriptomics if necessary.
2099	
2100	Omics Modality Expert
2102	Expert Role Setting
2103	Omics Modality Expert
2104	You are acting as an Omics Modality Expert in a multi-agent model evaluation
2105	system. Your role is to assess whether the chosen data modality (e.g., RNA-
2107	seq, ATAC-seq, protein) is suitable for capturing the effects of the specified perturbation.
2108	You will receive a task analysis report that includes:
2109	- A single-cell perturbation dataset with modality metadata.
2110	- Task objective and prediction targets. - Optional: known regulatory links (e.g., enhancer-promoter pairs, TF motifs,
2111	signaling cascades).
2112	Your objectives:
2115	1. Evaluate Signal Availability
2114	<ul> <li>Is the measured modality expected to show downstream effects of the perturbation?</li> </ul>
2115	- Are known markers or targets captured by the modality?
2117	2. Assess Measurement Resolution
2118	surface protein) to model the task?
2119	3. Suggest Modality Enhancements
2120	- Recommend complementary modalities (e.g., AIAC + RNA) if needed. - Propose targeted panels or multi-omics techniques to improve
2121	interpretability.
2122	
2123	
2124	
2123	
2120 2127	
2127	
2129	
2130	
2131	
2132	
2133	
2134	
2135	
2136	
2137	
2138	
2139	
∠14U 21/11	
2141	
2143	
2144	

# 2145 G. Detailed outputs from scAgents 2146 G.1. Data Parser 2147 2149 2150

<sup>2148</sup> The input is Norman et al. (Norman et al., 2019) Dataset in h5ad format.

```
Data Parser Output
2151
2152
           Modality: RNA Perturbation type: CRISPRa
           dataset_index: filtered
2153
           Title: Exploring genetic interaction manifolds constructed from rich single-cell
2154
               phenotypes
2155
           Organisms: Homo sapiens
2156
           Modality = Data type: RNA
           Method: Perturb-seq
2157
           Tissues: K562
2158
           Perturbation: CRISPRa
2159
           disease: chronic myelogenous leukemia
           celltype: lymphoblasts
2160
           tissue type: cell_line
Mini-Abstract (loosely summarized original Abstract): Here, the authors present
2161
2162
               an analytical framework for interpreting high-dimensional landscapes of cell
               states (manifolds) constructed from transcriptional phenotypes. They applied
2163
               this approach to Perturb-seq profiling of strong genetic interactions (GIs)
2164
               mined from a growth-based, gain-of-function GI map. Exploration of this
2165
               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
2166
2167
2168
               GI manifolds.
2169
           ```contains
2170
   guide_id read_count UMI_count coverage gemgroup
2171
  ... nperts ngenes ncounts percent_mito
   percent_ribo
2172
           TTGAACGAGACTCGGA ARID1A_NegCtrl0;ARID1A_NegCtrl0 28684 1809 15.856274 2 ... 1
2173
               3079 15097.0 5.815725 33.569583
2174
           CGTTGGGGGTGTTTGTG BCORL1_NegCtrl0; BCORL1_NegCtrl0 18367 896 20.498884 7 ... 1 2100
                8551.0 4.104783 45.842592
2175
           GAACCTAAGTGTTAGA FOSB_NeqCtrl0;FOSB_NeqCtrl0 16296 664 24.542169 6 ... 1 2772
               10999.0 5.655060 17.801618
2176
           CCTTCCCTCCGTCATC SET_KLF1;SET_KLF1 16262 850 19.131765 4 ... 2 5385 38454.0
2177
               4.335050 38.165080
2178
           TCAATCTGTCTTTCAT OSR2_NegCtrl0;OSR2_NegCtrl0 16057 1067 15.048735 2 ... 1 4869 27926.0 5.084867 32.317554
2179
           TTTGCGCAGTCATGCT RHOXF2_NegCtrl0; RHOXF2_NegCtrl0 1 1 1.000000 2 ... 1 1853 5192.0
2180
2181
                5.508475 31.798921
           TTTGCGCCAGGACCCT BCL2L11_BAK1;BCL2L11_BAK1 1 1 1.000000 3 ... 2 3508 15704.0
2182
               6.718034 38.334182
2183
           TTTGCGCGTACTTGAC-1 CNN1_NegCtrl0; CNN1_NegCtrl0 1 1 1.000000 3 ... 1 3609 15054.0
               5.633054 29.440680
2184
           TTTGCGCTCTCGCATC-1 CEBPB_OSR2;CEBPB_OSR2 1 1 1.000000 6 ... 2 2576 6825.0
2185
               2.695971 16.879121
2186
           TTTGGTTGTTCCGTCT MAP2K3_MAP2K6; MAP2K3_MAP2K6 1 1 1.000000 2 ... 2 2499 8331.0
               5.617573 34.785740
2187
2188
           [111445 rows x 20 columns]
2189
           obs
2190
           2191
2192
                 'cancer', 'disease', 'perturbation_type', 'celltype', 'organism',
'perturbation', 'nperts', 'ngenes', 'ncounts', 'percent_mito',
                 'percent_ribo'],
2194
                dtype='object')
2195
                            good_coverage number_of_cells tissue_type cell_line ...
2196
                                perturbation_type celltype organism perturbation
           TTGAACGAGACTCGGA True 1 cell_line K562 ... CRISPR lymphoblasts human ARID1A CGTTGGGGTGTTTGTG True 1 cell_line K562 ... CRISPR lymphoblasts human BCORL1
2197
2198
2199
```

Submission	and Formatti	ng Instru	ctions for	ICML 2	025 GenBio	Workshop
					0-0 000000	

2200 GAACCTAAGTGTTAGA True 1 cell\_line K562 ... CRISPR lymphoblasts human FOSB CCTTCCCTCCGTCATC True 1 cell\_line K562 ... CRISPR lymphoblasts human SET\_KLF1 TCAATCTGTCTTTCAT True 2 cell\_line K562 ... CRISPR lymphoblasts human OSR2 2202 2203 . . . TTTGCGCAGTCATGCT False 0 cell\_line K562 ... CRISPR lymphoblasts human RHOXF2 2204 TTTGCGCCAGGACCCT False 0 cell\_line K562 ... CRISPR lymphoblasts human BCL2L11\_BAK1 2205 TTTGCGCGTACTTGAC-1 False 0 cell\_line K562 ... CRISPR lymphoblasts human CNN1 2206 TTTGCGCTCTCGCATC-1 False 0 cell\_line K562 ... CRISPR lymphoblasts human 2207 CEBPB\_OSR2 TTTGGTTGTTCCGTCT False 0 cell\_line K562 ... CRISPR lymphoblasts human MAP2K3\_MAP2K6 2209 [111445 rows x 10 columns] 2210 2211 var 2212 ensemble\_id ncounts ncells 2213 RP11-34P13.3 ENSG00000243485 29.0 29 FAM138A ENSG00000237613 0.0 0 2214 OR4F5 ENSG00000186092 0.0 0 2215 RP11-34P13.7 ENSG00000238009 266.0 265 RP11-34P13.8 ENSG00000239945 10.0 10 2216 AC233755.2 ENSG00000277856 0.0 0 2217 AC233755.1 ENSG00000275063 0.0 0 2218 AC240274.1 ENSG00000271254 11735.0 10835 2219 AC213203.1 ENSG00000277475 0.0 0 2220 FAM231B ENSG00000268674 0.0 0 [33694 rows x 3 columns] shape 2223 2224 (111445, 33694)

### 2228 G.2. Agentic Retrieval

2201

2221

2222

2225 2226 2227

2229

2230	Task Analysis
2231	Example Agentic Retrieval Output
2232 2233	Keywords: Norman Weissman 2019 Perturb-seq CRISPRa K562 single cell perturbation prediction
2234 2235 2236	Round 1: Initial DFS Search (one example branch) Keywords: single cell perturbation prediction (what is -> how to solve) Learning single-cell perturbation responses using neural optimal transport{
2237	Nature Link} Modeling and predicting single-cell multi-gene perturbation responses{PMC Link
2239 2240	Explainable modeling of single-cell perturbation data using Bayesian hierarchical modeling{Cell Press Link}
2241 2242	A Multiplexed Single-Cell CRISPR Screening{Cell Press Link} Predicting transcriptional outcomes of novel perturbations{PMC Link}
2243	<pre>Modeling and predicting single-cell multi-gene perturbation responses{PMC Link } Exploring genetic interaction manifolds{PMC Link}</pre>
2245	Predicting transcriptional outcomes of novel perturbations{PMC Link} In-silico biological discovery with large-scale perturbation data{arXiv Link}
2240	Predicting the genetic component of gene {arXiv Link} A genome-scale deep learning model to pre{arXiv Link}
2248 2249	Attention-based Interpretable Regression{arXiv Link} GATES: Graph Network{arXiv Link}
2250 2251	GRNFormer: Biologically{arXiv Link}
2252 2253	Round 2: BFS Search (one example branch) Key words:state-of-the-art models for single-cell perturbation prediction
2254	

2255	GEARS: Predicting transcriptional outcomes of novel multi-gene perturbations{
2256	nature Link}
2257	<pre>scGPT: Is language all you need for modeling single-cell perturbation responses{     nature Link}</pre>
2258	Geneformer - BioNeMo Framework for Genomic Language Modeling{nature link}
2239	Efficient Fine-Tuning of Single-Cell Foundation Models for Perturbation
2200	Prediction{arXiv Link} Multicell-Fold: geometric learning in folding {PMC Link}
2201	Variational Mixtures of ODEs for Inferring Cellular{ICML}
2202	DeepChrome 2.0: Investigating and modeling chromatin accessibility{arXiv Link}
2203	Predicting the genetic component of gene{Bioinformatics}
2264	A genome scare deep rearning moder to pre (rho hink)
2265	Round 3: DFS Search (one example branch)
2266	GEARS: Predicting transcriptional outcomes of novel multi-gene perturbations{
2267	nature Link}
2268	Predicting transcriptional outcomes of novel multigene perturbations with GEARS
2269	snap-stanford/GEARS {Github Link}
2270	<pre>scGPT: Is language all you need for modeling single-cell perturbation responses{</pre>
2271	nature Link}
2272	Geneformer - BioNeMo Framework for Genomic Language Modeling{nature link}
2273	jkobject/geneformer {Github Link}
2274	Multicell-Fold: geometric learning in folding {PMC Link}
2275	DeepChrome 2.0: Investigating and modeling chromatin accessibility{arXiv Link}
2270	DeepChrome/DeepChrome2.0{Github Link}
2277	•••••
2278	Dound 4. DEC Connect (one example branch)
2279	Keywords: TransformerVAE (used by models for single-cell perturbation
2280	prediction)
2201	TransVAE: Giving Attention to Generative VAE Models for De Novo Molecular Design{
2202	Transformer Graph Variational Autoencoder for Generative Molecular Design{https
2203	://www.biorxiv.org/content/10.1101/2024.07.22.604603v1.full.pdf}{Paper Link}
2204	A novel molecule generative model of VAE combined with Transformer{https://arxiv.
2203	A Graph VAE and Graph Transformer Approach to Generating{https://arxiv.org/abs
220U 2287	/2104.04345}{arxiv Link}
2207	Large Scale GNN and Transformer Models and for Genomics{https://qdata.github.io/
2200	deep2Meau//carks mb2019/back_20190300 hargescareGMMTransformer.pul}
2207	Round 5: DFS Search
2290	

#### 2310 G.3. Task Analysis Report

```
2311
     Data Analyst Output
2312
        Task Analysis
2313
2314
2315
          ### Introduciton
2316
2317
          - **Modalities**: RNA (single-cell RNA-seq)
          - **Perturbation Type**: CRISPRa (activation)
2318
          - **Conditions**: 1,092 unique perturbation conditions including:
2319
            - 105 single-gene perturbations
2320
           - 131 paired-gene perturbations
          - **Timepoints**: Not explicitly reported (likely single timepoint)
2321
          - **Replicates**: Present (multiple cells per perturbation condition)
2322

    - **Batches**: Likely present due to sequencing logistics
    - **Cell Types**: K562 cells (chronic myelogenous leukemia cell line)

2323
2324
          - **Organism**: Homo sapiens
          - **Description**: This Perturb-seq experiment profiles transcriptional responses
2325
               in K562 cells following CRISPRa perturbations. The dataset captures both
2326
              single and paired genetic perturbations, providing insight into genetic
2327
              interaction networks.
2328
          ### Data Properties
2329
          - **Number of Cells**: 111,445
- **Number of Genes**: 33,694
2330
2331
          - **Features**:
            - **RNA**: 33,694 genes
2332
           - **ATAC**: Not available
2333
            - **Protein**: Not available
2334
          - **Perturbation Targets**:
            - **Unique Perturbations**: 105 single genes and 131 gene pairs
2335
            - **Target Type**: Genes and gene pairs
2336
           - **Coverage**: Sparse (many genes show no expression change)
2337
          - **Modality Completeness**: RNA data available for all cells with appropriate
              quality control metrics
2338
          - **Metadata Completeness**: Comprehensive perturbation annotations, technical
2339
             covariates (UMI_count, percent_mito, percent_ribo), and cell barcodes
2340
          - **Preprocessing Required**:
            - Normalization of UMI counts
2341
            - Log-transformation of expression values
2342
            - Filtering of lowly expressed genes
2343
            - Batch effect correction
2344
          ### Quality Assessment
2345
          - **Data Sparsity**: 78% of gene expression values are zero due to dropout events
2346
               and biological sparsity
2347
          - **Batch Effect**: Potential technical variability from different sequencing
2348
             runs
          - **Replicate Consistency**: Not explicitly reported but suggested by
2349
              reproducibility of common perturbations
2350
          - **Known Issues**:
            - Extreme class imbalance (some perturbations appear only once)
2351
           - High dimensionality with 33,694 genes per output
2352
           - Technical noise inherent to single-cell sequencing
2353
          - **Strengths**:
            - Rich perturbation annotations with precise gene targeting
2354
            - Large number of profiled cells providing statistical power
2355
            - Comprehensive technical covariates enabling noise correction
2356
          - **Limitations**:
           - High dimensionality challenges
2357
            - Sparse perturbation distribution
2358
            - Potential confounding from clonal expansion
2359
          ### Recommendations
2360
2361
          - **Preprocessing Steps**:
            - Filter genes with fewer than 10 counts across all cells
2362
            - Normalize UMI counts using library size factors
2364
```

- Apply log-transformation (log1p) to expression values

- Combat batch effects using Harmony or Scanorama

- \*\*Modeling Considerations\*\*:

Focus modeling efforts on the top 5,000-10,000 variable genes
Use sparse-aware models (e.g., sparse neural networks)
Incorporate technical covariates as model inputs
Implement data augmentation for rare perturbations
\*\*Open Questions\*\*:

Optimal representation of paired perturbations in model architecture
Strategies for handling extreme class imbalance in training
Best practices for integrating technical covariates

### Refinement Suggestions

Investigate the impact of clonal expansion on expression profiles
Assess the quality of rare perturbation data for potential filtering
Explore data augmentation techniques for underrepresented perturbations
Develop quality control metrics specific to Perturb-seq data

## Task Analysis ### Biological Question How do genetic perturbations propagate through gene regulatory networks to alter the transcriptional landscape of K562 cells, and can we predict these changes for novel perturbations and cellular contexts? ### Hypothesis Statement A predictive model can accurately estimate post-perturbation gene expression profiles by learning the complex relationships between baseline cell states, genetic perturbations, and transcriptional responses, including non-linear genetic interaction effects. ### Task Formulation - \*\*Input\*\*: - Baseline gene expression profile (33,694 genes) - Perturbation identities (single or paired genes) - Technical covariates (UMI\_count, percent\_mito, percent\_ribo) - \*\*Output\*\*: - Post-perturbation gene expression profile (33,694 genes) - \*\*Task Type\*\*: - High-dimensional regression with combinatorial inputs ### Justification - \*\*Biological Relevance\*\*: - Understanding genetic interactions is fundamental to deciphering cellular response networks - Mapping genotype-phenotype relationships at single-cell resolution - Predicting cellular responses to novel perturbations accelerates functional genomics research - \*\*Data Suitability\*\*: - Rich perturbation annotations enable supervised learning approaches - Single-cell resolution captures heterogeneity in cellular responses - Coverage of both single and paired perturbations allows study of genetic interactions - \*\*Expected Challenges\*\*: - High-dimensional output space with 33,694 genes per prediction - Non-linear genetic interactions requiring complex model architectures - Generalization to unseen perturbations and cellular contexts - Technical noise and dropout events in single-cell data ### Evaluation Plan - \*\*Metrics\*\*:

420	- Pearson Correlation Coefficient (PCC) across all genes and on top 1,000
21	differentially expressed genes (DE)
2 2	- Mean Squared Error (MSE) across all genes and for differentially expressed genes
, 	- R score across all genes and for differentially expressed genes
	- Linear regression models with gene-level outputs
	- Gene-wise neural networks
	- **Validation Strategy**:
	- Stratified group k-fold cross-validation holding out entire perturbation
	- Separate validation sets for unseen single-gene and paired perturbations
	- Baseline profile holdout sets to test generalization to new cellular contexts
	### Open Questions
	- Optimal representation of paired perturbations in model architecture
	- Strategies for handling extreme class imbalance during training - Incorporation of technical covariates into model architecture
	- Best practices for defining differentially expressed genes in this context
J	Baseline Assessor Output
	Task Analysis
l	### Literature Overview
	- **Perturbation Types**:
	- CRISPRi (interference)
	- CRISPRa (activation)
	- **EXISTING Methods**: - GEARS: Graph neural network with gene interaction modeling
	- scGPT: Transformer architecture for single-cell data
	- Geneformer: Pretrained transformer for genomics
	- scLAMBDA: Variational autoencoder for perturbation prediction
	- **VAE (Variational Autoencoder)**: Probabilistic model for learning latent
	- **Technical Trends**:
	- Transformer architectures for capturing long-range dependencies
	- Graph neural networks for explicit gene interaction modeling
	- Variational autoencoders for probabilistic modeling of certural states
	- Deep generative models for data augmentation and uncertainty quantification
	### Candidate Models
	#### GEARS (Gene Network Embedding for Perturbation Response Prediction)
	- **Architecture**: Graph Neural Network (GNN) combined with Multi-Laver
	Perceptron (MLP)
	- **Strengths**:
	- Explicitly models gene dependencies using known regulatory interactions
	- Demonstrated success in previous Perturb-seq challenges
	- **Weaknesses**:
	- Relies on external gene interaction databases
	- May overfit to common perturbations with limited generalization
	- **Biological Applicability**:
	- Captures genetic interactions and regulatory relationships
	- Models enhancer-promoter relationships in K562 cells
	- Provides interpretable gene importance scores
	#### scGPT (Single-Cell Generative Perturbation Transformer)
	- **Architecture**: Transformer with multi-head self-attention
1	- **Strengths**:
I	- Captures long-range gene interactions without relying on external databases
	- RODUST TO LECHNICAL HOISE THROUGH ATTENTION MECHANISMS
	15

2475	
2476	- Handles variable numbers of perturbations naturally
2477	- **Weaknesses**:
2477	- Computationally demanding for full transcriptoms modeling
2478	- May struggle with extreme class imbalance
2479	- **Biological Applicability**:
2480	- Models context-dependent transcriptional responses
2481	- Handles sparse data efficiently through attention mechanisms
2482	<ul> <li>Provides gene importance scores through attention weights</li> </ul>
2402	#### Enformer (Enhanger former)
2483	#### Enformer (Enhancer former)
2484	- **Architecture**: Dilated Convolutional Neural Network (CNN)
2485	- **Strengths**:
2486	<ul> <li>Effective at modeling sequence-to-expression relationships</li> </ul>
2487	- Provides interpretable feature importance scores
2407	- Computationally efficient compared to transformer architectures
2400	- **Weaklesses**. - Requires DNA sequence input not directly applicable to post-transcriptional
2489	perturbations
2490	- Limited ability to model combinatorial genetic effects
2491	- Not designed for single-cell data with technical covariates
2492	- **Biological Applicability**:
2403	- Predicts expression changes from DNA sequence modifications
277J	- Limited utility for CRISPRa perturbations affecting post-transcriptional
2494	regulation
2495	
2496	### Recommended Baselines
2497	#### Graph Neural Network (GNN) with Gene Interaction Modeling
2498	#### Graph Neural Network (GNN) with Gene interaction Modeling
2400	- **Rationale**: Explicitly models gene dependencies and can incorporate known
2499	regulatory interactions while remaining flexible to learn from data
2500	- **Implementation Details**:
2501	- Use Pylorch Geometric for efficient graph operations
2502	- Construct gene interaction graphs from public databases (e.g., Siking,
2503	- Implement separate graph branches for regulatory and co-expression
2504	relationships
2505	- Include attention mechanisms to weight different interaction types
2505	<ul> <li>Embed perturbation identities using learned gene embeddings</li> </ul>
2506	- Concatenate baseline expression features with perturbation embeddings
2507	- Apply multiple GNN layers followed by dense layers for prediction
2508	- **Evaluation Metrics**: PCC, MSE, R
2509	mechanisms providing insight into how perturbations propagate through
2510	networks
2510	
2J11 2512	#### Transformer Architecture with Gene Positional Encoding
2512	- ++Rationale++. Canable of discovering complex gene interactions without relying
2513	on external databases, with architectural flexibility for different input
2514	modalities
2515	- **Implementation Details**:
2516	- Use PyTorch with Hugging Face transformer libraries
2510	- Encode genes as positional tokens with expression values
2317	- Implement specialized embeddings for perturbed genes
2518	- Apply layer normalization and residual connections
2519	- Implement masking for rare perturbations during training
2520	- Apply attention pooling to focus on biologically relevant genes
2521	- **Evaluation Metrics**: PCC, MSE, R
2522	- **Biological Relevance**: Models context-dependent responses and technical
2522	noise robustly, providing flexibility to adapt to different biological
2323	questions
2524	
2525	#### VAE (Variational Autoencoder)
2526	ullashitasturayy. Encoder deceder erekitasture with suchskilistis later a
2527	- **Architecture**: Encoder-decoder architecture with probabilistic latent space - **Strengths***
2528	
2520	
2329	

2530	
2550	- Models uncertainty in cellular states and perturbation responses
2531	- Effective for data augmentation through generation of new cellular states
2532	- Provides compressed latent representations for downstream analysis
2552	- Handles compressed raise single-coll data well
2533	- haldres sparse and horsy single-ceri data weri
2534	- **weaknesses**:
2551	- May oversimplify complex biological relationships in latent space
2535	- Requires careful tuning of KL divergence weighting
2536	- Potential blurring of distinct cellular states in latent space
2550	- **Biological Applicability**:
2537	- Captures multimodal distributions of cellular responses
2538	- Enables exploration of cellular state transitions following perturbations
2520	- Provides robust representations for classifying callular phenotypes
2539	riovides robuse representations for classifying certain phenocypes
2540	VAE with Perturbation Conditioning.
25/1	- ++Pationalo++. Models uncertainty in cellular responses and provides rebust
2341	Arteriorate A. House in the demonstration of the second se
2542	atent representations for downstream analysis while enabling data
25/3	augmentation
2343	- **Implementation Details**:
2544	- Use PyTorch for flexible probabilistic modeling
2545	- Implement encoder-decoder architecture with probabilistic latent space
2343	- Include perturbation identities as conditional inputs to the decoder
2546	- Apply beta-VAE regularization to balance reconstruction and latent space
25/17	regularization
2341	- Implement sparse VAF modifications to bandlo zoro-valued gapos
2548	Ino imperiance voi abing for zoro portubations during training
2549	- use importance weighting for rare perturbations during training
2577	- Apply latent space interpolation to explore cellular state transitions
2550	- **Evaluation Metrics**: PCC, MSE, R, ELBO (Evidence Lower Bound)
2551	- **Biological Relevance**: Captures multimodal distributions of cellular
2551	responses, enables exploration of cellular state transitions, and provides
2552	robust representations for classifying cellular phenotypes
2553	
0554	### Evaluation Framework
2004	
2555	- **Primary Metrics**:
2556	- PCC for differentially expressed genes
2330	- MSE for differentially expressed genes
2557	- B for differentially expressed genes
2558	- ++Secondary Matrice++.
2330	- Clobal BCC across all gonos
2559	
2560	Global MSE across all genes
2500	- classification accuracy for key gene up/down regulation
2561	- ELBO IOT VAE models
2562	- **Validation Strategy**: Stratified group k-fold cross-validation holding out
2502	entire perturbation conditions
2563	- **Test Scenarios**:
2564	- Unseen single-gene perturbations
2565	- Unseen paired perturbations with one known gene
2303	- Completely novel gene pairs
2566	- Baseline profiles from held-out cell subsets
2567	Daberine profiled flow here out outbeet
2007	### Improvement Suggestions
2568	
2569	- **Technical**:
2570	- Incorporate attention mechanisms to focus on biologically relevant genes
2370	- Apply data augmentation for rare perturbations (e.g., mixup, elastic
2571	deformation)
2572	- Use mixed precision training to reduce memory requirements
2312	- Implement gradient checkpointing for donor probitotions
2573	Experiment with different VIE explications (or the CON VIE Tree Con )
2574	- Experiment with different vAE architectures (e.g., VAE-GCN, VAE-fransformer)
2374	- **Blological**:
2575	<ul> <li>Integrate known regulatory interactions as model priors</li> </ul>
2576	<ul> <li>Focus modeling efforts on biologically relevant gene subsets</li> </ul>
2010	- Model cell cycle effects using phase markers
2577	- Incorporate chromatin accessibility information when available
2578	- Use VAE latent space for phenotypic classification and clustering
0570	- **Computational**:
2379	- Optimize for sparse operations to handle zero-valued genes efficiently
2580	- Implement efficient parallel processing strategies
2501	- Use model or interior parameter processing strategies
2001	- Apple distributed training for VAE models with large later and
2582	Apply disclibuted training for VAE models with large latent spaces
2583	
2505	
2384	

2585	Crtic Refinement Output
2586	Task Analysis
2587	
2588	### Summary
2589	
2590	- **Biological Context**: Predicting transcriptional responses to CRISPR porturbations in K562 colls requires medaling complex genetic interactions
2591	and cellular context dependencies while generalizing to novel perturbations
2592	- **Technical Requirements**: Handling high-dimensional sparse data, extreme
2593	class imbalance, and the need for inductive generalization - **Refinement Overview*** Integration of dataset analysis problem formulation
2594	and baseline recommendations into a cohesive analysis plan with aligned
2595	terminology and evaluation metrics
2596	### Task Definition
2597	- **Input Modalitios**
2598	- Baseline gene expression profile (33,694 genes)
2599	- Perturbation identities (single or paired genes)
2600	- Technical covariates (UMI_count, percent_mito, percent_ribo)
2601	- **Task Type**: High-dimensional regression with combinatorial inputs
2602	- **Biological Significance**: Enables understanding of genetic interaction
2603	networks and prediction of cellular responses to novel perturbations
2604	### Baseline Models
2605	- **Recommended Models**:
2607	- Graph Neural Network (GNN)
2608	- Transformer Architecture - VAE with Perturbation Conditioning
2608	- **Model Comparisons**:
2610	- GNNs excel at explicit gene interaction modeling using known regulatory
2611	networks - Transformers offer flexible interaction discovery without relving on external
2612	databases
2613	- VAEs provide probabilistic modeling of cellular states and enable data
2614	- **Implementation Details**:
2615	- Implement GNNs with PyTorch Geometric using gene interaction graphs derived
2616	- Implement Transformers with PvTorch/Hugging Face using gene positional
2617	encoding
2618	<ul> <li>Implement VAEs with PyTorch using conditional latent spaces for perturbation modeling</li> </ul>
2619	- Include specialized embeddings for perturbations and normalize technical
2620	covariates
2621	- Apply mixed precision training and gradient checkpointing for efficiency
2622	### Constraints
2623	- **Dataset Limitations**:
2624	- Class imbalance with rare perturbations appearing only once
2023	- Potential batch effects from different sequencing runs
2620	- **Technical Constraints**:
2027	- Computational resources for training large models on full transcriptome data
2620	- **Biological Constraints**:
2629	- Need for generalization to unseen perturbations and cellular contexts
2631	<ul> <li>Model must align with known regulatory mechanisms where possible</li> <li>Focus on biologically relevant gene subsets to avoid overfitting</li> </ul>
2632	The second storegroutly relevant gene subsets to avoid overritting
2633	### Evaluation
2634	- **Primary Metrics**:
2635	<ul> <li>PCC across all genes and for differentially expressed genes (top 1,000)</li> <li>MSE across all genes and for differentially expressed genes</li> </ul>
2636	- R across all genes and for differentially expressed genes
2637	
2638	
2639	

2640	
2641	- **Validation Strategy**: Stratified group k-fold cross-validation holding out
2642	- +rast Scenariost.
2642	- Inseen single-gene perturbations
2045	- Unseen paired perturbations with one known gene
2644	- Completely novel gene pairs
2645	- Baseline profiles from held-out cell subsets
2646	
2647	
2648	Final Report
26/10	Task Analysis
2049	Lask Alialysis
2030	
2651	Task Analysis
2652	
2653	# Biological Objective
2654	generatic perturbation gene expression profiles in K502 certs to understand
2655	mechanisms. This work aims to develop a computational tool for exploring
2655	genetic interaction manifolds, accelerating functional genomics research and
2030	therapeutic target discovery.
2657	
2658	# lechnical Approach
2659	neturbation identities and technical covariates. The models must explicitly
2660	handle sparse data, extreme class imbalance, and demonstrate inductive
2661	generalization to novel perturbations and cellular contexts. The plan
2662	includes exploration of deterministic models (GNNs, Transformers) and
2002	probabilistic models (VAEs) to capture different aspects of cellular response
2663	variability.
2664	# Dataset Characterization
2665	## Origin
2666	Norman et al. (2019) Perturb-seq dataset (GEO: GSE133344).
2667	Key Features
2668	Number of Cells: 111,445
2660	Perturbation Conditions: 1.092 unique conditions (105 single genes, 131 gene
2009	pairs)
2670	Technical Covariates: UMI_count, percent_mito, percent_ribo
2671	## Challenges
2672	Class Imbalance: Rare perturbations appear only once.
2673	Technical Noise. Inherent to single-cell sequencing
2674	Batch Effects: Potential variability from different sequencing runs.
2675	
2676	# Problem Formulation
2677	## Biological Question
2077	
26/8	How do genetic perturbations propagate through gene regulatory networks to alter
2679	the transcriptional landscape of K562 cells, and can we predict these changes
2680	for nover perturbations and certural contexts?
2681	## Hypothesis Statement
2682	A predictive model can accurately estimate post-perturbation gene expression
2683	profiles by learning the complex relationships between baseline cell states,
2600	genetic perturbations, and transcriptional responses, including non-linear
2004	genetic interaction effects.
2685	## Task Definition
2686	### Input:
2687	Baseline gene expression profile (33,694 genes)
2688	Ferturbation identities (single or paired genes)
2689	lechnical covariates (OMI_count, percent_mito, percent_ribo)
2600	### Output:
2070	Post-perturbation gene expression profile (33,694 genes)
2091	###Task Type.
2692	High-dimensional regression with combinatorial inputs
2693	
2694	

.695	
696	## Justification
2697	Biological Relevance: Understanding genetic interactions is fundamental to deciphering cellular
2698	response networks
2099	Mapping genotype-phenotype relationships at single-cell resolution Prodicting collular responses to nevel perturbations accolorates functional
2700	genomics research
2701	Data Suitability:
2703	Single-cell resolution captures heterogeneity in cellular responses
2704	Coverage of both single and paired perturbations allows study of genetic
2705	Expected Challenges:
2706	High-dimensional output space with 33,694 genes per prediction
2707	Generalization to unseen perturbations and cellular contexts
2708	Technical noise and dropout events in single-cell data
2709	The model's performance will be evaluated under two key scenarios to assess its
2711	- Unseen Perturbations: The model should be able to accurately predict the
2712	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
2713	learned knowledge to novel genetic manipulations.
2714	- Unseen Cell Contexts: The model should be capable of predicting the response to
2715	observed during the training phase. This evaluates the model's robustness to
2717	the inherent heterogeneity within the K562 cell population.
2718	# Baseline Model Analysis
2719	**SOTA*** GEARS achieves best Pearson correlation in combinatorial prediction
2720 2721	tasks but violates the "no external database" constraint .
2722	Below are detailed critiques of each baselines shortcomings in the context of the
2723 2724	recommendationsgrounded in recent literaturefor how to overcome them. Each point is supported by highquality citations.
2725	1. SC-GPT
2727	**Shortcomings:**
2728 2729	1). **Discrete Perturbation Tokens:** SC-GPT treats each perturbation (e.g. a specific dualguide combination) as a unique token. It cannot form embeddings
2730 2731	<ul> <li>2). **No Zero-Inflated Modeling:** SC-GPTs Gaussian or cross-entropy losses dont account for dropoutdriven zeros common in scRNA-seq, causing biased</li> </ul>
2732	predictions for low-UMI cells 3). **Parameter Bloat for Dense Output:** Extending SC-GPTs languagemodel head to
2733	35 kdimensional gene outputs inflates parameters, hindering training
2734	efficiency and generalization
2736	2. GeneFormer **Short.comings:**
2737	1). **Single-Gene Focus:** GeneFormer has been validated primarily on singlegene
2738	<pre>knockouts, lacking mechanisms to **compose** multiple guide embeddings for combinatorial CRISPRi</pre>
2739	2). **Static Graph Priors:** It uses a fixed genegene network that doesnt adapt
2740	dynamic response modeling
2741 2742 2743	3. **Scalability Issues:** Fullgraph attention over 35 k genes is intractable, so practical implementations subsample to 25 k genesdiscarding potentially important UPR regulators
2744	3. DEEP (Plain MLP)
2745	**Shortcomings:**
2746 2747	<ol> <li>**Ignores Gene Covariance:** Treats each gene independently, missing co- regulation patterns (e.g., ATF6XBP1 axis in UPR)</li> </ol>
2748	
2749	

2750 2751	<ol> <li>**Overfitting Risk:** Millions of parameters on 35 k inputs with limited replicates per combination leads to memorization, not generalization to</li> </ol>
2752 2753 2754	unseen guide sets 3). **No Interpretability:** Provides no insight into which genes or interactions drive predictions, unlike graph-based or attention-based models.
2755 2756 2757 2758 2759 2760 2761 2762 2763	<ol> <li>GEARS         **Shortcomings:**         1). **External Knowledge Dependency:** GEARS integrates a genegene knowledge graph (e.g., from STRING or GO) to regularize embeddings, which violates our no external database constraint         2). **Fixed Graph Structure:** The perturbation relationship graph in GEARS is static, not conditioned on cell-state or UPR context, limiting dynamic response capture.         3). **Heavy GNN Overhead:** Graph neural network message passing on 35 k nodes multiple perturbations incurs high memory and compute costs, impractical for largescale CRISPRi screens.     </li> </ol>
2764 2765	
2766	Recommendations for Improvement
2767	1. **Factorized Perturbation Embeddings**
2769 2770 2771	- **Approach:** Learn a separate embedding \$e_{g}\$ for each guide \$g\$. Represent a perturbation set \$P\$ by a **learned nonlinear composition**.
2772 2773	<ul> <li>**Benefit:** Zeroshot support for unseen guide combinations via embedding arithmetic, as demonstrated by CPA and scGen</li> </ul>
2774	2. **Zero-Inflated Negative Binomial (ZINB) Loss**
2775 2776 2777	<ul> <li>- **Approach:** Replace MSE with a **ZINB loss** that models both dropout probability and overdispersion per gene.</li> </ul>
2778	<ul> <li>**Benefit:** Accounts for scRNA-seq technical noise, improving prediction in low-UMI cells (e.g., ~162 median UMI)</li> </ul>
2780	3. **Learned Dynamic Graph Priors**
2781 2782 2783	<ul> <li>**Approach:** Instead of a fixed PPI graph, **learn genegene affinity weights** from data using a **Gaussian kernel** on baseline coexpression, then refine during training.</li> </ul>
2784 2785	<ul> <li>**Benefit:** Captures UPR pathway rewiring under CRISPRi; avoids external databases (matches constraints)</li> </ul>
2786	4. **Contrastive Pretraining**
2787	- **Approach:**
2789 2790	Pretrain an encoder on (x,0) vs. (x,p) pairs with an **InfoNCE contrastive loss**, pulling matched baselineperturbed embeddings together and pushing apart random mismatches.
2791 2792	<ul> <li>**Benefit:** Disentangles baseline state from perturbation effect; enhances generalization to unseen contexts like new guide sets or cellcycle states</li> </ul>
2793	5. **Neural ODE or OT Trajectory Module**
2795 2796 2797	<pre>- **Approach:** For multiguide dynamics, model latent drift via a **neural ODE     ** \$dz/dt=f(z,hP)dz/dt = f(z,h_P)\$ or **optimaltransport regularization**     linking \$(x,p)\to(x,p')\$.</pre>
2798 2799	<ul> <li>- **Benefit:** Enforces smooth interpolation/extrapolation between perturbation levels (01 guide,,4 guides), capturing graded UPR responses</li> </ul>
2800	6. **Perturbation-Guided Attention Decoder**
2801 2802	<pre>- **Approach:** Use a **crossattention** layer where query = latent state, key /value = perturbation embedding \$h_P\$, then an attentionaugmented decoder</pre>
2803	

2805	to weight genes meet imported by UDD factors (e.g. ATE6 targets)
2806	to weight genes most impacted by ork factors (e.g., Airo targets).
2807	- **Benefit:** Focuses model capacity on biologically salient genes; improves
2808	interpretability and reduces parameter count
2809	. Implementing these enhancements, fasterized compositional embeddings. 7IND
2810	loss: datadriven graph priors: contrastive pretraining: latent ODE/OT modules
2811	; and perturbationguided attentionwill address the key limitations of SC-GPT,
2812	GeneFormer, DEEP, GEARS. This integrated architecture can robustly predict
2813	perturbations in the UPR pathway, while respecting the no external database
2814	constraint and modeling lowsignal, highsparsity regimes.
2815	
2816	
2817	
2818	
2819	
2820	
2821	
2822	
2823	
2824	
2825	
2826	
2827	
2828	
2029	
2830	
2831	
2833	
2834	
2835	
2836	
2837	
2838	
2839	
2840	
2841	
2842	
2843	
2844	
2845	
2846	
2847	
2848	
2849	
2850	
2001	
20 <i>32</i> 2853	
2033 2854	
2054	
2055	
2850	
2858	
2859	

33,000 genes. It's quite large and very sparse. I propose starting with lognormalization to handle the skewness in the data. We could also use PCA for dimensionality reduction to capture the main sources of variation. Maybe around 100-200 principal components would be a good starting point. What do

\*\*Data Preprocessing Expert\*\*: The dataset contains over 1000,000 cells and

\*\*Model Architecture Expert\*\*: I agree with the preprocessing steps. Given the complexity of gene regulatory networks, a graph neural network (GNN) might be suitable to capture gene interactions. However, constructing the gene interaction graph is challenging. We could start with a simple graph based on prior knowledge like protein-protein interaction networks. But our task is to predict the effects of unseen perturbations, so the model needs to be flexible enough to handle new genes not in the training set. Maybe combining a GNN with gene embeddings could work. The embeddings can capture functional

\*\*Training and Optimization Expert\*\*: Training a GNN on this scale will be

computationally intensive. We need an efficient training strategy. I suggest using a learning rate scheduler to adjust the learning rate during training. Also, to prevent overfitting, especially with the high dimensionality,

dropout layers should be incorporated. Given the potential class imbalance in the perturbation types, we might need to use weighted loss functions or some

negatively impact the predictions. Also, the embeddings might not capture all the necessary biological information. Maybe we should consider a more datadriven approach for constructing the gene interactions instead of relying

\*\*Pathway Analyst\*\*: From a biological standpoint, it's crucial to model not just direct effects of perturbations but also downstream effects. Genes don't operate in isolation. A perturbation to one gene can have cascading effects on many others. The model should ideally capture these dynamic changes. Also, cell-to-cell variability is significant in single-cell data. The model needs

\*\*Deep Learning Expert\*\*: I think incorporating attention mechanisms could help

particularly in capturing long-range dependencies. Maybe a hybrid model that

transformer can handle the gene expression data and the GNN can model the interactions. Also, using a multi-task learning approach where the model predicts both the gene expression changes and the perturbation type could

\*\*Critic\*\*: The initial plan covers several important aspects, but the model's
generalization to unseen perturbations remains a critical issue. The proposed
GNN with gene embeddings is a solid foundation, but the reliance on prior
knowledge for the gene interaction graph could be a limitation. The datadriven approach suggested by the Self-Critic could mitigate this. Also, the
training strategy needs to address the high dimensionality and sparsity
effectively. The inclusion of a contrastive learning component, as suggested

by the Deep Learning Specialist, could enhance the model's ability to distinguish different perturbation effects. However, the evaluation metrics

should be clearly defined to assess both the predictive accuracy and

53

the model focus on the most relevant genes for each perturbation. Transformers have been successful in other areas of bioinformatics,

combines transformer architecture with GNNs would be powerful. The

\*\*Self-Critic\*\*: The ideas so far are good, but I'm concerned about the model's ability to generalize. The GNN relies heavily on the quality of the gene interaction graph. If the graph is incomplete or inaccurate, it could

#### 2860 G.4. Graph based discussion

Round 1

#### **Graph Based Discussion Output**

you all think?

similarities between genes.

form of resampling.

solely on prior knowledge.

improve performance.

biological relevance.

to account for this inherent noise.

2	8	6	1	
2	8	6	2	ſ
2	8	6	3	F
2	0	6	7	L
2	0	6	4	
2	8	0	5	
2	8	6	6	
2	8	6	7	
2	8	6	8	
2	8	6	9	
2	8	7	0	
2	8	7	1	
2	8	7	2	
2	8	7	3	
2	8	7	4	
2	8	7	5	
2	8	7	6	
$\frac{-}{2}$	8	7	7	
$\frac{2}{2}$	8	, 7	8	
2	8	, 7	0	
2	0	2	י ה	
2	0	0	1	
2	ð o	ð o	1	
2	ð	ð	2	
2	8	8	3	
2	8	8	4	
2	8	8	5	
2	8	8	6	
2	8	8	7	
2	8	8	8	
2	8	8	9	
2	8	9	0	
2	8	9	1	
2	8	9	2	
2	8	9	3	
$\frac{2}{2}$	8	0	Δ	
イ つ	2	2 0	+ 5	
ム つ	Q Q	1 0	5 6	
2	0	7 0	0 7	
2	ð o	9 0	/	
2	8	9	8	
2	8	9	9	
2	9	0	0	
2	9	0	1	
2	9	0	2	
2	9	0	3	
2	9	0	4	
2	9	0	5	
2	9	0	6	
2	9	0	7	
2	9	0	8	
2	9	0	9	
$\frac{2}{2}$	0	1	0	
イ つ	0	1 1	1	C
2	7	1	1 つ	
2	9	1	2	
1	ų	1	-5	

2915	Graph Based Discussion Output
2916	Round 2
2917 2918	
2919	**Data Preprocessing Expert**: I've thought more about the preprocessing. Besides PCA, we could use a variational autoencoder (VAE) for dimensionality
2920	reduction. VAEs are generative models and might better capture the complex
2921	distributions in the data. Also, for handling the sparsity, adding a step to
2922	neighbors imputation method.
2923	++Model Architecture Evnert++. I like the idea of combining transformers and GNNs
2924	For the gene interaction graph in the GNN, perhaps we can construct it
2925	dynamically based on the data. Using the correlations between genes in the
2926	training data to build an adjacency matrix. This way, the graph is data- driven and doesn't rely solely on prior knowledge. Also, for the gene
2927	embeddings, we can initialize them with pre-trained embeddings from resources
2928	like Gene2Vec and then fine-tune them during training.
2929	**Training and Optimization Expert**: Regarding training, I suggest using the
2930	AdamW optimizer with a OneCycle learning rate schedule. This schedule helps
2931	In finding a good balance between the learning rate and the number of epochs. For handling class imbalance, focal loss could be an option. It focuses the
2932	training more on the hard-to-classify examples. Also, implementing gradient
2933	clipping will help prevent exploding gradients, which is important given the model's complexity
2934	MOULT D COMPTEXTLY.
2935	**Self-Critic**: While the hybrid model sounds promising, I'm still concerned
2936	training data. Maybe starting with a simpler model and gradually adding
2937	complexity would be a better approach. Also, how do we validate that the
2938	model's predictions are biologically meaningful? We need a solid evaluation strategy beyond just metrics like MSE and PCC
2939	
2940 2941	**Pathway Analyst**: I agree with the need for dynamic modeling. The effects of
2941	or using a time-series approach could help model the progression of changes.
2943	Also, including cell cycle phase as a covariate might improve predictions,
2944	as gene expression varies significantly across different phases.
2945	**Deep Learning Expert**: I've considered the multi-task learning approach more.
2946	another branch classifying the perturbation type. This forces the model to
2947	learn features that are useful for both tasks. Also, using a contrastive loss
2948	component could help the model distinguish between different perturbation
2949	particularly useful for capturing long-range dependencies in gene expression
2950	data, which could be crucial for understanding the downstream effects of
2951	perturbations.
2952	**Critic**: The revised plan addresses several of my previous concerns. The use
2953	or a vam for preprocessing and a data-driven gene interaction graph are positive steps. The multi-task learning and contrastive loss components add
2954	robustness. However, the evaluation strategy still needs clarification. We
2955	should define specific metrics for both predictive performance (MSE, PCC, R2)
2956	to handle combinatorial perturbations should be explicitly tested.
2957	
2738 2050 -	
2939 2960	Graph Based Discussion Output
2961	Round 3
2962	
2963	**Data Preprocessing Expert**: I've further refined the preprocessing steps.
2964	After careful consideration, I believe using a combination of log-
2965	normalization, imputation for missing values, and PCA for dimensionality reduction will provide a solid foundation for the model. If we also added data
2966	augmentation techniques like adding Gaussian noise and random masking to
2967	improve model robustness. The log-normalization step is crucial for handling
2968	the skewness in the gene expression data, while the PCA helps in reducing
2969	

2970	noise and computational complexity
2971	noise and completeronar complexity.
2972	**Model Architecture Expert**: I've updated the model architecture to better
2973	expression data, a separate embedding layer for perturbations, and a
2974	transformer for capturing long-range dependencies. The dynamic graph
2975	construction based on data correlations addresses the concern about prior
2976	model should have a good starting point for capturing gene functionalities.
2977	The VAE is particularly useful for denoising and capturing the complex
2978	distributions in the data.
2979	$\star\star$ Training and Optimization Expert $\star\star\colon$ I've adjusted the training strategy to
2980	include three phases: 1) VAE pretraining for 20 epochs, 2) Transformer
2981	for 50 epochs. The learning rate scheduling will be adjusted per phase to
2982	optimize training. I've also incorporated focal loss to handle class
2983	imbalance and gradient clipping to prevent exploding gradients. The focal
2984	hard-to-classify examples.
2985	Colf Criticult I/m more confident new The plan to start with a simpler model
2986	and incrementally add complexity is reasonable. I've also developed a
2987	comprehensive evaluation strategy that includes comparisons with known
2988	biological pathways and interactions to ensure the predictions make sense
2989	evaluation strategy will include both quantitative metrics and qualitative
2990	assessments to ensure biological relevance.
2991	**Pathway Analyst**: I'm satisfied with the incorporation of biological
2992	considerations. The model's ability to capture dynamic changes and the
2993	inclusion of cell cycle phase as a covariate will enhance the biological
2994	insights into the effects of CRISPRi perturbations. The inclusion of cell
2995	cycle phase as a covariate is particularly important as it helps in
2997	accounting for the inherent biological noise in single-cell data.
2998	**Deep Learning Expert**: I've implemented the multi-task learning and
2999	contrastive loss components. The model can now predict both gene expression changes and perturbation types, which should improve performance. I've also
3000	added a recurrent component to model the progression of changes over time.
3001	The multi-task learning approach ensures that the model learns features that
3002	are relevant for multiple tasks, improving its overall generalization capabilities. The transformer's self-attention mechanism allows the model to
3003	weigh the importance of different genes dynamically, which is particularly
3004	useful for capturing the complex regulatory networks in single-cell data.
3005	**Critic**: The final plan addresses most of the previous concerns. The
3006	preprocessing steps are comprehensive, the model architecture is innovative
3007	and well-justified, and the training strategy is robust. The evaluation metrics are clearly defined, including both predictive performance (MSE/PCC/
3008	R2) and biological relevance (MSE_DE/PCC_DE/R2_DE). The model's ability to
3009	handle combinatorial perturbations through the dynamic graph and contrastive
3010	toss is a sciong point. I betteve this plan has a high likelihood of success.
3011	
3012	
3013	
3014	
3015	
3016	
3017	
3018	
3019	
3020	
3021	
2022	
2024	
3024	

3025 G.5. Method Design Research Plan

3026	Graph Based Discussion Output
3028	
3029	# Data Preprocessing
3030	Steps:
3031	## Load Data Load the gene expression matrix and metadata from the Norman et al. (2019)
3032	dataset using appropriate libraries like scanpy or anndata.
3034	## Filter Low-Quality Cells and Genes
3035	Remove cells with low RNA content or high mitochondrial content.
3036	Filter out genes with low expression across cells.
3037	## Normalize Data Normalize the gene expression counts to account for differences in sequencing
3038 3039	depth between cells. This can be done using methods like library size normalization.
3040	## Log-Transformation
3041	Apply log-transformation to the gene expression data to handle skewness.
3043 3044	## Batch Effect Correction Use methods like harmony or Combat to remove batch effects from the data.
3045	## Feature Selection
3046	Identify highly variable genes (HVGs) to reduce dimensionality. This can be done using statistical methods to select genes with the highest variability across
3047	cells.
3048	## PCA Dimensionality Reduction
3049	Apply PCA to further reduce the dimensionality of the gene expression data. This below in centuring the main sources of variation in a lower-dimensional space
3050	heips in capturing the main sources of variation in a fower dimensional space.
3051	## Perturbation Encoding
3053	Extract the perturbation information from the dataset's metadata.
3054	Alternatively, use learned embeddings for each gene in the perturbation.
3055	## Control Sample Handling
3056 3057	Identify and use control cells (unperturbed) to establish baseline gene expression profiles.
3058	## Data Augmentation
3059 3060	Add Gaussian noise and apply random masking to the data to improve model robustness.
3061	## Data Splitting
3062	perturbations are held out for validation, and testing to evaluate the model's
3063	ability to generalize to unseen perturbations and cell contexts.
3064 3065	# Model Design
3066	## Overview The proposed model is a hybrid neural network architecture designed to predict
3067	post-perturbation gene expression profiles in single cells. It integrates
3068	three key components: Variational Autoencoder (VAE) for learning robust low-dimensional
3069	representations of gene expression data.
3070	Graph Neural Network (GNN) for modeling gene-gene interactions using a
3071	Transformer for capturing long-range dependencies and complex patterns in the
3072	fused gene-perturbation representations.
3073	and noise characteristic of single-cell data while providing biological
3075	interpretability through attention mechanisms and graph structures.
3076	## Key Components
3077	1. Variational Autoencoder (VAE) Encoder
3078	
3079	

3080	
3081	Purpose: To compress high-dimensional gene expression data into a compact latent
3082	representation while capturing the underlying data distribution.
3083	Architecture:
308/	Input layer with size matching the number of principal components (e.g., 128).
2005	Two hidden layers with decreasing dimensions (e.g., 256 128 64).
3085	Supplier and the second
3086	in scRNL-sed data
3087	Biological Interpretability: The learned latent space can be analyzed for
3088	biological patterns.
3080	
2000	2. Perturbation Embedding Layer
3090	Purpose. To convert perturbation identities (single or paired genes) into dense
3091	vector representations.
3092	Architecture:
3093	Input layer with size matching the number of unique perturbations (one-hot
3094	encoded).
3005	Embedding layer projecting into a lower-dimensional space (e.g., 32 dimensions).
2006	Feasibility: Embedding layers are standard for categorical data and reduce
3096	Biological Interpretability. Similar perturbations will have similar embeddings
3097	allowing analysis of perturbation relationships.
3098	
3099	3. Graph Neural Network (GNN)
3100	Purpose. To model gene regulatory networks (GRNs) and capture gene-gene
2101	interactions.
2102	Architecture:
3102	Dynamically construct a gene interaction graph using correlations from the
3103	training data.
3104	Use Graph Convolutional Layers to propagate information across the graph.
3105	Feasibility: GNNs are effective for capturing relational data and have been
3106	applied to gene networks.
2107	relationships
3107	Teracionanipa.
3108	4. Transformer Encoder
3109	Durpass, To conture complex, long range dependencies in the fund cons
3110	nerturbation representations
3111	Architecture:
3112	Input layer combining VAE latent vectors and perturbation embeddings.
2112	Multi-head self-attention mechanisms to weigh the importance of different genes.
3113	
211/	Feed-forward neural networks for non-linear transformations.
5114	Feed-forward neural networks for non-linear transformations. Feasibility: Transformers excel at modeling sequential and complex data patterns.
3114	Feed-forward neural networks for non-linear transformations. Feasibility: Transformers excel at modeling sequential and complex data patterns. Biological Interpretability: Attention weights can reveal important gene
3114 3115 3116	Feed-forward neural networks for non-linear transformations. Feasibility: Transformers excel at modeling sequential and complex data patterns. Biological Interpretability: Attention weights can reveal important gene interactions. 5 Fusion and Output Lawors
3114 3115 3116 3117	Feed-forward neural networks for non-linear transformations. Feasibility: Transformers excel at modeling sequential and complex data patterns. Biological Interpretability: Attention weights can reveal important gene interactions. 5. Fusion and Output Layers Purpose: To combine information from the VAE and perturbation branches and
3114 3115 3116 3117 3118	<ul> <li>Feed-forward neural networks for non-linear transformations.</li> <li>Feasibility: Transformers excel at modeling sequential and complex data patterns.</li> <li>Biological Interpretability: Attention weights can reveal important gene interactions.</li> <li>5. Fusion and Output Layers</li> <li>Purpose: To combine information from the VAE and perturbation branches and produce the final prediction.</li> </ul>
3114 3115 3116 3117 3118 3110	Feed-forward neural networks for non-linear transformations. Feasibility: Transformers excel at modeling sequential and complex data patterns. Biological Interpretability: Attention weights can reveal important gene interactions. 5. Fusion and Output Layers Purpose: To combine information from the VAE and perturbation branches and produce the final prediction. Architecture:
3115 3116 3117 3118 3119	Feed-forward neural networks for non-linear transformations. Feasibility: Transformers excel at modeling sequential and complex data patterns. Biological Interpretability: Attention weights can reveal important gene interactions. 5. Fusion and Output Layers Purpose: To combine information from the VAE and perturbation branches and produce the final prediction. Architecture: Concatenate VAE latent vectors and perturbation embeddings.
3114 3115 3116 3117 3118 3119 3120	Feed-forward neural networks for non-linear transformations. Feasibility: Transformers excel at modeling sequential and complex data patterns. Biological Interpretability: Attention weights can reveal important gene interactions. 5. Fusion and Output Layers Purpose: To combine information from the VAE and perturbation branches and produce the final prediction. Architecture: Concatenate VAE latent vectors and perturbation embeddings. Dense layers to transform the fused representation into gene expression
3114         3115         3116         3117         3118         3119         3120         3121	<ul> <li>Feed-forward neural networks for non-linear transformations.</li> <li>Feasibility: Transformers excel at modeling sequential and complex data patterns.</li> <li>Biological Interpretability: Attention weights can reveal important gene interactions.</li> <li>5. Fusion and Output Layers</li> <li>Purpose: To combine information from the VAE and perturbation branches and produce the final prediction.</li> <li>Architecture:</li> <li>Concatenate VAE latent vectors and perturbation embeddings.</li> <li>Dense layers to transform the fused representation into gene expression predictions.</li> </ul>
3114 3115 3116 3117 3118 3119 3120 3121 3122	<ul> <li>Feed-forward neural networks for non-linear transformations.</li> <li>Feasibility: Transformers excel at modeling sequential and complex data patterns.</li> <li>Biological Interpretability: Attention weights can reveal important gene interactions.</li> <li>5. Fusion and Output Layers</li> <li>Purpose: To combine information from the VAE and perturbation branches and produce the final prediction.</li> <li>Architecture:</li> <li>Concatenate VAE latent vectors and perturbation embeddings.</li> <li>Dense layers to transform the fused representation into gene expression predictions.</li> <li>Feasibility: Concatenation is a simple yet effective fusion strategy.</li> <li>Pielogical Interpretability: The final output can be directly compared to</li> </ul>
3114 3115 3116 3117 3118 3119 3120 3121 3122 3123	Feed-forward neural networks for non-linear transformations. Feasibility: Transformers excel at modeling sequential and complex data patterns. Biological Interpretability: Attention weights can reveal important gene interactions. 5. Fusion and Output Layers Purpose: To combine information from the VAE and perturbation branches and produce the final prediction. Architecture: Concatenate VAE latent vectors and perturbation embeddings. Dense layers to transform the fused representation into gene expression predictions. Feasibility: Concatenation is a simple yet effective fusion strategy. Biological Interpretability: The final output can be directly compared to observed gene expression changes
3114 3115 3116 3117 3118 3119 3120 3121 3122 3123 3124	<ul> <li>Feed-forward neural networks for non-linear transformations.</li> <li>Feasibility: Transformers excel at modeling sequential and complex data patterns.</li> <li>Biological Interpretability: Attention weights can reveal important gene interactions.</li> <li>5. Fusion and Output Layers</li> <li>Purpose: To combine information from the VAE and perturbation branches and produce the final prediction.</li> <li>Architecture:</li> <li>Concatenate VAE latent vectors and perturbation embeddings.</li> <li>Dense layers to transform the fused representation into gene expression predictions.</li> <li>Feasibility: Concatenation is a simple yet effective fusion strategy.</li> <li>Biological Interpretability: The final output can be directly compared to observed gene expression changes.</li> </ul>
3114 3115 3116 3117 3118 3119 3120 3121 3122 3123 3124 3125	<pre>Feed-forward neural networks for non-linear transformations. Feasibility: Transformers excel at modeling sequential and complex data patterns. Biological Interpretability: Attention weights can reveal important gene interactions. 5. Fusion and Output Layers Purpose: To combine information from the VAE and perturbation branches and produce the final prediction. Architecture: Concatenate VAE latent vectors and perturbation embeddings. Dense layers to transform the fused representation into gene expression predictions. Feasibility: Concatenation is a simple yet effective fusion strategy. Biological Interpretability: The final output can be directly compared to observed gene expression changes. ## Feasibility and Biological Interpretability</pre>
3114 3115 3116 3117 3118 3119 3120 3121 3122 3123 3124 3125	<pre>Feed-forward neural networks for non-linear transformations. Feasibility: Transformers excel at modeling sequential and complex data patterns. Biological Interpretability: Attention weights can reveal important gene interactions. 5. Fusion and Output Layers Purpose: To combine information from the VAE and perturbation branches and produce the final prediction. Architecture: Concatenate VAE latent vectors and perturbation embeddings. Dense layers to transform the fused representation into gene expression predictions. Feasibility: Concatenation is a simple yet effective fusion strategy. Biological Interpretability: The final output can be directly compared to observed gene expression changes. ### Feasibility and Biological Interpretability #### Feasibility</pre>
3114         3115         3116         3117         3118         3119         3120         3121         3122         3123         3124         3125         3126	<pre>Feed-forward neural networks for non-linear transformations. Feasibility: Transformers excel at modeling sequential and complex data patterns. Biological Interpretability: Attention weights can reveal important gene interactions. 5. Fusion and Output Layers Purpose: To combine information from the VAE and perturbation branches and produce the final prediction. Architecture: Concatenate VAE latent vectors and perturbation embeddings. Dense layers to transform the fused representation into gene expression predictions. Feasibility: Concatenation is a simple yet effective fusion strategy. Biological Interpretability: The final output can be directly compared to observed gene expression changes. ### Feasibility and Biological Interpretability #### Feasibility Computational Efficiency: The model uses well-established and optimized</pre>
3114         3115         3116         3117         3118         3119         3120         3121         3122         3123         3124         3125         3126         3127	<pre>Feed-forward neural networks for non-linear transformations. Feasibility: Transformers excel at modeling sequential and complex data patterns. Biological Interpretability: Attention weights can reveal important gene interactions. 5. Fusion and Output Layers Purpose: To combine information from the VAE and perturbation branches and produce the final prediction. Architecture: Concatenate VAE latent vectors and perturbation embeddings. Dense layers to transform the fused representation into gene expression predictions. Feasibility: Concatenation is a simple yet effective fusion strategy. Biological Interpretability: The final output can be directly compared to observed gene expression changes. ## Feasibility and Biological Interpretability ### Feasibility Computational Efficiency: The model uses well-established and optimized components (VAE, Transformer, GNN) that can be efficiently trained on GPU.</pre>
3114         3115         3116         3117         3118         3119         3120         3121         3122         3123         3124         3125         3126         3127         3128	<pre>Feed-forward neural networks for non-linear transformations. Feasibility: Transformers excel at modeling sequential and complex data patterns. Biological Interpretability: Attention weights can reveal important gene interactions. 5. Fusion and Output Layers Purpose: To combine information from the VAE and perturbation branches and produce the final prediction. Architecture: Concatenate VAE latent vectors and perturbation embeddings. Dense layers to transform the fused representation into gene expression predictions. Feasibility: Concatenation is a simple yet effective fusion strategy. Biological Interpretability: The final output can be directly compared to observed gene expression changes. ## Feasibility and Biological Interpretability ### Feasibility Computational Efficiency: The model uses well-established and optimized components (VAE, Transformer, GNN) that can be efficiently trained on GPU. Scalability: The architecture is designed to handle the high dimensionality of</pre>
3114         3115         3116         3117         3118         3119         3120         3121         3122         3123         3124         3125         3126         3127         3128         3129	<pre>Feed-forward neural networks for non-linear transformations. Feasibility: Transformers excel at modeling sequential and complex data patterns. Biological Interpretability: Attention weights can reveal important gene interactions. 5. Fusion and Output Layers Purpose: To combine information from the VAE and perturbation branches and produce the final prediction. Architecture: Concatenate VAE latent vectors and perturbation embeddings. Dense layers to transform the fused representation into gene expression predictions. Feasibility: Concatenation is a simple yet effective fusion strategy. Biological Interpretability: The final output can be directly compared to observed gene expression changes. ## Feasibility and Biological Interpretability ### Feasibility Computational Efficiency: The model uses well-established and optimized components (VAE, Transformer, GNN) that can be efficiently trained on GPU. Scalability: The architecture is designed to handle the high dimensionality of single-cell data while maintaining reasonable computational requirements.</pre>
3114         3115         3116         3117         3118         3119         3120         3121         3122         3123         3124         3125         3126         3127         3128         3129         3130	<pre>Feed-forward neural networks for non-linear transformations. Feasibility: Transformers excel at modeling sequential and complex data patterns. Biological Interpretability: Attention weights can reveal important gene interactions. 5. Fusion and Output Layers Purpose: To combine information from the VAE and perturbation branches and produce the final prediction. Architecture: Concatenate VAE latent vectors and perturbation embeddings. Dense layers to transform the fused representation into gene expression predictions. Feasibility: Concatenation is a simple yet effective fusion strategy. Biological Interpretability: The final output can be directly compared to observed gene expression changes. ## Feasibility Computational Efficiency: The model uses well-established and optimized components (VAE, Transformer, GNN) that can be efficiently trained on GPU. Scalability: The architecture is designed to handle the high dimensionality of single-cell data while maintaining reasonable computational requirements. Implementation: The model can be implemented using popular deep learning</pre>
3114         3115         3116         3117         3118         3119         3120         3121         3122         3123         3124         3125         3126         3127         3128         3129         3130	<pre>Feed-forward neural networks for non-linear transformations. Feasibility: Transformers excel at modeling sequential and complex data patterns. Biological Interpretability: Attention weights can reveal important gene interactions. 5. Fusion and Output Layers Purpose: To combine information from the VAE and perturbation branches and produce the final prediction. Architecture: Concatenate VAE latent vectors and perturbation embeddings. Dense layers to transform the fused representation into gene expression predictions. Feasibility: Concatenation is a simple yet effective fusion strategy. Biological Interpretability: The final output can be directly compared to observed gene expression changes. ## Feasibility and Biological Interpretability ### Feasibility Computational Efficiency: The model uses well-established and optimized components (VAE, Transformer, GNN) that can be efficiently trained on GPU. Scalability: The architecture is designed to handle the high dimensionality of single-cell data while maintaining reasonable computational requirements. Implementation: The model can be implemented using popular deep learning frameworks like PyTorch.</pre>
3114         3115         3116         3117         3118         3119         3120         3121         3122         3123         3124         3125         3126         3127         3128         3129         3130         3131	<pre>Feed-forward neural networks for non-linear transformations. Feasibility: Transformers excel at modeling sequential and complex data patterns. Biological Interpretability: Attention weights can reveal important gene interactions. 5. Fusion and Output Layers Purpose: To combine information from the VAE and perturbation branches and produce the final prediction. Architecture: Concatenate VAE latent vectors and perturbation embeddings. Dense layers to transform the fused representation into gene expression predictions. Feasibility: Concatenation is a simple yet effective fusion strategy. Biological Interpretability: The final output can be directly compared to observed gene expression changes. ## Feasibility Computational Efficiency: The model uses well-established and optimized components (VAE, Transformer, GNN) that can be efficiently trained on GPU. Scalability: The architecture is designed to handle the high dimensionality of single-cell data while maintaining reasonable computational requirements. Implementation: The model can be implemented using popular deep learning frameworks like PyTorch. ### Biological Interpretability</pre>
3114         3115         3116         3117         3118         3119         3120         3121         3122         3123         3124         3125         3126         3127         3128         3129         3130         3131         3132	<pre>Feed-forward neural networks for non-linear transformations. Feasibility: Transformers excel at modeling sequential and complex data patterns. Biological Interpretability: Attention weights can reveal important gene interactions. 5. Fusion and Output Layers Purpose: To combine information from the VAE and perturbation branches and produce the final prediction. Architecture: Concatenate VAE latent vectors and perturbation embeddings. Dense layers to transform the fused representation into gene expression predictions. Feasibility: Concatenation is a simple yet effective fusion strategy. Biological Interpretability: The final output can be directly compared to observed gene expression changes. ## Feasibility Computational Efficiency: The model uses well-established and optimized components (VAE, Transformer, GNN) that can be efficiently trained on GPU. Scalability: The architecture is designed to handle the high dimensionality of single-cell data while maintaining reasonable computational requirements. Implementation: The model can be implemented using popular deep learning frameworks like PyTorch. ### Biological Interpretability</pre>
3114         3115         3116         3117         3118         3119         3120         3121         3122         3123         3124         3125         3126         3127         3128         3129         3130         3131         3132         3133	<pre>Feed-forward neural networks for non-linear transformations. Feasibility: Transformers excel at modeling sequential and complex data patterns. Biological Interpretability: Attention weights can reveal important gene interactions. 5. Fusion and Output Layers Purpose: To combine information from the VAE and perturbation branches and produce the final prediction. Architecture: Concatenate VAE latent vectors and perturbation embeddings. Dense layers to transform the fused representation into gene expression predictions. Feasibility: Concatenation is a simple yet effective fusion strategy. Biological Interpretability: The final output can be directly compared to observed gene expression changes. ## Feasibility Computational Efficiency: The model uses well-established and optimized components (VAE, Transformer, GNN) that can be efficiently trained on GPU. Scalability: The architecture is designed to handle the high dimensionality of single-cell data while maintaining reasonable computational requirements. Implementation: The model can be implemented using popular deep learning frameworks like PyTorch. ### Biological Interpretability</pre>

```
3135
          Attention Mechanisms: The Transformer's self-attention layers can highlight
3136
              important gene interactions, providing insights into regulatory mechanisms.
3137
          Perturbation Embeddings: The learned embeddings can be analyzed to understand
3138
              relationships between different perturbations.
          Gene Regulatory Networks: The GNN component can reveal key regulatory
3139
             interactions within the gene network.
3140
          Latent Space Analysis: The VAE's latent space can be visualized to identify
              patterns corresponding to biological conditions.
3141
3142
          # Training Strategy
3143
          ## Overview
3144
          The training strategy is designed to optimize the model's ability to predict post
3145
              -perturbation gene expression profiles while ensuring generalization to
              unseen perturbations and cell contexts. The strategy incorporates several key
3146
              components including a custom loss function, advanced optimization techniques, and mechanisms for monitoring and preventing overfitting.
3147
3148
          ## Key Components
3149
3150
          ### 1. Loss Function
          **Components**:
3151
          - **Reconstruction Loss**: Measures the difference between the predicted and true
3152
               gene expression profiles using Mean Squared Error (MSE).
3153
          - **KL Divergence Loss**: Ensures the learned latent space distribution from the
              VAE component is close to a prior distribution (typically a standard normal
3154
              distribution).
3155
          **Implementation**:
'`'python
3156
3157
          def loss_function(output, x, mu, logvar, beta=1.0):
3158
             recon_loss = F.mse_loss(output, x)
             kl_loss = -0.5 * torch.sum(1 + logvar - mu.pow(2) - logvar.exp())
3159
             kl_loss = kl_loss / x.size(0) # Normalize by batch size
3160
             total_loss = recon_loss + beta * kl_loss
          return total_loss, recon_loss, kl_loss
3161
3162
          ### 2. Optimizer and Learning Rate Scheduler
3163
          **Components**:
3164
          - **AdamW Optimizer**: Used for its adaptive learning rate capabilities and
3165
              decoupled weight decay.
          - **OneCycleLR Scheduler**: Dynamically adjusts the learning rate over the
3166
              training period to find an optimal balance between learning rate and
3167
              convergence speed.
3168
          **Implementation**:
3169
          ```python
3170
          optimizer = torch.optim.AdamW(model.parameters(), lr=1e-3, weight_decay=1e-5)
3171
          scheduler = torch.optim.lr_scheduler.OneCycleLR(
             optimizer,
3172
             max_lr=1e-3,
3173
             steps_per_epoch=len(train_loader),
             epochs=num_epochs,
pct_start=0.2,
3174
3175
             anneal_strategy='cos'
3176
          ),,,
3177
3178
          ### 3. Regularization Techniques
          **Components**:
3179
          - **Gradient Clipping**: Prevents exploding gradients by clipping gradient values
3180
               during backpropagation.
3181
          - **Early Stopping**: Monitors validation loss and stops training if it doesn't
              improve for a specified number of epochs.
3182
3183
          **Implementation**:
'''python
3184
          def train_epoch(model, train_loader, optimizer, device, beta):
3185
             model.train()
3186
             total_loss = 0
             for batch in train_loader:
3187
3188
3189
```

Submission and Formatting Instructions for ICML 2025 GenBio Workshop

```
3190
                x, pert = batch
3191
                x, pert = x.to(device), pert.to(device)
                optimizer.zero_grad()
3193
                output, mu, logvar = model(x, pert)
3194
3195
                loss, _, _ = loss_function(output, x, mu, logvar, beta)
3196
                loss.backward()
3197
                torch.nn.utils.clip_grad_norm_(model.parameters(), max_norm=1.0)
3198
                optimizer.step()
3199
                total_loss += loss.item()
3200
             return total_loss / len(train_loader)
          def validate(model, val_loader, device, beta):
             model.eval()
             total_loss = 0
3204
             with torch.no_grad():
3205
                for batch in val_loader:
                   x, pert = batch
3206
                   x, pert = x.to(device), pert.to(device)
3208
                   output, mu, logvar = model(x, pert)
3209
                   loss, _, _ = loss_function(output, x, mu, logvar, beta)
total_loss += loss.item()
3210
3211
             return total_loss / len(val_loader)
3212
          def train_model(model, train_loader, val_loader, optimizer, scheduler, device,
3213
              num_epochs=100, beta=1.0, patience=5):
3214
             best_loss = float('inf')
3215
             patience_counter = 0
             best_model_weights = None
3216
3217
             for epoch in range (num_epochs):
                train_loss = train_epoch(model, train_loader, optimizer, device, beta)
3218
                val_loss = validate(model, val_loader, device, beta)
3219
                scheduler.step()
3220
3221
                print(f'Epoch {epoch+1}/{num_epochs} - Train Loss: {train_loss:.4f} - Val
                    Loss: {val_loss:.4f}')
3222
3223
                if val_loss < best_loss:
3224
                   best_loss = val_loss
                   patience_counter = 0
3225
                   best_model_weights = model.state_dict()
3226
                else:
3227
                   patience_counter +=1
                    if patience_counter >= patience:
3228
                       print ('Early stopping triggered')
3229
                      break
             model.load_state_dict(best_model_weights)
3231
          ,,,return model
3232
3233
          ### 4. Training Loop with Monitoring
3234
          **Components**:
          - **Epoch Training**: Trains the model on the entire training dataset for one
3235
             pass.
3236
          - **Validation**: Evaluates the model on the validation set after each epoch.
           **Checkpointing**: Saves the model weights if the validation loss improves.
3237
3238
          **Implementation**:
          ''python
3239
          # Example usage
3240
          if _____ == "____
                            _main__":
3241
             # Assuming adata is your preprocessed AnnData object
             # Convert your data to PyTorch tensors
3242
3243
3244
```

Submission and Formatting Instructions for ICML 2025 GenBio Workshop

```
X = torch.tensor(data_augmented, dtype=torch.float32)
          perturbations = torch.tensor(perturbations, dtype=torch.float32)
           # Create dataset and dataloader
          dataset = TensorDataset(X, perturbations)
           train_size = int(0.8 * len(dataset))
           test_size = len(dataset) - train_size
           train_dataset, test_dataset = torch.utils.data.random_split(dataset, [
              train_size, test_size])
           train_loader = DataLoader(train_dataset, batch_size=32, shuffle=True)
          val_loader = DataLoader(test_dataset, batch_size=32, shuffle=False)
           # Initialize model, optimizer, and device
           device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
          model = HybridModel(input_dim=pca_dim, pert_dim=pert_dim).to(device)
           optimizer = torch.optim.AdamW(model.parameters(), lr=1e-3, weight_decay=1e-5)
           scheduler = torch.optim.lr_scheduler.OneCycleLR(
              optimizer,
             max_lr=1e-3,
              steps_per_epoch=len(train_loader),
              epochs=100,
             pct_start=0.2,
              anneal_strategy='cos'
           )
           # Train the model
           trained_model = train_model(model, train_loader, val_loader, optimizer,
              scheduler, device, beta=0.1)
        • • •
        ## Expert Recommendations Incorporated
        ### Addressing Class Imbalance
         **Weighted Loss**: Assign higher weights to underrepresented classes in the
           loss function.
        - **Oversampling**: Increase the sampling rate of underrepresented perturbations
           during training.
        ### Enhancing Generalization
        - **Data Augmentation**: Apply techniques like Gaussian noise addition and random
            masking to improve model robustness.
        - **VAE Regularization**: Use the KL divergence term in the loss function to
           ensure the latent space is well-regularized.
        ### Improving Training Stability
        - **Gradient Clipping**: Prevents gradients from exploding during backpropagation.
        - **Learning Rate Scheduling**: Uses a OneCycle schedule to dynamically adjust
           the learning rate.
        ### Biological Interpretability
        - **Attention Mechanisms**: The transformer component provides attention weights
           that can be analyzed to understand gene interactions.
        - **Perturbation Embeddings**: The learned embeddings can be visualized to
           explore relationships between different perturbations.
        This comprehensive training strategy ensures the model is trained effectively
           while addressing the specific challenges of single-cell gene expression data,
            including class imbalance, high dimensionality, and the need for
           generalization to unseen perturbations. The incorporation of expert
           recommendations further enhances the model's performance and biological
           relevance.
329 G.6. Example model code
        graph TD
```

```
3294 graph ID
3295 A[Raw Gene Expression Data] --> B[Filter Low-Quality Cells and Genes]
B --> C[Normalize Data]
3296 C --> D[Log-Transformation]
3297 D --> E[Batch Effect Correction]
3298 E --> F[Feature Selection]
F --> G[PCA Reduction]
```

```
H[Perturbation Data] --> I[One-Hot Encoding]
        G --> J[VAE Encoder]
        I --> K[Perturbation Embedding]
        J --> L[Fusion Layer]
        K --> L
        L --> M[Transformer Encoder]
        M --> N[Output Layer]
        N --> O[Predicted Expression Profile]
        // Data Preprocessing
3309 function preprocess_data(adata, pca_dim):
        // Filter low-quality cells and genes
        filter_cells(adata, min_genes=200)
        filter_genes(adata, min_cells=3)
        // Normalize data
        normalize_total(adata, target_sum=1e4)
        log1p(adata)
        // Identify highly variable genes
        highly_variable_genes(adata, n_top_genes=3000)
        data = adata[:, adata.var['highly_variable']].X
        // Batch effect correction
        combat(adata, key='batch')
        // Standardize data
        scaler = StandardScaler()
        data = scaler.fit_transform(data)
        data = clip_values(data, -10, 10)
        data = scale(data)
        // PCA reduction
        pca = PCA(n_components=pca_dim)
        data_pca = pca.fit_transform(data)
        // Perturbation encoding
        perturbations = one_hot_encode(adata.obs['perturbation'])
        // Data augmentation
        data_augmented = add_gaussian_noise(data_pca)
        data_augmented = apply_random_masking(data_augmented)
        return data_augmented, perturbations
     // Model Training
     function train_model(model, train_loader, optimizer, device, beta):
        model.train()
        total_loss = 0
        for each batch in train_loader:
           x, pert = batch
           x, pert = move_to_device(x, pert, device)
           optimizer.zero_grad()
           output, mu, logvar = model(x, pert)
           recon_loss = mean_squared_error(output, x)
           kl_loss = compute_kl_loss(mu, logvar)
loss = recon_loss + beta * kl_loss
           loss.backward()
           optimizer.step()
           total_loss += loss.item()
        return total_loss / number_of_batches(train_loader)
     '''python
    class VAEEncoder(nn.Module):
        def __init__(self, input_dim, latent_dim, hidden_dim):
```

```
super().__init__()
self.fc1 = nn.Linear(input_dim, hidden_dim)
           self.fc_mu = nn.Linear(hidden_dim, latent_dim)
           self.fc_logvar = nn.Linear(hidden_dim, latent_dim)
        def forward(self, x):
    h = F.gelu(self.fcl(x))
           mu = self.fc_mu(h)
           loqvar = self.fc_logvar(h)
           std = torch.exp(0.5 * logvar)
           eps = torch.randn_like(std)
           z = mu + eps * std
           return z, mu, logvar
3365 class VAEDecoder(nn.Module):
        def __init__(self, latent_dim, output_dim, hidden_dim):
           super().__init__()
            self.fc1 = nn.Linear(latent_dim, hidden_dim)
           self.fc2 = nn.Linear(hidden_dim, output_dim)
        def forward(self, z):
           h = F.gelu(self.fcl(z))
           return self.fc2(h)
     class PerturbationEmbedding(nn.Module):
        def __init__(self, pert_dim, emb_dim):
           super().__init__()
           self.embedding = nn.Linear(pert_dim, emb_dim)
        def forward(self, pert):
           return self.embedding(pert)
3378 class HybridAttentionModel(nn.Module):
             __init___(self, input_dim, pert_dim, hidden_dim=512, n_layers=2, n_heads=8,
        def
            dropout=0.1, attention_dropout=0.1, ffn_dropout=0.1, activation='gelu',
            use_transformer=True, use_vae=False, vae_latent_dim=64, vae_hidden_dim=256,
use_pert_emb=False, pert_emb_dim=32, vae_beta=1.0):
            super(HybridAttentionModel, self).__init__()
           self.input_dim = input_dim
            self.pert_dim = pert_dim
           self.hidden dim = hidden dim
           self.use_transformer = use_transformer
           self.use_vae = use_vae
           self.vae_beta = vae_beta
           self.use_pert_emb = use_pert_emb
           if use_vae:
               self.vae_encoder = VAEEncoder(input_dim, vae_latent_dim, vae_hidden_dim)
               self.vae_decoder = VAEDecoder(vae_latent_dim, input_dim, vae_hidden_dim)
               expr_out_dim = vae_latent_dim
           else:
               self.expression_encoder = nn.Sequential(
                  nn.Linear(input_dim, hidden_dim),
                  nn.LayerNorm(hidden_dim),
                  nn.GELU(),
                  nn.Dropout (dropout),
                  nn.Linear(hidden_dim, hidden_dim),
                  nn.LayerNorm(hidden_dim),
                  nn.GELU(),
                  nn.Dropout (dropout)
               )
               expr_out_dim = hidden_dim
           if use_pert_emb:
               self.pert_encoder = PerturbationEmbedding(pert_dim, pert_emb_dim)
               pert_out_dim = pert_emb_dim
           else:
               self.pert_encoder = nn.Sequential(
                  nn.Linear(pert_dim, hidden_dim),
                  nn.LayerNorm(hidden_dim),
                  nn.GELU(),
                  nn.Dropout (dropout)
               pert_out_dim = hidden_dim
```

```
if use_transformer:
      encoder_layer = nn.TransformerEncoderLayer(
         d_model=expr_out_dim + pert_out_dim,
         nhead=n_heads,
         dim_feedforward=hidden_dim * 4,
         dropout=ffn_dropout,
         activation=activation,
         batch_first=True,
         norm_first=True
      )
      self.transformer = nn.TransformerEncoder(encoder_layer, num_layers=n_layers)
   else:
      self.self_attention = nn.MultiheadAttention(expr_out_dim + pert_out_dim,
          n_heads, dropout=attention_dropout, batch_first=True)
   self.fusion = nn.Sequential(
      nn.Linear(expr_out_dim + pert_out_dim, hidden_dim),
      nn.LayerNorm(hidden_dim),
      nn.GELU(),
      nn.Dropout (dropout)
   )
   self.output = nn.Sequential(
      nn.Linear(hidden_dim, hidden_dim),
      nn.LayerNorm(hidden_dim),
      nn.GELU(),
      nn.Dropout (dropout),
      nn.Linear(hidden_dim, input_dim)
   )
   self.perturbation_head = nn.Sequential(
      nn.Linear(hidden_dim, hidden_dim // 2),
      nn.LayerNorm(hidden_dim // 2),
      nn.GELU(),
      nn.Dropout (dropout),
      nn.Linear(hidden_dim // 2, pert_dim)
   )
   self.apply(self._init_weights)
def __init_weights(self, module):
   if isinstance(module, nn.Linear):
      nn.init.xavier_uniform_(module.weight)
      if module.bias is not None:
         nn.init.zeros_(module.bias)
   elif isinstance(module, nn.LayerNorm):
      nn.init.ones_(module.weight)
      nn.init.zeros_(module.bias)
def forward(self, x, pert):
   vae_kl = 0
   vae_recon = None
   if self.use_vae:
      z, mu, logvar = self.vae_encoder(x)
      vae_recon = self.vae_decoder(z)
      vae_kl = -0.5 * torch.sum(1 + logvar - mu.pow(2) - logvar.exp(), dim=1).mean()
      expr_feat = z
   else:
      expr_feat = self.expression_encoder(x)
   pert_feat = self.pert_encoder(pert)
   fusion_input = torch.cat([expr_feat, pert_feat], dim=1).unsqueeze(1)
   if self.use_transformer:
      x_trans = self.transformer(fusion_input).squeeze(1)
   else:
      x_trans, _ = self.self_attention(fusion_input, fusion_input, fusion_input)
      x_trans = x_trans.squeeze(1)
   fused = self.fusion(x_trans)
   output = self.output(fused)
   pert_pred = self.perturbation_head(fused)
```

3466       return output, pert_pred, vae_tecon, vae_k1         467	3465							
3467       3468       3469       3470       3471       3472       3473       3474       3475       3476       3477       3478       3481       3482       3483       3484       3485       3486       3487       3488       3489       3489       3490       3491       3492       3493       3494       3495       3496       3497       3498       3490       3491       3492       3493       3494       3495       3496       3497       3498       3490       3491       3492       3493       3494       3495       3496       3497       3498       3498       3499       3501       3502       3503       3504       3505       3506       3507       3508       3509       3510       3514 <th>3466</th> <th> return</th> <th>output,</th> <th>pert_pred,</th> <th>vae_recon,</th> <th>vae_kl</th> <th></th> <th></th>	3466	 return	output,	pert_pred,	vae_recon,	vae_kl		
348         3470         3471         3472         3473         3474         3475         3476         3477         3478         3480         3481         3482         3483         3484         3485         3486         3487         3488         3488         3489         3490         3491         3492         3493         3494         3495         3496         3497         3498         3499         3491         3492         3493         3494         3495         3500         3501         3502         3503         3504         3505         3506         3507         3508         3509         3501         3512         3513         3514         3515         3516 <tr< td=""><td>3467</td><td></td><td></td><td></td><td></td><td></td><td></td><td></td></tr<>	3467							
3470         3471         3472         3473         3474         3475         3476         3477         3478         3479         3481         3482         3483         3484         3485         3486         3487         3488         3489         3490         3491         3492         3493         3494         3495         3500         3501         3502         3503         3504         3505         3506         3507         3508         3509         3500         3501         3502         3503         3504         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516 <t< td=""><td>3468</td><td></td><td></td><td></td><td></td><td></td><td></td><td></td></t<>	3468							
3470         3471         3472         3473         3474         3475         3476         3477         3478         3479         3480         3481         3482         3483         3484         3485         3486         3487         3488         3488         3489         3490         3491         3492         3493         3494         3495         3502         3503         3504         3505         3506         3507         3508         3509         3501         3511         3512         3513         3514         3515         3516         3517         3518	3469							
3471         3472         3473         3474         3475         3476         3477         3478         3479         3481         3482         3483         3484         3485         3486         3487         3488         3489         3490         3491         3492         3493         3494         3493         3494         3493         3494         3493         3494         3495         3496         3497         3498         3499         3500         501         502         3503         3504         505         505         506         507         3508         509         510         511         512         513         514         515         516         5	3470							
3472         3473         3474         3475         3476         3477         3478         3479         3480         3481         3482         3483         3484         3485         3486         3487         3488         3489         3490         3491         3492         3493         3494         3495         3496         3499         3501         3502         3503         3504         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518	3471							
3473         3474         3475         3476         3477         3478         3479         3480         3481         3482         3483         3484         3485         3486         3487         3488         3489         3490         3491         3492         3493         3494         3493         3494         3495         3496         3497         3498         3499         3501         3502         3503         3504         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3516         3517         3518	3472							
3474         3475         3476         3477         3478         3480         3481         3482         3483         3484         3483         3484         3483         3484         3485         3486         3487         3488         3489         3490         3491         3492         3493         3494         3495         3496         3491         3492         3493         3494         3495         3496         3497         3498         3499         3500         3501         3502         3503         3504         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516 <t< td=""><td>3473</td><td></td><td></td><td></td><td></td><td></td><td></td><td></td></t<>	3473							
3475         3476         3477         3478         3480         3481         3482         3483         3484         3485         3486         3487         3488         3489         3490         3491         3492         3493         3494         3495         3500         3501         3502         3503         3504         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518	3474							
3476         3477         3478         3479         3481         3482         3483         3484         3485         3486         3487         3488         3489         3490         3491         3492         3493         3494         3495         3496         3497         3498         3499         3500         3501         3502         3503         3504         3505         3506         3507         3508         3509         3510         3512         3513         3514         3515         3516         3517         3518	3475							
3477         3478         3480         3481         3482         3483         3484         3485         3486         3487         3488         3489         3490         3491         3492         3493         3494         3495         3496         3497         3498         3499         3500         3501         3502         3503         3504         3505         3506         3507         3508         3509         3510         3511         3512         3514         3515         3516         3517         3518         3519         3511         3512         3513         3514         3515         3516         3517         3518         3519         3519         3519 <t< th=""><th>3476</th><th></th><th></th><th></th><th></th><th></th><th></th><th></th></t<>	3476							
3477       3478       3479       3480       3481       3482       3483       3484       3485       3486       3487       3488       3489       3490       3491       3492       3493       3494       3495       3496       3497       3498       3499       3501       3502       3503       3504       3505       3506       3507       3508       3509       3510       3511       3512       3513       3514       3515       3516       3517       3518	3477							
3479         3480         3481         3482         3483         3484         3485         3486         3487         3488         3489         3490         3491         3492         3493         3494         3495         3496         3497         3498         3499         3501         3502         3503         3504         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518	3477							
3449         3480         3481         3482         3483         3484         3485         3486         3487         3488         3489         3490         3491         3492         3493         3494         3495         3496         3497         3498         3499         3500         5501         3502         3503         3504         3505         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518	2470							
3480         3481         3482         3483         3484         3485         3486         3487         3488         3489         3490         3491         3492         3493         3494         3495         3496         3497         3498         3499         3500         3501         3502         3503         3504         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518	2400							
3481         3482         3483         3484         3485         3486         3487         3488         3489         3490         3491         3492         3493         3494         3495         3496         3497         3498         3499         3500         3501         3502         3503         3504         3505         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518	3480							
3482         3483         3484         3485         3486         3487         3488         3489         3490         3491         3492         3493         3494         3495         3496         3497         3498         3500         3501         3502         3503         3504         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518	3481							
3483         3484         3485         3486         3487         3488         3489         3490         3491         3492         3493         3494         3495         3496         3497         3498         3499         3500         3501         3502         3503         3504         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518         3519	3482							
3484         3485         3487         3488         3489         3490         3491         3492         3493         3494         3495         3496         3497         3498         3499         3500         3501         3502         3503         3504         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518	3483							
3485         3486         3487         3488         3489         3491         3492         3493         3494         3495         3496         3497         3498         3499         3500         3501         3502         3503         3504         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518	3484							
3480         3487         3488         3490         3491         3492         3493         3494         3495         3496         3497         3498         3500         3501         3502         3503         3504         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518	3485							
3487         3488         3490         3491         3492         3493         3494         3495         3496         3497         3498         3499         3500         3501         3502         3503         3504         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518	3486							
3488         3489         3490         3491         3492         3493         3494         3495         3496         3497         3498         3499         3500         3501         3502         3503         3504         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518         3519	3487							
3489         3490         3491         3492         3493         3494         3495         3496         3497         3498         3499         3500         3501         3502         3503         3504         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518         3519	3488							
3490         3491         3492         3493         3494         3495         3496         3497         3498         3499         3500         3501         3502         3503         3504         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518         3519	3489							
3491         3492         3494         3495         3496         3497         3498         3499         3500         3501         3502         3503         3504         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518         3519	3490							
3492         3493         3494         3495         3496         3497         3498         3499         3500         3501         3502         3503         3504         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518         3519	3491							
3493         3494         3495         3496         3497         3498         3499         3500         3501         3502         3503         3504         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518         3519	3492							
3494         3495         3496         3497         3498         3499         3500         3501         3502         3503         3504         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518         3519	3493							
3495         3496         3497         3498         3499         3500         3501         3502         3503         3504         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518         3519	3494							
3496         3497         3498         3499         3500         3501         3502         3503         3504         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518         3519	3495							
3497         3498         3499         3500         3501         3502         3503         3504         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518         3519	3496							
3498         3499         3500         3501         3502         3503         3504         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518         3519	3497							
3499         3500         3501         3502         3503         3504         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518         3519	3498							
3500         3501         3502         3503         3504         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518         3519	3499							
3501         3502         3503         3504         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518         3519	3500							
3501         3502         3503         3504         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518         3519	3501							
3503         3504         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518         3519	3502							
3504         3505         3506         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518         3519	3503							
3505 3506 3507 3508 3509 3510 3511 3512 3513 3514 3515 3516 3517 3518 3519	3504							
3505 3506 3507 3508 3509 3510 3511 3512 3513 3514 3515 3516 3517 3518 3519	3505							
3500         3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518         3519	2506							
3507         3508         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518         3519	2507							
3500         3509         3510         3511         3512         3513         3514         3515         3516         3517         3518         3519	2500							
3509         3510         3511         3512         3513         3514         3515         3516         3517         3518         3519	3308							
3510         3511         3512         3513         3514         3515         3516         3517         3518         3519	3509							
3511         3512         3513         3514         3515         3516         3517         3518         3519	3510							
3512         3513         3514         3515         3516         3517         3518         3519	3511							
3513 3514 3515 3516 3517 3518 3519	3512							
3514 3515 3516 3517 3518 3519	3513							
3515 3516 3517 3518 3519	3514							
3516 3517 3518 3519	3515							
3517 3518 3519	3516							
3518 3519	3517							
3519	3518							
	3519							

# 3520 H. Detailed outputs from other research Agents

-	Task Descintion 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).
	Dataset Description:
	- Source: Norman et al., 2019
	- Cell Type: Human K562 leukemia cells
	- Perturbations: CRISPRI targeting 105 single genes and 131 gene pairs
	perturbed conditions
	Please search the following websets for more detailed information:
	- https://www.ncbi.nlm.nih.gov/geo/guery/acc.cgi?acc=GSE133344
	- https://www.science.org/doi/10.1126/science.aax4438
	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 Parturbations, Predict effects of gene parturbations not present during
	training
	2. Unseen Cell Contexts: Predict responses in cells with gene expression profiles
	not observed during training
	Evaluation Matrice.
	- 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.
	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.
	Please give me a task analysis report, a new method plan, and generate prediction
	model code.
)ne	n & I DeenResearch Outnut
γpc 1	Deer Deserverkh Output
	Deepkesearchn Output
	## Tack Analysis, Data Characteristics and Challenges
	## TASK MHATASTS: DALA CHALACLELISTICS 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

- 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.

3572 3573 3574

3564

3565

3566

3567

3568

3569

3570

3575	
3576	- **Perturbation Types:** Perturbations consist of **single-gene knockdowns** (
2577	CRISPRi) and **two-gene combinations**, enabling the study of genetic
5577	interactions. For each target gene, multiple sgRNAs were used (to control for
3578	surve-specific effects), and for two-gene compos, certs co-express two
3579	Synnio (one por cargoo).
3580	**Modeling Challenges:** Developing an accurate predictive model for this task
3581	entails several significant challenges:
3582	- **Single-Cell Variability:** Even for the same perturbation, different cells
3593	can respond with variability in magnitude of gene expression changes. Each
2501	cell has an intrinsic baseline expression profile reflecting its state (cell
2505	cycle stage, stochastic gene expression, etc.).
3585	- **Generalization to Unseen Perturbations:** A critical use-case is predicting
3586	outcomes for perturbations that were **never seen in training**. This
3587	includes:
3588	- *New target genes:* The model may be asked to predict the effect of
3589	knocking down a gene that was not perturbed in the training data. This
3590	example for that genes perturbation. The model must leverage other
3501	information (e.g. the genes baseline expression or relationships to
2502	other genes) to infer its likely impact.
2502	- *New gene combinations:* Similarly, the model might see single-gene
3393	perturbations during training, but not a particular pair in combination.
3594	before). As reported in prior work, when one or both perturbed genes
3595	are unseen, prediction difficulty increases markedly. The model needs a
3596	strategy to generalize **inductively**, combining knowledge of
3597	individual genes to predict combined effects.
3598	- ++Accuracy vs Interpretability.++ We prioritize predictive accuracy. This
3599	means we are willing to employ a complex, black-box model (deep neural
3600	network) with many parameters if it yields better predictive performance,
2601	even if its not easy to interpret biologically. This choice acknowledges that
2602	the gene regulatory relationships underlying the perturbation responses are
3002	complex; a simpler interpretable model (like linear regression per gene)
3603	as deep model. However, the complexity of the model also raises the risk
3604	of overfitting and requires careful regularization and validation.
3605	
3606	In summary, the task is essentially learning a high-dimensional function: **i**(*
3607	the underlying gene network responses from data, handle combinatorial
3608	perturbation effects, and generalize to new genes and cell states. These
3609	challenges will guide the modeling decisions discussed below.
3610	
2611	
2612	## Proposed Modeling Approach for Accuracy and Generalization
3612	### Model Architecture and Input Representation
3613	
3614	inputs (1) the baseline gene expression profile of an upperturbed coll and
3615	(2) the identity of the target gene (s) perturbed and outputs the predicted
3616	post-perturbation gene expression profile. The architecture is designed to
3617	flexibly capture non-linear gene interactions and make use of prior gene
3618	relationship knowledge to generalize to unseen perturbations. The key
3610	components are:
3620	- **Baseline Expression Encoder:** A front-end network that processes the
2621	baseline expression vector (dimension = $*G*$ genes). Directly using the raw
3021	high-dimensional vector as input to a fully-connected network is feasible but
3622	autoencoder or dimensionality reduction layer) that compresses the baseline
3623	gene expression profile into a **lower-dimensional latent representation**.
3624	For example, a few fully-connected layers with ReLU activation can reduce the
3625	~20k-dimensional input to a dense ~512-dimensional embedding. This latent
3626	vector is intended to capture the cells overall state or context (e.g., if
3627	the cell is in a high-profileration state, or has high expression of certain pathways etc.) By encoding the baseline, the model can later modulate
3628	pachways, ecc.,. by chooding the baseline, the model can later modulate
3620	
JU47	

3630	norturbation offects depending on these latent features
3631	percurbación erreccs depending on chese facenc feacures.
3632	- **Perturbation Encoder:** We represent the **perturbation identity** in a way
3633	that the model can easily utilize and generalize. Each target gene (from the
3634	- a one-hot vector (of length equal to the number of target genes) if a
3635	single gene is perturbed, or a multi-hot vector if multiple genes are
3636	perturbed (e.g. for a two-gene perturbation, the vector has ones in the
3637	positions corresponding to the two targeted genes). This binary
3638	lookup table or a small fully-connected network) to produce a **
3630	perturbation embedding**. The embedding is a continuous vector (e.g.
3640	128-dimensional) that represents the effect of the perturbation in a
36/1	
2642	maintain a trainable embedding vector for each gene in the target set,
2642	and for a combination perturbation, we combine the embeddings of the
2643	individual genes (e.g. by summation or an attention mechanism). Using a
3044	each genes characteristic perturbation impact. For a multi-gene
3645	perturbation, a simple summation assumes independence of effects, while
3646	a more sophisticated combination (see below) can capture interactions.
3647	- **Combination Module:** The baseline context and perturbation effect must be
3648	integrated. We concatenate the baseline latent vector and the perturbation
3649	embedding vector into a single combined latent representation. This combined
3650	vector (or rengen 640 in our example, if paseline latent is 512 and perturbation embedding 128) now contains information about where the cell
3651	started and what perturbation was applied. This is passed through further
3652	layers (a **fusion network**) to compute the output. For instance, a
3653	multilayer perceptron (MLP) with one or two hidden layers (e.g., 512 neurons,
3654	interactions between cell state and perturbation e.g., the effect of
3655	perturbing gene X might depend on the level of gene $\check{ ext{Y}}$ in the baseline state,
3656	which a multiplicative interaction in the MLP can learn.
3657	- **Output Layer (Prediction Head):** The final layer of the network produces a
3658	vector of length *G* (the number of genes), which is the predicted post-
3659	fact that many genes dont change, we design the output to predict a **change
3660	(delta) from baseline** for each gene rather than an absolute expression. In
3661	practice, the network can output expression for each gene, and then this is
3662	added to the baseline input to yield the final predicted expression:
3663	baseline} + \Delta {\text{model}} (E {\text{baseline}}, P).
3664	This formulation makes it easier for the model to output zero for genes
3665	that should remain unchanged, and focuses on learning the deviations.
3666	It also implicitly grounds the prediction in the baseline level (so if a gene is high at baseline and not affected by the perturbation the
3667	model can just output a near-zero change).
3668	uNer linear Interaction Medeling
3669	- **NON-IINEAR INTERACTION MODELING:** White a simple concatenation of embeddings treats multi-gene perturbations roughly as an additive combination of single
3670	effects, we can enhance the model to capture **genegene interaction effects
3671	**. One idea is to use an **attention mechanism or gating** in the
3672	perturbation encoder: for example, if two genes A and B are perturbed,
3672	network that can learn a pairwise interaction term. Another approach is to
2674	include pairwise products of gene embeddings in the combined feature (
2675	allowing the network to learn a unique contribution for the pair *A&B* beyond
2676	A + B). Given that Norman *et al.* tested primarily pairwise perturbations,
2677	of genes in the training set to capture any deviation from additivity.
30//	However, to generalize to unseen pairs, a better strategy is to learn a **
36/8	function** for combining embeddings (like attention) rather than a fixed
3679	lookup for each pair.
3680	- **Incorporating Prior Knowledge (for Generalization):** To improve inductive
3681	generalization to unseen genes, we can draw inspiration from GEARS and
3682	similar methods, we could initialize of regularize the gene perturbation
3683	
3684	

3685	omboddings using oxtornal knowlodge.
3686	- Use a **gene co-expression network** (computed from the baseline single-
3687	cell data or external data) as a graph, and pass gene embeddings
3688	through a Graph Neural Network (GNN) layer. This encourages genes that
3689	have similar roles or expression patterns to have embeddings that
3600	produce similar effects. Thus, if an unseen gene has a similar co-
2601	expression profile to a seen gene, the model might infer similar perturbation outcomes
3091	- Use **pathway or GO (Gene Ontology) information** to place genes in a
3692	relational graph (as GEARS did with a GO-derived graph for perturbation
3693	embeddings). Two genes in the same pathway might be expected to
3694	produce related downstream effects; by training on the known genes, the
3695	Technically, this can be done by adding a loss term that encourages
3696	the learned embedding to correspond to the genes position in the
3697	knowledge graph, or by a GNN that propagates influence from neighbors
3698	during training.
3600	- These additions make the model more complex but aim to imbue it with **
2700	vield similar perturbation responses and (ii) genes in related
3700	pathways affect overlapping sets of genes when perturbed. We will
3701	prioritize implementing a simpler version (like using co-expression PCA
3702	or clusters to initialize embeddings), and note that full graph-based
3703	learning could further improve generalization if needed.
3704	- **Residual Connections and Regularization:** We will include skip-connections
3705	wherever helpful (for example, the baseline input could be fed not only into
3706	the encoder but also concatenated directly to a later layer, or the output
3707	head could directly see the raw baseline as well, ensuring the model can
3708	dropout in the MLP layers, L2 weight decay, or even an auxiliary loss to
3709	reconstruct the baseline (autoencoder style) can be employed to prevent
3710	overfitting and encourage the model to learn meaningful latent features
3710	rather than memorizing training examples.
2712	In essence, the architecture is a **conditioned deep neural network**: it
2712	conditions on the cells initial state and the perturbation, and produces an
3/13	output state. This is somewhat analogous to an encoderdecoder model where the
3714	encoder is the baseline expression and the condition or control signal is the perturbation identity. Because accuracy is paramount, we allow a fairly
3715	large model with enough capacity to capture complex gene regulatory responses.
3716	
3717	### Training Strategy for Accuracy and Generalization
3718	
3719	With the architecture in place, we next focus on **training methodology**, as it
3720	greatly affects model generalization and performance.
3721	- **Training Data Construction:** We will pair each perturbed cells data with a
3722	baseline profile as input. Since in the actual experiment we typically do not
3723	nave a *pre-percurpation* measurement of the same Cell, We have to simulate a baseline. We can use the expression profiles of control cells (non-
3724	targeting sgRNA) as proxies for baseline input. For each perturbed cell in
3725	the training set, we can randomly sample a control cells expression as the
2726	baseline input. This effectively assumes that any control cell is an example
3720	of an unperturbed state the perturbed cell *could* have come from. Over many samples, the model will learn to man from an average baseline state to the
3727	perturbed outcome. We can further refine this pairing by matching on cell
3728	state: e.g., cluster the baseline cells by expression and pick a baseline
3729	from the same cluster as the perturbed cells profile (minus the perturbation
3730	effect) to provide a closer starting point. However, random pairing with a
3731	overfit a one-to-one mapping.
3732	control a one co one mapping.
3733	- *Unseen cell context generalization:* By exposing the model to many
3734	allierent baseline samples paired with a given perturbation outcome (
3735	for the same perturbation. This should improve the models robustness to
3736	any particular baseline context and enable generalization to new
2727	baseline profiles. Essentially, the model sees that the same
2121	perturbation can apply to various starting expression patterns.
3/38	
5139	

3740	
3741	- **Loss Functions:** The primary loss will be **Mean Squared Error (MSE)**
3742	between the predicted and actual post-perturbation expression vectors. To
3743	ensure we adequately learn the important changes, we can modify the loss in
2744	two ways:
5744	- Compute a weighted MSE that gives higher weight to genes that are truly
3/45	differentially expressed in that training example. For instance, if we
3746	know gene j changed significantly in the real perturbed cell (compared
3747	to baseline or compared to controls), we can upweight the error on gene
3748	j for that sample. This forces the model to focus on fitting the genes
3749	that move, rather than being dominated by the many near-zero changes.
3750	- Alternatively, we can rai trands then fine ture the model on just the
2751	top-k DE genes for each perturbation (or using a loss like contrastive
3751	that emphasizes getting the direction of change correct).
3752	
3753	In practice, a simpler approach is to stick with standard MSE on all genes
3754	the model doesn't ignore those signals. If we see the model predicting
3755	trivial (no-change) solutions, we will adjust the loss weighting.
3756	
3757	- **Optimizer and Regularization:** We will use **Adam optimizer** (adaptive
3758	realizing face) which is well-surface for training deep networks on possibly noisy data. A relatively small learning rate ( $\alpha$ , $\beta$ , $\beta$ , $\beta$ , $\beta$ , $\beta$ ) will be
3750	used and well monitor validation loss for convergence. Early stopping on the
5739	validation MSE/PCC can prevent overfitting. Dropout layers (e.g. dropout rate
3760	0.2–0.5) can be inserted in the MLP to regularize. Weight decay (L2) will
3761	help keep embedding weights from growing excessively (important if some genes
3762	cause large effects, we dont want the weights to explode to fit one-hot
3763	Signals).
3764	- **Generalization Enhancements:** To specifically improve performance on **
3765	unseen genes**:
3766	- We will pre-train the **perturbation encoder embeddings**: one idea is to
2767	<code>train</code> a simple autoencoder or matrix factorization on the <code>*observed</code>
3707	perturbation effects $\star$ in training. For example, compute the average
3768	expression change each perturbation causes on each gene (from data),
3769	matrix to get a low-dimensional embedding for genes and for
3770	perturbations. If a gene is unseen, we can still derive a gene
3771	embedding from its baseline expression pattern using the gene embedding
3772	from baseline encoder. This is a bit complex, but essentially we want
3773	the model to know something about gene function even if not perturbed.
3774	Incorporating gene ontology or known regulator information for unseen
2775	genes could also inform an initial embedding (e.g., if an unseen gene
5775	similar perturbation embedding).
3776	
3777	By combining a carefully-designed architecture with these training strategies,
3778	depending to new perturbations and contexts** In particular success would
3779	be a model that, given an arbitrary K562 baseline profile and a novel target
3780	gene, can predict a biologically plausible outcome capturing which pathways
3781	might be up or down-regulated by that genes perturbation, even without
3782	having seen it before. The inclusion of baseline input enables context-
3783	specific predictions (so the same gene perturbation could have different effects in different baseline states if the model learns those interactions)
2704	and the use of gene embeddings with potential graph-based inductive bias
2705	enables inductive generalization across genes.
5/85	North we present the implementation of this redship Duranch struct and for ODV
3786	Next, we present the implementation of this model in Pylorch, structured for GPU
3787	
3788	## DuTorch Implementation (Model & Training Code)
3789	## Fyroren imprementation (Model & fraining code)
3790	Below is the complete PyTorch code for the proposed model and an example training
3791	loop. This code assumes that the dataset has been preprocessed into tensors:
3702	'X_train' (baseline expressions), 'Y_train' (perturbed expressions), and '
2702	r_train (perturbation indicator vectors). The model is designed to run on a
3793	
3794	

```
3795
             GPU for efficiency.
3796
          '''python
3797
          import torch
3798
          import torch.nn as nn
3799
          import torch.optim as optim
3800
          # Define device for GPU acceleration
3801
          device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
3802
          print("Using device:", device)
3803
          # Example shapes (for illustration; replace with actual dataset dimensions)
3804
          num_genes = X_train.shape[1] # dimensionality of gene expression (e.g. ~20000)
          num_target_genes = P_train.shape[1] # dimensionality of perturbation vector (e.g.
3805
               100)
3806
          # 1. Define the model architecture
          class PerturbationPredictor(nn.Module):
3808
             def __init__(self, input_dim_expr, input_dim_pert,
3809
                       hidden_dim_expr=512, hidden_dim_pert=128, hidden_dim_comb=512):
                ....
3810
                input_dim_expr: number of genes (baseline expression vector length)
3811
                input_dim_pert: number of possible targeted genes (length of perturbation
                indicator vector)
3812
3813
                super(PerturbationPredictor, self)._
                                                      _init__()
3814
                # Baseline expression encoder network
3815
                self.expr_encoder = nn.Sequential(
                   nn.Linear(input_dim_expr, hidden_dim_expr),
3816
                   nn.ReLU(),
3817
                   nn.Linear(hidden_dim_expr, hidden_dim_expr),
3818
                   nn.ReLU()
3819
                # Perturbation encoder network
3820
                self.pert_encoder = nn.Sequential(
3821
                   nn.Linear(input_dim_pert, hidden_dim_pert),
                   nn.ReLU(),
3822
                   nn.Linear(hidden_dim_pert, hidden_dim_pert),
3823
                   nn.ReLU()
3824
                # Fusion network to combine baseline and perturbation representations
3825
                self.comb_fc1 = nn.Linear(hidden_dim_expr + hidden_dim_pert,
3826
                   hidden_dim_comb)
                self.comb_act = nn.ReLU()
self.comb_fc2 = nn.Linear(hidden_dim_comb, input_dim_expr) # outputs delta
3827
3828
                    for each gene
3829
             def forward(self, baseline_expr, perturb_vec):
3830
                # baseline_expr: Tensor of shape (batch_size, input_dim_expr)
3831
                # perturb_vec: Tensor of shape (batch_size, input_dim_pert)
3832
                # Encode baseline and perturbation
                z_expr = self.expr_encoder(baseline_expr) # shape: (batch, hidden_dim_expr)
3833
                z_pert = self.pert_encoder(perturb_vec) # shape: (batch, hidden_dim_pert)
3834
                # Combine encoded representations
                z_comb = torch.cat([z_expr, z_pert], dim=1) # shape: (batch,
3835
                   hidden_dim_expr+hidden_dim_pert)
3836
                h = self.comb_act(self.comb_fcl(z_comb)) # hidden layer
3837
                delta = self.comb_fc2(h) # predicted change in expression
3838
                # Add the predicted change (delta) to the baseline to get final prediction
                pred_expr = baseline_expr + delta
3839
                return pred_expr
3840
          # Instantiate the model and move to GPU
3841
          model = PerturbationPredictor(num_genes, num_target_genes).to(device)
3842
          print(model)
3843
          # 2. Set up loss function and optimizer
3844
          criterion = nn.MSELoss() # Mean Squared Error loss
3845
          optimizer = optim.Adam(model.parameters(), lr=1e-3)
3846
          # 3. Training loop (example)
3847
          num_epochs = 20
3848
3849
```

```
3850
          batch size = 256
3851
          # If data is large, use DataLoader for batching (assuming X_train, P_train,
3852
              Y_train are TensorDatasets)
3853
          train_dataset = torch.utils.data.TensorDataset(X_train, P_train, Y_train)
3854
          train_loader = torch.utils.data.DataLoader(train_dataset, batch_size=batch_size,
              shuffle=True)
3855
3856
          for epoch in range(1, num_epochs+1):
3857
             model.train()
             total_loss = 0.0
3858
             for X_batch, P_batch, Y_batch in train_loader:
3859
                 # Move data to device
                X_batch = X_batch.to(device)
3860
                P_batch = P_batch.to(device)
                Y_batch = Y_batch.to(device)
                 # Forward pass
                pred = model(X_batch, P_batch)
                 loss = criterion(pred, Y_batch)
3864
                 # Backpropagation
3865
                optimizer.zero_grad()
3866
                loss.backward()
                optimizer.step()
                total_loss += loss.item() * X_batch.size(0)
3868
             avg_loss = total_loss / len(train_dataset)
             if epoch % 5 == 0 or epoch == 1:
                print(f"Epoch {epoch:02d}: Training MSE = {avg_loss:.4f}")
3870
3871
          # 4. Evaluation on test set
          model.eval()
3872
          with torch.no_grad():
3873
             X_test = X_test.to(device)
3874
             P_test = P_test.to(device)
             Y_test = Y_test.to(device)
3875
             pred_test = model(X_test, P_test)
3876
             mse_full = criterion(pred_test, Y_test).item()
3877
          # Compute Pearson correlation coefficient (PCC) for each test sample
          pred_np = pred_test.cpu().numpy()
3878
          Y_np = Y_test.cpu().numpy()
3879
          baseline_np = X_test.cpu().numpy()
          pcc_list = []
          topk = 20
3881
          topk_pcc_list = []
3882
          for i in range(Y_np.shape[0]):
3883
             true_expr = Y_np[i]
             pred_expr = pred_np[i]
3884
             # PCC for all genes
3885
             cov = np.cov(true_expr, pred_expr, bias=True)
             # cov matrix 2x2: cov[0,1] is covariance between true and pred
             pcc = cov[0,1] / (np.std(true_expr) * np.std(pred_expr) + 1e-8)
3887
             pcc_list.append(pcc)
3888
             # PCC for top-k differentially expressed genes
3889
             # Identify top-k genes by absolute change in true expression vs baseline
             base_expr = baseline_np[i]
3890
             diff = np.abs(true_expr - base_expr)
3891
             topk_idx = np.argsort(diff)[-topk:]
3892
             if topk > 0:
                 true_top = true_expr[topk_idx]
3893
                pred_top = pred_expr[topk_idx]
3894
                cov_top = np.cov(true_top, pred_top, bias=True)
topk_pcc = cov_top[0,1] / (np.std(true_top) * np.std(pred_top) + 1e-8)
3895
                 topk_pcc_list.append(topk_pcc)
3896
          # Calculate mean metrics
3897
          mean_pcc = float(np.mean(pcc_list))
          mean_topk_pcc = float(np.mean(topk_pcc_list))
          print(f"Test MSE (all genes): {mse_full:.4f}")
3899
          print(f"Test mean PCC (all genes): {mean_pcc:.3f}")
print(f"Test mean PCC (top-{topk} DE genes): {mean_topk_pcc:.3f}")
3900
3901
3902
          **Explanation of the Code:**
3903
3904
```

3905	
3906	- We define a 'PerturbationPredictor' model class that implements the
3907	architecture described. The baseline expression encoder ('expr_encoder') and
3908	perturbation encoder (`pert_encoder`) are simple feed-forward networks. These
3909	could be extended or replaced with more complex sub-networks (e.g., adding
3010	aropout, or using a graph convolution in 'pert_encoder' if incorporating a
2011	comb fc1' and 'comb fc2') with a Bell in between The output of 'comb fc2' is
3911	a vector of length equal to number of genes, representing the predicted
3912	change in expression. We then add this to the baseline ('pred_expr = baseline
3913	+ delta`) to get the final prediction.
3914	- We move the model to GPU (`.to(device)`). The training loop iterates over mini-
3915	batches of data (using a DataLoader for efficiency). For each batch, we do a
3916	loss against the true perturbed expression. We then backpropagate ()loss
3017	backward() ) and update weights with Adam optimizer. We print the training
2010	MSE every few epochs for monitoring.
3918	- In the evaluation section, after training, we compute metrics on $X_{test}$ ,
3919	P_test, Y_test`. We calculate the overall MSE on test (as `mse_full`). Then
3920	between predicted and true expression. We do this manually using numpy:
3921	computing the covariance and standard deviations (note: one could also use '
3922	scipy.stats.pearsonr' or torchmetrics for this). We also compute **top-20 PCC
3923	$\star\star$ by finding the 20 genes with largest absolute change in the true data (
3924	compared to baseline) and computing PCC on that subset. Finally, we print the
2025	average PCC across all test cells, and the average top-20 PCC.
2026	**Note:** In practice, you would likely use vectorized operations or library
3926	functions to compute PCC for all cells at once. Also, you might evaluate the
3927	model separately on the different generalization scenarios (e.g., one test
3928	set for unseen genes, one for seen genes, etc.) and compare performance. The
3929	code above provides a template that can be adapted for those purposes.
3930	By executing this model training pipeline, we expect to obtain a model that can
3931	accurately predict gene expression changes in K562 cells after CRISPRi
3932	perturbations. Its performance will be reflected in low MSE (especially on
2022	perturbed genes) and high correlation with actual profiles, even for
2922	demonstrates a powerful approach to **in silico** predict cellular response
3934	to genetic perturbations, leveraging modern deep learning (PyTorch with GPU)
3935	for both precision and generalization. The combination of baseline context
3936	and perturbation identity as inputs allows the model to capture conditional
3937	dependencies (how a perturbations effect can vary by cell state), addressing
3938	one of the key chartenges in predictive biology.
3939	
3940	Perplexity DeepResearch
3941	DeepResearchh Output
3942	
3943	
3944	Task Analysis and Technical Challenges
3945	T Distantest of CDICDDi Dauturbations
3946	I. BIOLOGICAL CONTEXT OF CRISPRE PERTUrbations
3947	The Norman et al. dataset captures transcriptional consequences of CRISPRi-
3948	mediated gene repression through single-cell profiling of:
30/0	- 105 individual gene knockouts
2050	- 131 pairwise gene interactions[^1][^4]
3950	- ~90,000 cells with baseline (unperturbed) and perturbed states [^4]
3951	
3952	key plological considerations:
3953	1. **Non-linear genetic interactions**: Pairwise perturbations exhibit emergent
3954	effects not predictable from individual gene knockouts[^2]
3955	2. **Cell-state heterogeneity**: Native transcriptional variability in K562
3056	populations influences perturbation responses [2] [4]
3057	distributions differ between perturbed/control cells[^4][^6]
2059	
3938	
3959	
3960	II. Machine Learning Challenges
----------------------	--
3962	1. **High-dimensional output space**: Predicting 20,000+ gene expressions
3963	simultaneously
3964	<pre>2. **Sparse perturbation coverage**: Only 0.0% of possible gene pairs tested in original study[^1]</pre>
3965	3. **Context-dependent effects**: Perturbation responses vary with baseline
3966	4. **Combinatorial explosion**: Quadratic growth of possible gene interactions
3967	with target count
3969	**Evaluation Protocol Design**
3970	**Train/Test Splits**:
3971	- **Unseen perturbations**: Hold out 20% of single-gene and 30% of pairwise
3972 3973 3074	<pre>- **Unseen contexts**: Stratify by baseline PCA clusters, hold out entire cell- state subpopulations</pre>
3975	
3976	Methodological Framework: GenePerturbNet
3977	I. Architectural Overview
3978	**Core Components**:
3980	1. **Baseline Encoder**: Variational autoencoder mapping unperturbed expression
3981	<pre>\$x\$ to latent \$z\$.</pre>
3982 3983	2. **Perturbation Transformer**: Attention network processing target genes \$g\$ into perturbation vector \delta.
3984 3985	3. **Interaction Module**: Graph neural network modeling gene-gene interaction effects.
3986 3987	4. **Response Decoder**: generating perturbed expression profile
3988	II. Training Strategy
3989	**Loss Function:**
3990 3991 3002	<pre>\$\mathcal{L} = ELBO(x,y) + \lambda_1 MSE(y, \hat{y}) + \lambda_2 \mathcal{L}_{     contrast}\$</pre>
3993	Where:
3994	- **Contrastive Loss**:
3995 3996	<pre>\$\mathcal{L}_{contrast} = -\log\frac{\exp(s(y,\hat{y})/\tau)}{\sum_{y'}\exp (s(y',\hat{y})/\tau)}\$</pre>
3997	- \$s()\$: Cosine similarity
3998	- \$\tau\$: Temperature parameter
4000 4001	IV. Curriculum Learning Schedule: 1. Phase 1: Pretrain on single-gene perturbations
4001	2. Phase 2: Fine-tune with pairwise interactions
4003	3. Phase 3: Joint optimization with contrastive loss
4004	
4005	Implementation and Benchmarking
4007	I. Code Implementation
4008	'''python
4009	import scvi
4010	import torch from torch import nn
4011	from scvi import REGISTRY_KEYS
4013	from scvi.module.base import BaseModuleClass
4014	

```
4015
4016
          class GenePerturbNet(scvi.model.SCVI):
4017
             def __init__(self, adata, n_hidden=256, n_latent=32,
                       perturbation_embedding=64, **kwargs):
4018
                super().__init__(adata, **kwargs)
4019
                self.pert_embed = nn.Embedding(adata.n_vars, perturbation_embedding)
                self.transformer = nn.TransformerEncoderLayer(
4020
                    d_model=perturbation_embedding, nhead=4
4021
4022
                self.gnn = GeneInteractionGNN(perturbation_embedding)
4023
                self.decoder = nn.Sequential(
4024
                   nn.Linear(n_latent + perturbation_embedding, n_hidden),
4025
                   nn.ReLU(),
                   nn.Linear(n_hidden, adata.n_vars)
4026
                )
4027
             def forward(self, x, perturbations, *args):
    qz_m, qz_v, z = self.encoder(x)
4028
4029
                pert_emb = self.pert_embed(perturbations)
4030
                delta = self.transformer(pert_emb)
4031
4032
                z_pert = z + delta.mean(dim=1)
4033
                px_scale, px_r, px_rate, px_dropout = self.decoder(z_pert)
4034
                return px_rate, qz_m, qz_v, z_pert
4035
          class GeneInteractionGNN (nn.Module):
4036
             def __init__(self, embed_dim):
4037
                super().__init__()
4038
                self.conv1 = nn.Conv1d(embed_dim, embed_dim*2, kernel_size=3)
                self.attn = nn.MultiheadAttention(embed_dim, num_heads=4)
4039
4040
             def forward(self, gene_pairs):
                x = self.conv1(gene_pairs)
x, _ = self.attn(x, x, x)
4041
4042
                return x.mean(dim=1)
4043
          ...
4044
4045
          II. Training Protocol
4046
          '''python
4047
          adata = sc.read("norman_2019.h5ad")
          scvi.model.setup_anndata(adata, layer="counts",
4048
                             categorical_covariate_keys=["perturbation"])
4049
4050
          model = GenePerturbNet(adata, n_latent=32, perturbation_embedding=64)
          model.train(max_epochs=100, batch_size=256,
4051
                  plan_kwargs={"contrastive_weight":0.5})
4052
          control_cells = adata[adata.obs["perturbation"] == "control"]
4053
          perturbed_pred = model.predict(control_cells, targets=["geneA", "geneB"])
4054
4055
          III. Conclusion and Future Directions
4056
4057
          This framework demonstrates that deep learning architectures explicitly modeling
              genetic interactions can significantly improve perturbation prediction
4058
              accuracy. Key recommendations:
4059
          1. **Data Expansion**: Incorporate epigenetic data to better model CRISPRi
4060
             efficiency
4061
          2. **Dynamic Modeling**: Add temporal component to capture perturbation kinetics
4062
          3. **Transfer Learning**: Adapt model to other cell types using few-shot learning
4063
          The provided codebase offers a foundation for extending perturbation prediction
4064
              to novel experimental designs and therapeutic applications.
4065
4066
4067
4068
     Gemini DeepResearch
4069
```

4070	DeepResearchh Output
4071	
4072	Task Analysis
4073	
4074	# Dataset Description
4075	The foundation of this predictive modeling task lies in the comprehensive dataset
4076	generated by Norman et al. and published in Science in 2019. This study
4077	cell phenotypes using a Perturb-seq approach.15 The experiments were
4078	conducted on human K562 leukemia cells, a commonly used model system in
4079	hematological research. The study employed CRISPRi to target a total of 236 genetic perturbations encompassing 105 single genes and 131 pairs of genes
4080	The scale of the dataset is substantial, comprising approximately 90,000
4081	single-cell RNA-seq profiles, which includes both control cells and cells
4082	subjected to various single and combinatorial CRISPRi perturbations.
4083	The dataset's focus on both single and paired gene perturbations is crucial for
4084	training a model capable of predicting combinatorial effects, which are often
4085	information about the effect of individual gene silencing. Paired
4086	perturbations allow the model to learn about genetic interactions, where the
4087	combined effect of perturbing two genes is different from the sum of their individual effects. This is essential for conturing the per-linear neture of
4088	biological systems.
4089	
4090	# Problem Definition
4091	The primary objective is to develop a predictive model that takes as input the
4092	baseline gene expression profile of an unperturbed K562 cell and the identity
4093	of the target gene(s) for CRISPRi perturbation. The desired output of the
4094	model is the predicted gene expression profile of the same cell after the specified perturbation has been applied
4095	specified percubación has been appried.
4090	The model's performance will be evaluated under two key scenarios to assess its generalizability:
4098	1. **Unseen Perturbations:** The model should be able to accurately predict the
4099	effects of CRISPRi targeting genes or gene pairs that were not included in
4100	learned knowledge to novel genetic manipulations.
4101	2. **Unseen Cell Contexts:** The model should be capable of predicting the
4102	response to a perturbation in cells with baseline gene expression profiles
4103	robustness to the inherent heterogeneity within the K562 cell population.
4105	The evaluation scenarios highlight the need for a model that can generalize
4106	beyond the specific perturbations and cell states seen during training, which
4107	is a significant challenge in piological systems due to their complexity and inherent variability. A model trained only on observed perturbations might
4108	not accurately predict the effects of novel perturbations or responses in
4109	slightly different cellular contexts. Evaluating on unseen data is crucial
4110	Tor assessing the model o robustness and rear-world applicability.
4111	# Challenges and Considerations
4112	Developing an accurate predictive model for gene expression profiles following
4113	CRISPRi perturbation presents several significant challenges stemming from
4114	the inherent complexity of biological systems and the nature of single-cell data.
4115	
4116	Ine gene expression data obtained from scRNA-seq experiments is inherently high- dimensional, with measurements for thousands of genes in each individual cell
4117	This high dimensionality necessitates the use of sophisticated modeling
4118	techniques capable of handling a large number of input features. Furthermore,
4119	the underlying biological networks that govern gene expression are incredibly complex, involving intricate regulatory relationships between
4120	genes, transcription factors, and other cellular components.
4121	The relationship between gene perturbations and the resulting changes in gene
4122	expression is often non-linear.1 Genetic interactions, such as synergistic or
4123	
4124	

i.

25 26 27 28 29	epistatic effects, mean that the combined effect of perturbing multiple genes cannot simply be predicted by summing the individual effects of each perturbation.1 The model must be able to capture these complex, non-additive relationships to make accurate predictions, especially for the 131 gene pairs targeted in the Norman et al. (2019) dataset.
30 31 32 33 34 35	While CRISPRi is generally considered to have fewer off-target effects compared to CRISPRko, the possibility of unintended gene repression at secondary genomic loci still exists. The specificity of CRISPR-Cas9 systems relies heavily on the design of the single guide RNA (sgRNA), and mismatches between the sgRNA and off-target DNA sequences can lead to unintended gene silencing. Although the Norman et al. (2019) study likely employed well-designed sgRNAs, the potential for some level of off-target activity should be considered.
36 37 38 39 40	Single-cell RNA-seq data is also characterized by inherent noise and technical variability.11 Factors such as cell lysis efficiency, reverse transcription efficiency, and sequencing depth can introduce variability in the measured gene expression levels. The predictive model needs to be robust to this noise and focus on capturing the true biological signals induced by the CRISPRi perturbations.
41 42 43 44 45	Finally, clonal expansion, a phenomenon where certain perturbed cells proliferate more rapidly than others, can occur during single-cell CRISPR screens. This can lead to an overrepresentation of certain perturbations in the sequencing data, potentially biasing the results and affecting the training of the predictive model. Computational methods exist to identify and potentially adjust for clonal cell expansion.
46 47 48 49 50 51	Accurately modeling the intricate web of gene regulatory interactions and the non -linear effects of perturbations requires sophisticated computational approaches that can capture these complexities beyond simple additive models. Biological systems are highly interconnected, and the effect of perturbing one gene can propagate through the network, 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.
52 53	New Method Plan
54 55 56 57 58	To address the challenges outlined above and develop an accurate predictive model for gene expression profiles following CRISPRi, a deep learning-based approach is proposed. Specifically, a Graph Neural Network (GNN) architecture, inspired by the success of models like GEARS, appears to be a promising candidate.
59	**Proposed Model Architecture**
50 51 52 53 54	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.
65 66 67 68 69	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.
70	The model architecture will comprise the following key components:
71 72 73 74 75	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.
76 77 78	2. **Perturbation Embedding Layer:** The identity of the perturbed gene(s) will also be encoded into an embedding vector. For single perturbations, a
10	

4180	dedicated embedding will be learned for each targeted gene. For paired
4181 4182	perturbations, the embeddings of the two target genes can be combined (e.g.,
4182	chrough summation of concatenation, to represent the combined perturbation.
4184	3. **Graph Neural Network (GNN):** A gene regulatory network (GRN) will be constructed where genes are represented as nodes and edges represent
4185	regulatory relationships between them. This GRN can be derived from publicly
4186	available databases or inferred from the unperturbed single-cell expression
4187	The GNN will then propagate information across the network, allowing each
4188	gene's representation to be informed by its neighbors and their interactions.
4189	The perturbation embedding will be incorporated into the GNN, potentially by modifying the node features of the perturbed gene(s) or by influencing the
4190	meaning passing process.
4191	4. **Cell State Encoding Laver:** The baseline gene expression profile of the
4192	unperturbed cell will be passed through a separate neural network (e.g., a
4193	multi-layer perceptron) to learn a low-dimensional representation of the cell
4194	context and will be used to condition the prediction of the perturbed state.
4195	5. **Prediction Laver:** The output of the GNN (representing the perturbed gene
4190 4107	embeddings) and the cell state encoding will be combined (e.g., through
4198	concatenation followed by another neural network) to predict the gene
4199	same dimensionality as the input gene expression profile, representing the
4200	predicted expression levels for each gene in the cell.
4201	The rationale behind choosing this architecture is that it allows for the
4202	integration of prior biological knowledge about gene-gene interactions
4203	downstream effects of a perturbation. Furthermore, the use of embeddings
4204	allows the model to learn meaningful representations of genes and
4205	perturbations, potentially enabling better generalization to unseen
4206	*
4207 4208	
4208	used to train the model. The steps involved will include:
4210	1. **Data Loading and Normalization:** The processed gene expression matrices
4211	will be loaded using appropriate libraries like Scanpy or AnnData. The gene
4212	expression counts will be normalized to account for differences in sequencing depth between cells. Log transformation (e.g., using a natural logarithm
4213	after adding a pseudocount) will be applied to stabilize the variance of gene
4214	expression levels.
4215	2. **Perturbation Information Encoding:** The perturbation information,
4210 4217	specifying the targeted gene(s) for each cell, will be extracted from the dataset's metadata. For single gene perturbations, the gene name will be used
4217	For paired gene perturbations, both gene names will be used. These gene
4219	names will then be mapped to their corresponding indices or identifiers in the gene expression matrix. The perturbation information will be opened as
4220	input to the model, potentially using one-hot encoding initially, where a
4221	binary vector indicates which genes are targeted. Alternatively, learned
4222	embeddings for each gene could be used to represent the perturbation.
4223	3. **Control Sample Handling:** Cells labeled as control (unperturbed) will be
4224	control profiles will be crucial for training the model to predict the
4225	changes in expression induced by the perturbations.
4226	4. **Feature Selection:** Given the high dimensionality of the gene expression
4227	data, feature selection techniques may be employed to focus on the most
4228	HVGs) across the cell population and use only these genes as input to the
4229 4230	model.16 This can reduce the dimensionality of the input, potentially
4231	improving model training and performance.
4232	5. **GRN Construction (if applicable):** If a GNN is used, a gene regulatory
4233	network with need to be constructed. This could involve using publicly
4234	

4235	
1236	
4230	
4237	
4238	
4230	
4239	
4240	
4241	
4242	
1212	
4243	
4244	
4245	
1210	
4240	
4247	
4248	
1210	
4249	
4250	
4251	
1252	
+232	
4253	
4254	
1255	
4255	
4256	
4257	
4258	
1250	
4259	
4260	
4261	
1262	
4202	
4263	
4264	
4265	
4205	
4266	
4267	
4268	
1200	
4209	
4270	
4271	
4271 4272	
4271 4272	
4271 4272 4273	
4271 4272 4273 4274	
4271 4272 4273 4274 4275	
4271 4272 4273 4274 4275	
4271 4272 4273 4274 4275 4276	
4271 4272 4273 4274 4275 4275 4276 4277	
4271 4272 4273 4274 4275 4276 4277 4278	
4271 4272 4273 4274 4275 4276 4277 4278 4278	
4271 4272 4273 4274 4275 4276 4277 4278 4279	
4271 4272 4273 4274 4275 4276 4277 4278 4279 4280	
4271 4272 4273 4274 4275 4276 4277 4278 4279 4280 4281	
4271 4272 4273 4274 4275 4276 4277 4278 4279 4280 4281 4282	
4271 4272 4273 4274 4275 4276 4277 4278 4279 4280 4281 4282 4282	
4271 4272 4273 4274 4275 4276 4277 4278 4279 4280 4281 4282 4283	
4271 4272 4273 4274 4275 4276 4277 4278 4279 4280 4281 4282 4283 4284	
4271 4272 4273 4274 4275 4276 4277 4278 4279 4280 4281 4282 4283 4284 4285	
4271 4272 4273 4274 4275 4276 4277 4278 4279 4280 4281 4282 4283 4284 4285 4285	
4271 4272 4273 4274 4275 4276 4277 4278 4279 4280 4281 4282 4283 4284 4285 4284	
4271 4272 4273 4274 4275 4276 4277 4278 4279 4280 4281 4282 4283 4284 4285 4284 4285 4286 4287	

4289

available databases of known gene interactions or inferring a network from the unperturbed single-cell expression data using methods like co-expression analysis or network inference algorithms.

- \*\*Training Strategy\*\*
- The training of the predictive model will involve the following steps:
- 1. \*\*Data Splitting:\*\* The dataset will be split into training, validation, and test sets. To address the evaluation scenarios of unseen perturbations, the split will be performed at the level of perturbations. This means that all cells corresponding to certain perturbations (both single and paired) will be held out in the test set and will not be seen by the model during training. A separate validation set, also containing held-out perturbations, will be used for hyperparameter tuning and model selection.16 Carefully designing the data splitting strategy is crucial to ensure that the model is truly evaluated on unseen perturbations and cell contexts, avoiding information leakage from the training set.
- 2. \*\*Loss Function:\*\* The model will be trained to minimize the difference between the predicted gene expression profiles and the observed gene expression profiles. The Mean Squared Error (MSE) will be used as the primary loss function, as it directly measures the average squared difference between the predicted and observed values.
- 3. \*\*Optimizer:\*\* An appropriate optimization algorithm, such as Adam, will be used to update the model's parameters during training. A learning rate schedule, which gradually reduces the learning rate over time, may be employed to improve convergence and prevent overfitting.
- 4. \*\*Hyperparameter Tuning:\*\* The model architecture and training process will have several hyperparameters (e.g., the dimensionality of the embeddings, the number of layers in the neural networks, the learning rate). These hyperparameters will be tuned using the validation set. Techniques like grid search or random search can be used to explore different combinations of hyperparameters and select the configuration that yields the best performance on the validation set.
- 5. \*\*Overfitting Prevention:\*\* Techniques such as dropout, weight regularization
   (e.g., L1 or L2 regularization), and early stopping (monitoring the
   performance on the validation set and stopping training when it starts to
   degrade) will be used to prevent the model from overfitting to the training
   data and improve its generalization ability.

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

- The proposed GNN-based architecture offers several advantages for handling the challenges of unseen perturbations and cell contexts:
- \*\*Generalization to Unseen Genes (for GNN):\*\* If a GNN is used and the GRN includes genes that are not targeted in the training set, the model may still be able to make informed predictions about the effects of perturbing these unseen genes by leveraging their relationships with other genes in the network.36 The model can learn general principles of how perturbations propagate through the network, allowing it to extrapolate to new nodes (genes ). Predicting the effects of completely novel perturbations (genes not seen during training) is a significant challenge.
- \*\*Cell State Conditioning:\*\* By explicitly encoding the baseline gene expression profile of the unperturbed cell, the model can condition its prediction on the specific context of that cell. This allows the model to capture some of the inherent heterogeneity within the cell population and potentially make more accurate predictions for cells with unseen baseline profiles.
- \*\*Learned Embeddings:\*\* The use of learned embeddings for genes and perturbations can help the model to capture semantic relationships between different genes and perturbations. If the embedding space is learned effectively, the model may be able to generalize to unseen perturbations that are functionally similar to those seen during training, even if the specific genes were not encountered before.

Submission and Formatting Instructions for ICML 2025 GenBio Workshop

÷

4290	
4291	While more advanced techniques like meta-learning or domain adaptation could potentially further improve the model's ability to handle unseen
4292	perturbations and cell contexts, the proposed GNN architecture with cell
4203	state conditioning and learned embeddings provides a strong foundation for
4293	addressing these challenges.
4294	
4295	Generate Prediction Model Code
4296	The prediction model will be implemented using the Duthen programming language
4297	and several key libraries commonly used in machine learning and single-cell
4298	data analysis.
4299	· · Implementation Dataila.
4300	**Implementation Details** - **Programming Language*** Python
4301	- **Key Libraries:**
4302	- **PyTorch** or **TensorFlow/Keras:** For implementing the neural network
4303	architecture, including the GNN and other layers.
4304	- **Scanpy** or **AnnData:** For efficient handling and preprocessing of the single-coll PNA-sog data
4305	- **NumPv:** For numerical computations and array manipulations.
4306	- **SciPy:** For scientific computing, including statistical functions.
4207	- **scikit-learn:** For evaluation metrics (MSE, PCC) and potential utility
4307	functions.
4308	**Code Structure and Key Functions**
4309	The codebase will be organized into several modules or classes to ensure
4310	modularity and maintainability:
4311	1. ** `DataLoader ` Class:**
4312	- **Purpose:** Responsible for loading the Norman et al. (2019) dataset from
4313	the specified files (e.g., MTX, TSV).
4314	<ul> <li>- **Functionality:**</li> <li>- Reads the gene expression matrix and metadata</li> </ul>
4315	- Performs normalization and log transformation.
4316	- Encodes perturbation information.
4317	- Splits the data into training, validation, and test sets based on
4318	perturbations. - Provides data loaders for officient batching during training
4319	2. ** 'Model' Class:**
4320	- **Purpose:** Defines the neural network architecture, including the gene
4320	embedding layer, perturbation embedding layer, GNN (if chosen), cell state
4222	encoding layer, and prediction layer.
4322	- Initializes the model parameters.
4323	- Implements the forward pass of the network, taking baseline expression
4324	and perturbation information as input and outputting the predicted
4325	perturbed expression.
4326	- **Purpose:** Handles the training process of the model.
4327	- **Functionality:**
4328	- Takes the model, data loaders, loss function, and optimizer as input.
4329	<ul> <li>Implements the training loop, including forward pass, loss calculation, backpropagation, and parameter undated</li> </ul>
4330	- Monitors performance on the validation set and implements early stopping
4331	if needed.
4332	- Saves the trained model parameters.
4333	4. ** 'Evaluator' Class:**
4334	- **Functionality:**
4335	- Loads the trained model parameters.
1336	- Iterates through the test data.
1327	- Makes predictions for each test sample.
4220	predictions with the observed data.
4338	- Reports the evaluation results.
4339	5. ** 'Predictor' Class:**
4340	- **Purpose: ** Allows for making predictions on new, unseen data (e.g., a
4341	- **Functionality:**
4342	- Loads the trained model parameters.
4343	
4344	

4345	- Takes a baseline expression profile and perturbation information as input.
4340 4347	- Preprocesses the input data as needed.
4348	- Uses the trained model to generate the predicted gene expression profile after the perturbation.
4349 4350	**Instructions for Usage**
4351	1. **Install Libraries:** Install the necessary Python libraries using pip:
4352	''Bash pip install torch torchvision torchaudio scanpy numpy scipy scikit-learn
4353	(or )topsorflow) and )koras) if using TopsorFlow)
4355	(or tensorriow and keras in using remotive).
4356	2. **Download Data:** Download the processed data files for GSE133344 from GEO or Figshare and place them in a designated directory.
4358 4359	3. **Prepare Input Data:** Ensure that the input data for prediction (baseline gene expression profile and perturbation identity) is in the expected format (e.g., a NumPy array or a Pandas DataFrame).
4360 4361	4. **Run the Code:** Execute the Python script containing the model implementation, training, and evaluation steps. The script should:
4362 4363 4364	<ul> <li>Initialize the 'DataLoader' with the path to the data files.</li> <li>Split the data into training, validation, and test sets.</li> <li>Initialize the 'Model' with the appropriate architecture and hyperparameters.</li> </ul>
4365	- Initialize the `Trainer` with the model, data loaders, loss function, and
4366 4367	optimizer. - Train the model using the 'Trainer'.
4368	<ul> <li>Initialize the 'Evaluator' with the trained model and test data loader.</li> <li>Evaluate the model and print the results.</li> </ul>
4369	- (Optionally) Use the 'Predictor' to make predictions on new data.
4370	**Evaluation and Results**
4371 4372 4373	A comprehensive evaluation of the developed predictive model is crucial to assess its performance and generalizability.
4373	**Experimental Setup**
4375 4376 4377 4378 4379 4380 4381	The Norman et al. (2019) dataset will be split into training, validation, and test sets, ensuring that perturbations are held out for both validation and testing to evaluate the model's ability to predict the effects of unseen genetic manipulations. The split will likely involve a significant portion of the data for training, with smaller portions allocated for validation (for hyperparameter tuning and model selection) and testing (for final performance evaluation).16 The specific hyperparameters of the chosen model architecture (e.g., number of layers, dimensionality of embeddings, learning rate) will be determined through experimentation and tuning on the validation set.
4382 4383 4384 4385	To make a prediction for a given perturbation, the model will take as input the baseline gene expression profile of an 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).
4386	**Performance Metrics**
4387 4388 4389	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 gape expression profiles with
4390 4391	the actual observed profiles in the test set for the held-out perturbations.
4392	and unseen paired-gene perturbations to assess the model's ability to handle
4393	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.
4394	as these are often the most biologically relevant. Visualizations, such as
4395	scatter plots of predicted vs. observed gene expression for representative perturbations, can provide further insights into the model's predictive
4396 4397	capabilities.
4398	
4399	

#### I. LLM-as-a-judge details

#### I.1. methods

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

Outputs(Examples can be found in Appendix G) from scAgents(with 5 different LLM APIconfigurations: Claude 3.7, o1, DeepSeek R1, Qwen-plus, Llama 3,1) were anonymized to prevent bias. For each evaluation round, a set of 8 outputs was randomly selected and their order randomized to ensure fairness. This process was repeated 5 times, resulting in 5 distinct evaluation rounds with different output sequences. In each round, the LLM evaluated the 8 outputs individually, assigning a score from 1 to 10 based on predefined criteria. The LLM was unaware of the source of each output, ensuring unbiased assessments. 

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

#### I.2. prompts

L	LM As Judge-Gene
Т	ask Analysis
	You are an expert evaluator specializing in data-driven analysis of CRISPR-based single-cell perturbation experiments. Your background includes:
	<ul> <li>In-depth knowledge of single-cell omics data modalities (e.g., RNA-seq, ATAC-seq, CITE-seq)</li> <li>Experience in characterizing perturbation types and experimental settings</li> <li>Familiarity with agent-based literature retrieval and scientific reasoning</li> <li>Ability to assess baseline model performance in biological prediction tasks</li> <li>Understanding of automated systems for scientific task decomposition</li> </ul>
	Please evaluate the following Task Analysis report using rigorous and objective scientific standards. You may receive multiple reports in randomized order across five rounds. **Evaluate each report independently**, without assuming knowledge of other submissions.
	EVALUATION CRITERIA Each criterion should be scored on a scale of 1~10, with clear justification based on the report content. Use full-score ranges (1~10) where appropriate.
	<ol> <li>Analyse Dataset (1~10):         <ul> <li>Clarity and correctness in summarizing dataset properties (modality, perturbation type, species, cell type distribution, etc.)</li> <li>Relevance of identified features for downstream modeling</li> <li>Ability to standardize and interpret metadata across modalities</li> <li>Quality of data summaries and diagnostic insights (e.g., sparsity, heterogeneity)</li> </ul> </li> </ol>
	<ul> <li>2. Analyse Task Type (1<sup>1</sup>0):</li> <li>Accuracy in identifying the biological question and mapping it to a computational prediction task</li> <li>Insightfulness in selecting the right task framing (e.g., classification vs regression, single-cell vs population-level)</li> <li>Alignment of task framing with perturbation mechanism and data granularity</li> <li>Ability to distinguish this task from superficially similar ones</li> </ul>
	<ul> <li>3. Analyse Baseline Defects (1~10):</li> <li>Thoroughness in identifying limitations of current baseline models</li> <li>Correctness in linking model weaknesses to data/task-specific challenges (e.g. model mismatch with modality, lack of interpretability)</li> <li>Thoughtfulness in proposing key evaluation gaps or unaddressed risks</li> <li>Clarity in explaining why the baseline is insufficient and what improvement directions are needed</li> </ul>
	FORMAT FOR YOUR EVALUATION: 1. NUMERICAL SCORES

1155	
4433	Analyse Dataset: [Score]/10
4456	Analyse Task Type: [Score]/10
4457	Analyse Baseline Defects: [Score]/10
1150	
4438	2. DETAILED JUSTIFICATION
4459	Provide specific and concise reasoning for each score, referencing relevant parts
4460	of the analysis. Address both strengths and limitations within each
1161	criterion.
4461	
4462	3. KEY STRENGTHS
1162	[List major strengths of the Task Analysis report]
4403	[Reference specific elements that demonstrate scientific rigor or originality]
4464	
4465	4. AREAS FOR IMPROVEMENT
1166	[Identify specific aspects that could be clarified or strengthened]
4400	[Offer constructive, actionable suggestions for refinement]
4467	
4468	Drovide a concise event Consider:
1160	Provide a concise overall assessment. Consider:
4409	- Does the lask Analysis provide a strong foundation for follow-up modeling?
4470	- Are the dataset and task leatures well-characterized and actionable?
4471	- Are the limitations of baseline models accurately diagnosed and explained?
4 4 7 0	REMINDERS.
4472	- Maintain scientific neutrality and avoid assumptions not grounded in the
4473	provided text
4474	- Consider both biological and computational aspects equally
	- Consider both biological and computational aspects equally.
4475	- Provide constructive reedback almed at improving scientific understanding.
4476	- Use current SOTA practices in perturbation modeling and single-cell analysis as
1177	your reference frame.
44//	- Assume the audience is a mix of computational biologists, experimentalists, and
4478	system developers.
4479	TASK ANALYSIS REPORT TO EVALUATE.
1100	[Paste your report here] / [Will be provided in the next message]
4480	[rable your report here] / [will be provided in the none mesoage]
4481	
4482	LIMAs Judge Cone
4482	LLM As Judge-Gene
4482 4483	LLM As Judge-Gene
4482 4483 4484	LLM As Judge-Gene       Method Design
4482 4483 4484 4485	LLM As Judge-Gene       Method Design
4482 4483 4484 4485 4485	LLM As Judge-Gene       Method Design       You are an expert evaluator specializing in CRISPR-based single-cell perturbation
4482 4483 4484 4485 4485	LLM As Judge-Gene         Method Design         You are an expert evaluator specializing in CRISPR-based single-cell perturbation         prediction models and experimental designs. Your background includes:
4482 4483 4484 4485 4486 4487	LLM As Judge-Gene         Method Design         You are an expert evaluator specializing in CRISPR-based single-cell perturbation prediction models and experimental designs. Your background includes:
4482 4483 4484 4485 4485 4486 4487 4488	LLM As Judge-Gene         Method Design         You are an expert evaluator specializing in CRISPR-based single-cell perturbation prediction models and experimental designs. Your background includes:         - Deep expertise in computational biology and single-cell omics
4482 4483 4484 4485 4486 4487 4488 4488	LLM As Judge-Gene         Method Design         You are an expert evaluator specializing in CRISPR-based single-cell perturbation prediction models and experimental designs. Your background includes:         - Deep expertise in computational biology and single-cell omics         - Practical experience with CRISPR-based perturbation experiments
4482 4483 4484 4485 4486 4487 4488 4489	LLM As Judge-Gene         Method Design         You are an expert evaluator specializing in CRISPR-based single-cell perturbation prediction models and experimental designs. Your background includes:         - Deep expertise in computational biology and single-cell omics         - Practical experience with CRISPR-based perturbation experiments         - Familiarity with multimodal single-cell datasets (e.g., gene expression, ATAC-
4482 4483 4484 4485 4486 4487 4488 4489 4489	LLM As Judge-Gene         Method Design         You are an expert evaluator specializing in CRISPR-based single-cell perturbation prediction models and experimental designs. Your background includes:         - Deep expertise in computational biology and single-cell omics         - Practical experience with CRISPR-based perturbation experiments         - Familiarity with multimodal single-cell datasets (e.g., gene expression, ATAC-seq, protein expression)
4482 4483 4484 4485 4486 4487 4488 4489 4490 4491	LLM As Judge-Gene         Method Design         You are an expert evaluator specializing in CRISPR-based single-cell perturbation prediction models and experimental designs. Your background includes:         - Deep expertise in computational biology and single-cell omics         - Practical experience with CRISPR-based perturbation experiments         - Familiarity with multimodal single-cell datasets (e.g., gene expression, ATAC-seq, protein expression)         - Advanced understanding of machine learning models for biological prediction
4482 4483 4484 4485 4485 4486 4487 4488 4489 4490 4491 4492	LLM As Judge-Gene         Method Design         You are an expert evaluator specializing in CRISPR-based single-cell perturbation prediction models and experimental designs. Your background includes:         - Deep expertise in computational biology and single-cell omics         - Practical experience with CRISPR-based perturbation experiments         - Familiarity with multimodal single-cell datasets (e.g., gene expression, ATAC-seq, protein expression)         - Advanced understanding of machine learning models for biological prediction tasks
4482 4483 4484 4485 4486 4487 4488 4489 4490 4491 4492	LLM As Judge-Gene         Method Design         You are an expert evaluator specializing in CRISPR-based single-cell perturbation prediction models and experimental designs. Your background includes:         - Deep expertise in computational biology and single-cell omics         - Practical experience with CRISPR-based perturbation experiments         - Familiarity with multimodal single-cell datasets (e.g., gene expression, ATAC-seq, protein expression)         - Advanced understanding of machine learning models for biological prediction tasks         - Knowledge of statistical validation methods and experimental reproducibility
4482 4483 4484 4485 4486 4487 4488 4489 4490 4491 4492 4493	LLM As Judge-Gene         Method Design         You are an expert evaluator specializing in CRISPR-based single-cell perturbation prediction models and experimental designs. Your background includes:         - Deep expertise in computational biology and single-cell omics         - Practical experience with CRISPR-based perturbation experiments         - Familiarity with multimodal single-cell datasets (e.g., gene expression, ATAC-seq, protein expression)         - Advanced understanding of machine learning models for biological prediction tasks         - Knowledge of statistical validation methods and experimental reproducibility standards
4482 4483 4484 4485 4486 4487 4488 4489 4490 4490 4491 4492 4493 4494	<pre>LLM As Judge-Gene Method Design You are an expert evaluator specializing in CRISPR-based single-cell perturbation     prediction models and experimental designs. Your background includes:     Deep expertise in computational biology and single-cell omics     Practical experience with CRISPR-based perturbation experiments     Familiarity with multimodal single-cell datasets (e.g., gene expression, ATAC-     seq, protein expression)     Advanced understanding of machine learning models for biological prediction     tasks     Knowledge of statistical validation methods and experimental reproducibility     standards     Awareness of recent state-of-the-art (SOTA) approaches in perturbation modeling</pre>
4482 4483 4484 4485 4486 4487 4488 4489 4490 4490 4491 4492 4493 4494	<pre>LLM As Judge-Gene Method Design You are an expert evaluator specializing in CRISPR-based single-cell perturbation     prediction models and experimental designs. Your background includes:     Deep expertise in computational biology and single-cell omics     Practical experience with CRISPR-based perturbation experiments     Familiarity with multimodal single-cell datasets (e.g., gene expression, ATAC-     seq, protein expression)     Advanced understanding of machine learning models for biological prediction     tasks     Knowledge of statistical validation methods and experimental reproducibility     standards     Awareness of recent state-of-the-art (SOTA) approaches in perturbation modeling</pre>
4482 4483 4484 4485 4486 4487 4488 4489 4490 4491 4492 4493 4494	<pre>LLM As Judge-Gene Method Design You are an expert evaluator specializing in CRISPR-based single-cell perturbation prediction models and experimental designs. Your background includes: Deep expertise in computational biology and single-cell omics Practical experience with CRISPR-based perturbation experiments Familiarity with multimodal single-cell datasets (e.g., gene expression, ATAC- seq, protein expression) Advanced understanding of machine learning models for biological prediction tasks Knowledge of statistical validation methods and experimental reproducibility standards Awareness of recent state-of-the-art (SOTA) approaches in perturbation modeling Please evaluate the following research plan using rigorous and objective</pre>
4482 4483 4484 4485 4486 4487 4488 4489 4490 4491 4492 4493 4494 4495 4496	LLM As Judge-GeneMethod DesignYou are an expert evaluator specializing in CRISPR-based single-cell perturbation prediction models and experimental designs. Your background includes:- Deep expertise in computational biology and single-cell omics - Practical experience with CRISPR-based perturbation experiments - Familiarity with multimodal single-cell datasets (e.g., gene expression, ATAC- seq, protein expression)- Advanced understanding of machine learning models for biological prediction tasks - Knowledge of statistical validation methods and experimental reproducibility standards - Awareness of recent state-of-the-art (SOTA) approaches in perturbation modeling Please evaluate the following research plan using rigorous and objective scientific standards. You may receive multiple plans in randomized order
4482 4483 4484 4485 4486 4487 4488 4489 4490 4490 4491 4492 4493 4494 4495 4496 4497	<pre>LLM As Judge-Gene Method Design You are an expert evaluator specializing in CRISPR-based single-cell perturbation     prediction models and experimental designs. Your background includes:     Deep expertise in computational biology and single-cell omics     Practical experience with CRISPR-based perturbation experiments     Familiarity with multimodal single-cell datasets (e.g., gene expression, ATAC-     seq, protein expression)     Advanced understanding of machine learning models for biological prediction     tasks     Knowledge of statistical validation methods and experimental reproducibility     standards     Awareness of recent state-of-the-art (SOTA) approaches in perturbation modeling     Please evaluate the following research plan using rigorous and objective     scientific standards. You may receive multiple plans in randomized order     across five rounds. **Evaluate each plan independently**, without assuming </pre>
4482 4483 4484 4485 4486 4487 4488 4489 4490 4490 4491 4492 4493 4494 4495 4496 4497 4498	<pre>LLM As Judge-Gene Method Design You are an expert evaluator specializing in CRISPR-based single-cell perturbation prediction models and experimental designs. Your background includes: - Deep expertise in computational biology and single-cell omics - Practical experience with CRISPR-based perturbation experiments - Familiarity with multimodal single-cell datasets (e.g., gene expression, ATAC- seq, protein expression) - Advanced understanding of machine learning models for biological prediction tasks - Knowledge of statistical validation methods and experimental reproducibility standards - Awareness of recent state-of-the-art (SOTA) approaches in perturbation modeling Please evaluate the following research plan using rigorous and objective scientific standards. You may receive multiple plans in randomized order across five rounds. **Evaluate each plan independently**, without assuming knowledge of other submissions.</pre>
4482 4483 4484 4485 4486 4487 4488 4489 4490 4491 4492 4493 4494 4495 4496 4497	<pre>LLM As Judge-Gene Method Design You are an expert evaluator specializing in CRISPR-based single-cell perturbation prediction models and experimental designs. Your background includes: - Deep expertise in computational biology and single-cell omics - Practical experience with CRISPR-based perturbation experiments - Familiarity with multimodal single-cell datasets (e.g., gene expression, ATAC- seq, protein expression) - Advanced understanding of machine learning models for biological prediction tasks - Knowledge of statistical validation methods and experimental reproducibility standards - Awareness of recent state-of-the-art (SOTA) approaches in perturbation modeling Please evaluate the following research plan using rigorous and objective scientific standards. You may receive multiple plans in randomized order across five rounds. **Evaluate each plan independently**, without assuming knowledge of other submissions.</pre>
4482 4483 4484 4485 4486 4487 4488 4489 4490 4491 4492 4493 4494 4495 4496 4497 4498 4499	<pre>LLM As Judge-Gene Method Design You are an expert evaluator specializing in CRISPR-based single-cell perturbation prediction models and experimental designs. Your background includes: - Deep expertise in computational biology and single-cell omics - Practical experience with CRISPR-based perturbation experiments - Familiarity with multimodal single-cell datasets (e.g., gene expression, ATAC- seq, protein expression) - Advanced understanding of machine learning models for biological prediction tasks - Knowledge of statistical validation methods and experimental reproducibility standards - Awareness of recent state-of-the-art (SOTA) approaches in perturbation modeling Please evaluate the following research plan using rigorous and objective scientific standards. You may receive multiple plans in randomized order across five rounds. **Evaluate each plan independently**, without assuming knowledge of other submissions. EVALUATION CRITERIA EVALUATION CRITERIA</pre>
4482 4483 4484 4485 4486 4487 4488 4489 4490 4491 4492 4493 4494 4495 4496 4497 4498 4499 4500	<pre>LLLM As Judge-Gene Method Design You are an expert evaluator specializing in CRISPR-based single-cell perturbation prediction models and experimental designs. Your background includes: - Deep expertise in computational biology and single-cell omics - Practical experience with CRISPR-based perturbation experiments - Familiarity with multimodal single-cell datasets (e.g., gene expression, ATAC- seq, protein expression) - Advanced understanding of machine learning models for biological prediction tasks - Knowledge of statistical validation methods and experimental reproducibility standards - Awareness of recent state-of-the-art (SOTA) approaches in perturbation modeling Please evaluate the following research plan using rigorous and objective scientific standards. You may receive multiple plans in randomized order across five rounds. **Evaluate each plan independently**, without assuming knowledge of other submissions. EVALUATION CRITERIA Each criterion should be scored on a scale of 1~10, with clear justification</pre>
4482 4483 4484 4485 4486 4487 4488 4489 4490 4490 4491 4492 4493 4494 4495 4496 4497 4498 4499 4500 4501	<pre>LLM As Judge-Gene Method Design You are an expert evaluator specializing in CRISPR-based single-cell perturbation prediction models and experimental designs. Your background includes: - Deep expertise in computational biology and single-cell omics - Practical experience with CRISPR-based perturbation experiments - Familiarity with multimodal single-cell datasets (e.g., gene expression, ATAC- seq, protein expression) - Advanced understanding of machine learning models for biological prediction tasks - Knowledge of statistical validation methods and experimental reproducibility standards - Awareness of recent state-of-the-art (SOTA) approaches in perturbation modeling Please evaluate the following research plan using rigorous and objective scientific standards. You may receive multiple plans in randomized order across five rounds. **Evaluate each plan independently**, without assuming knowledge of other submissions. EVALUATION CRITERIA Each criterion should be scored on a scale of 1~10, with clear justification based on the content of the plan. Use full-score ranges (1~10) where</pre>
4482 4483 4484 4485 4486 4487 4488 4489 4490 4490 4491 4492 4493 4494 4495 4496 4497 4498 4499 4500 4501	LLM As Judge-GeneMethod DesignYou are an expert evaluator specializing in CRISPR-based single-cell perturbation prediction models and experimental designs. Your background includes:- Deep expertise in computational biology and single-cell omics - Practical experience with CRISPR-based perturbation experiments- Familiarity with multimodal single-cell datasets (e.g., gene expression, ATAC- seq, protein expression)- Advanced understanding of machine learning models for biological prediction tasks- Knowledge of statistical validation methods and experimental reproducibility standards- Awareness of recent state-of-the-art (SOTA) approaches in perturbation modelingPlease evaluate the following research plan using rigorous and objective scientific standards. You may receive multiple plans in randomized order across five rounds. **Evaluate each plan independently**, without assuming knowledge of other submissions.EVALUATION CRITERIA Each criterion should be scored on a scale of 1~10, with clear justification based on the content of the plan. Use full-score ranges (1~10) where appropriate.
4482 4483 4484 4485 4486 4487 4488 4489 4490 4490 4491 4492 4493 4494 4495 4496 4497 4498 4499 4500 4501 4502	<pre>LLM As Judge-Gene Method Design You are an expert evaluator specializing in CRISPR-based single-cell perturbation prediction models and experimental designs. Your background includes: - Deep expertise in computational biology and single-cell omics - Practical experience with CRISPR-based perturbation experiments - Familiarity with multimodal single-cell datasets (e.g., gene expression, ATAC- seq, protein expression) - Advanced understanding of machine learning models for biological prediction tasks - Knowledge of statistical validation methods and experimental reproducibility standards - Awareness of recent state-of-the-art (SOTA) approaches in perturbation modeling Please evaluate the following research plan using rigorous and objective scientific standards. You may receive multiple plans in randomized order across five rounds. **Evaluate each plan independently**, without assuming knowledge of other submissions. EVALUATION CRITERIA Each criterion should be scored on a scale of 1~10, with clear justification based on the content of the plan. Use full-score ranges (1~10) where appropriate. 1 Scientific Validity (1~10):</pre>
4482 4483 4484 4485 4486 4487 4488 4489 4490 4490 4491 4492 4493 4494 4495 4496 4497 4498 4499 4500 4501 4502 4503	<pre>LLM As Judge-Gene Method Design You are an expert evaluator specializing in CRISPR-based single-cell perturbation prediction models and experimental designs. Your background includes: - Deep expertise in computational biology and single-cell omics - Practical experience with CRISPR-based perturbation experiments - Familiarity with multimodal single-cell datasets (e.g., gene expression, ATAC- seq, protein expression) - Advanced understanding of machine learning models for biological prediction tasks - Knowledge of statistical validation methods and experimental reproducibility standards - Awareness of recent state-of-the-art (SOTA) approaches in perturbation modeling Please evaluate the following research plan using rigorous and objective scientific standards. You may receive multiple plans in randomized order across five rounds. **Evaluate each plan independently**, without assuming knowledge of other submissions. EVALUATION CRITERIA Each criterion should be scored on a scale of 1~10, with clear justification based on the content of the plan. Use full-score ranges (1~10) where appropriate. 1. Scientific Validity (1~10): - Biological relevance and mechanistic insight</pre>
4482 4483 4484 4485 4486 4487 4488 4489 4490 4490 4491 4492 4493 4494 4495 4496 4497 4498 4499 4500 4501 4502 4503 4504	<pre>LLM As Judge-Gene Method Design You are an expert evaluator specializing in CRISPR-based single-cell perturbation     prediction models and experimental designs. Your background includes:     Deep expertise in computational biology and single-cell omics     Practical experience with CRISPR-based perturbation experiments     Familiarity with multimodal single-cell datasets (e.g., gene expression, ATAC-     seq, protein expression)     Advanced understanding of machine learning models for biological prediction     tasks     Knowledge of statistical validation methods and experimental reproducibility     standards     Awareness of recent state-of-the-art (SOTA) approaches in perturbation modeling     Please evaluate the following research plan using rigorous and objective     scientific standards. You may receive multiple plans in randomized order     across five rounds. **Evaluate each plan independently**, without assuming     knowledge of other submissions.     EVALUATION CRITERIA     Each criterion should be scored on a scale of 1~10, with clear justification     based on the content of the plan. Use full-score ranges (1~10) where     appropriate.     Scrength of theoretical foundation</pre>
4482 4483 4484 4485 4486 4487 4488 4489 4490 4491 4492 4493 4494 4495 4494 4495 4496 4497 4498 4499 4500 4501 4502	<pre>LLM As Judge-Gene Method Design You are an expert evaluator specializing in CRISPR-based single-cell perturbation prediction models and experimental designs. Your background includes:     - Deep expertise in computational biology and single-cell omics     Practical experience with CRISPR-based perturbation experiments     Familiarity with multimodal single-cell datasets (e.g., gene expression, ATAC-     seq, protein expression)     Advanced understanding of machine learning models for biological prediction     tasks     Knowledge of statistical validation methods and experimental reproducibility     standards     Awareness of recent state-of-the-art (SOTA) approaches in perturbation modeling Please evaluate the following research plan using rigorous and objective     scientific standards. You may receive multiple plans in randomized order     across five rounds. **Evaluate each plan independently**, without assuming     knowledge of other submissions.  EVALUATION CRITERIA Each criterion should be scored on a scale of 1~10, with clear justification     based on the content of the plan. Use full-score ranges (1~10) where     appropriate. 1. Scientific Validity (1~10):     Biological relevance and mechanistic insight     Strength of theoretical foundation     Alignment with every experiments     Audition approaches in single-cell</pre>
4482 4483 4484 4485 4486 4487 4488 4489 4490 4491 4492 4493 4494 4495 4494 4495 4496 4497 4498 4499 4500 4501 4502 4503 4504 4505	<pre>LLM As Judge-Gene Method Design You are an expert evaluator specializing in CRISPR-based single-cell perturbation prediction models and experimental designs. Your background includes:     Deep expertise in computational biology and single-cell omics     Practical experience with CRISPR-based perturbation experiments     Familiarity with multimodal single-cell datasets (e.g., gene expression, ATAC-     seq, protein expression)     Advanced understanding of machine learning models for biological prediction     tasks     Knowledge of statistical validation methods and experimental reproducibility     standards     Awareness of recent state-of-the-art (SOTA) approaches in perturbation modeling Please evaluate the following research plan using rigorous and objective     scientific standards. You may receive multiple plans in randomized order     across five rounds. **Evaluate each plan independently**, without assuming     knowledge of other submissions. EVALUATION CRITERIA Each criterion should be scored on a scale of 1°10, with clear justification     based on the content of the plan. Use full-score ranges (1°10) where     appropriate. 1. Scientific Validity (1°10):     Biological relevance and mechanistic insight     Strength of theoretical foundation     Alignment with current scientific understanding in single-cell biology </pre>
4482 4483 4484 4485 4486 4487 4488 4489 4490 4491 4492 4493 4494 4495 4494 4495 4496 4497 4498 4499 4500 4501 4502 4503 4504 4505 4506	<pre>LLM As Judge-Gene Method Design You are an expert evaluator specializing in CRISPR-based single-cell perturbation     prediction models and experimental designs. Your background includes:     Deep expertise in computational biology and single-cell omics     Practical experience with CRISPR-based perturbation experiments     Familiarity with multimodal single-cell datasets (e.g., gene expression, ATAC-     seq, protein expression)     Advanced understanding of machine learning models for biological prediction     tasks     Knowledge of statistical validation methods and experimental reproducibility     standards     Awareness of recent state-of-the-art (SOTA) approaches in perturbation modeling Please evaluate the following research plan using rigorous and objective     scientific standards. You may receive multiple plans in randomized order     across five rounds. **Evaluate each plan independently**, without assuming     knowledge of other submissions. EVALUATION CRITERIA Each criterion should be scored on a scale of 1~10, with clear justification     based on the content of the plan. Use full-score ranges (1~10) where     appropriate. 1. Scientific Validity (1~10):     Biological relevance and mechanistic insight     Strength of theoretical foundation     Alignment with current scientific understanding in single-cell biology Integration with existing knowledge on perturbation responses </pre>
4482 4483 4484 4485 4486 4487 4488 4489 4490 4490 4491 4492 4493 4494 4495 4494 4495 4496 4497 4498 4499 4500 4501 4502 4503 4504 4505 4506 4507	<pre>LLM As Judge-Gene Method Design You are an expert evaluator specializing in CRISPR-based single-cell perturbation prediction models and experimental designs. Your background includes:     Deep expertise in computational biology and single-cell omics     Practical experience with CRISPR-based perturbation experiments     Familiarity with multimodal single-cell datasets (e.g., gene expression, ATAC-     seq, protein expression)     Advanced understanding of machine learning models for biological prediction     tasks     Knowledge of statistical validation methods and experimental reproducibility     standards     Awareness of recent state-of-the-art (SOTA) approaches in perturbation modeling Please evaluate the following research plan using rigorous and objective     scientific standards. You may receive multiple plans in randomized order     across five rounds. **Evaluate each plan independently**, without assuming     knowledge of other submissions. EVALUATION CRITERIA Each criterion should be scored on a scale of 1~10, with clear justification     based on the content of the plan. Use full-score ranges (1~10) where     appropriate. 1. Scientific Validity (1~10):     Biological relevance and mechanistic insight     Strength of theoretical foundation     Alignment with current scientific understanding in single-cell biology</pre>
4482 4483 4484 4485 4486 4487 4488 4489 4490 4490 4491 4492 4493 4494 4495 4494 4495 4496 4497 4498 4499 4500 4501 4502 4503 4504 4505	<pre>LLM As Judge-Gene Method Design You are an expert evaluator specializing in CRISPR-based single-cell perturbation prediction models and experimental designs. Your background includes:         - Deep expertise in computational biology and single-cell omics         - Practical experience with CRISPR-based perturbation experiments         - Familiarity with multimodal single-cell datasets (e.g., gene expression, ATAC-             seq, protein expression)         - Advanced understanding of machine learning models for biological prediction         tasks         - Knowledge of statistical validation methods and experimental reproducibility         standards         - Awareness of recent state-of-the-art (SOTA) approaches in perturbation modeling Please evaluate the following research plan using rigorous and objective         scientific standards. You may receive multiple plans in randomized order         across five rounds. **Evaluate each plan independently**, without assuming         knowledge of other submissions. EVALUATION CRITERIA Each criterion should be scored on a scale of 1~10, with clear justification         based on the content of the plan. Use full-score ranges (l^-10) where         appropriate. 1. Scientific Validity (l^-10):         - Biological relevance and mechanistic insight         - Strength of theoretical foundation         - Alignment with current scientific understanding in single-cell biology Integration with existing knowledge on perturbation responses 2. Technical Feasibility (l^-10): </pre>
4482 4483 4484 4485 4486 4487 4488 4489 4490 4490 4491 4492 4493 4494 4495 4494 4495 4496 4497 4498 4499 4500 4501 4502 4503 4504 4505 4506 4507 4508	<pre>LLM As Judge-Gene Method Design You are an expert evaluator specializing in CRISPR-based single-cell perturbation prediction models and experimental designs. Your background includes:     Deep expertise in computational biology and single-cell omics     Practical experience with CRISPR-based perturbation experiments     Familiarity with multimodal single-cell datasets (e.g., gene expression, ATAC-     seq, protein expression)     Advanced understanding of machine learning models for biological prediction     tasks     Knowledge of statistical validation methods and experimental reproducibility     standards     Awareness of recent state-of-the-art (SOTA) approaches in perturbation modeling Please evaluate the following research plan using rigorous and objective     scientific standards. You may receive multiple plans in randomized order     across five rounds. **Evaluate each plan independently**, without assuming     knowledge of other submissions. EVALUATION CRITERIA Each criterion should be scored on a scale of 1~10, with clear justification     based on the content of the plan. Use full-score ranges (1~10) where     appropriate. 1. Scientific Validity (1~10):     Biological relevance and mechanistic insight     Strength of theoretical foundation     Alignment with current scientific understanding in single-cell biology Integration with existing knowledge on perturbation responses 2. Technical Feasibility (1~10): </pre>
4482 4483 4484 4485 4486 4487 4488 4489 4490 4490 4491 4492 4493 4494 4493 4494 4495 4496 4497 4498 4499 4500 4501 4502 4503 4504 4505 4506 4507 4508 4509	<pre>LLM As Judge-Gene Method Design You are an expert evaluator specializing in CRISPR-based single-cell perturbation     prediction models and experimental designs. Your background includes:     Deep expertise in computational biology and single-cell omics     Practical experience with CRISPR-based perturbation experiments     Familiarity with multimodal single-cell datasets (e.g., gene expression, ATAC-     seq, protein expression)     Advanced understanding of machine learning models for biological prediction     tasks     Knowledge of statistical validation methods and experimental reproducibility     standards     Awareness of recent state-of-the-art (SOTA) approaches in perturbation modeling Please evaluate the following research plan using rigorous and objective     scientific standards. You may receive multiple plans in randomized order     across five rounds. **Evaluate each plan independently**, without assuming     knowledge of other submissions. EVALUATION CRITERIA Each criterion should be scored on a scale of 1^10, with clear justification     based on the content of the plan. Use full-score ranges (1^10) where     appropriate. 1. Scientific Validity (1^10):     Biological relevance and mechanistic insight     Alignment with current scientific understanding in single-cell biology Integration with existing knowledge on perturbation responses 2. Technical Feasibility (1^10): </pre>

4510 - Practicality of implementation 4511 - Computational resource requirements 4512 - Scalability to larger datasets or new tasks - Feasibility and clarity of data preprocessing or modeling pipeline 4513 4514 3. Innovation Level (1~10): - Novelty compared to current state-of-the-art approaches 4515 - Creative problem-solving or hypothesis generation 4516 - Potential for new biological or computational insights 4517 - Unique contributions in methodology or design 4518 4. Experimental Design (1~10): 4519 - Quality of proposed validation and evaluation methodology - Inclusion of appropriate controls and baselines 4520 - Statistical soundness (e.g., replicates, robustness) 4521 - Attention to data quality and reproducibility 4522 5. Impact Potential (1~10): 4523 - Relevance and contribution to advancing single-cell biology 4524 - Translational potential (e.g., drug discovery, therapeutic design) 4525 - Scalability to broader biological questions or contexts - Potential to inspire follow-up research or community adoption 4526 4527 FORMAT FOR YOUR EVALUATION: 4528 1. NUMERICAL SCORES Scientific Validity: [Score]/10 4529 Technical Feasibility: [Score]/10 4530 Innovation Level: [Score]/10 Experimental Design: [Score]/10 4531 Impact Potential: [Score]/10 4532 4533 2. DETAILED JUSTIFICATION Provide specific and concise reasoning for each score, referencing relevant parts 4534 of the research plan. Address both strengths and limitations within each 4535 criterion. 4536 3. KEY STRENGTHS 4537 [List major strengths of the research plan] 4538 [Reference specific elements that demonstrate scientific merit or originality] 4539 4. AREAS FOR IMPROVEMENT 4540 [Identify specific aspects that could be clarified or strengthened] 4541 4542 [Offer constructive, actionable suggestions for refinement] 4543 5. OVERALL RECOMMENDATION Provide a concise overall assessment. Consider: 4544 - Is this plan a strong candidate for further development or implementation? 4545 - What would be the next logical steps to improve scientific or translational 4546 value? 4547 - How does this plan compare (in quality and originality) to what is typically seen in high-impact research? 4548 4549 Additional Considerations for Gene Knockout: 4550 - Gene regulatory network modeling - Knockout effect propagation 4551 - Compensatory mechanism prediction 4552 - Pathway analysis integration 4553 REMINDERS: 4554 - Maintain scientific neutrality and avoid assumptions not grounded in the 4555 provided text. - Consider both biological and computational aspects equally. 4556 - Provide constructive feedback aimed at improving the research direction. 4557 - Use current state-of-the-art methods and literature in single-cell genomics and 4558 perturbation modeling as your reference frame. - Assume the audience is a mix of computational biologists, experimentalists, and 4559 translational researchers. 4560 RESEARCH PLAN TO EVALUATE: 4561 [Paste your plan here] / [Will be provided in the next message] 4562 4563 4564

### **I.3. example output**

4566	
4567	LLM As Judge Output-Gene
4568	Task Analysis-Claude 3.7
4569	1 NUMERICAL SCORES
4570	Analyse Dataset: 8/10
4571	Analyse Task Type: 8/10
4572	Analyse Baseline Defects: 8/10
4573	2. DETAILED JUSTIFICATION
4574	Analyse Dataset (8/10)
4575	The report provides a comprehensive and detailed summary of the dataset
4576	RNA-seq with CRISPR perturbations), and specific features like the scale
4577	(~170,000 cells), composition (1,092 unique perturbation conditions), and
4578	high-dimensional readouts (~20,000 genes per cell). The characterization of
4579	perturbation types (single-gene knockdowns and two-gene combinations) is
4580	interactions. The report effectively identifies key challenges such as single-
4581	cell variability and the need for generalization to unseen perturbations. The
4582	only minor limitation is the lack of specific diagnostic metrics or
4583	s complexity, such as dropout rates or batch effect magnitudes. The data
4584	summaries are thorough and provide critical diagnostic insights for model
4585	development.
4586	Analyse Task Type (8/10)
4587	The biological question is accurately identified as learning a high-dimensional
4588	function to predict post-perturbation gene expression profiles. The mapping
4589	to a computational prediction task (nign-dimensional regression with combinatorial inputs) is insightful and aligns well with the dataset's
4500	granularity and the nature of CRISPR perturbations. The justification for
4501	selecting this task type is robust, emphasizing the biological relevance of
4502	understanding genetic interactions and predicting cellular responses to novel
4592	alternative framings by highlighting the unique challenges of handling single-
4595	cell variability, generalization to unseen perturbations, and the need for a
4594	complex model architecture. The connection between the biological objective
4595	and the technical approach is exceptionally well-articulated. The only minor
4596	have strengthened the rationale further.
4597	
4598	Analyse Baseline Delects (8/10) The report thoroughly identifies several limitations of current baseline models
4599	with exceptional specificity. For example, it correctly highlights the
4600	challenge of generalization to unseen perturbations, including new target
4601	genes and gene combinations, as a critical weakness. The critique of simpler,
4602	non-linear genetic interactions is particularly well-supported. The
4603	evaluation gaps proposed are thoughtfully grounded in the dataset's specific
4604	challenges, and the explanation of why simpler models are insufficient is
4605	detailed and well-reasoned. The only minor limitation is that some improvement directions could be more explicitly connected to specific model
4606	architectures or recent advances in the field. However, the proposed emphasis
4607	on complex, deep learning models is well-justified given the task's
4608	requirements.
4609	3. KEY STRENGTHS
4610	Comprehensive Dataset Characterization: The report provides exceptional
4611	documentation of dataset properties, perturbation types, and specific
4612	Robust Task Framing. The biological guestion is precisely mapped to an
4613	appropriate computational task with clear justification for high-dimensional
4614	regression over alternative approaches.
4615	Detailed Challenge Analysis: The analysis of modeling challenges is exceptionally
4616	task and the limitations of simpler models in addressing it.
4617	Forward-Looking Model Selection: The rationale for prioritizing predictive
4610	accuracy and employing complex models demonstrates a sophisticated
401ð	
4019	

4620	understanding of the trade-offs between model complexity and interpretability
4621	in this domain.
4623	4. AREAS FOR IMPROVEMENT
4624	dropout patterns, or perturbation effects would enhance the dataset analysis
4625	section and provide clearer guidance for subsequent modeling efforts.
4626	architectures (e.g., hybrid models that balance interpretability and
4627	complexity) could strengthen the rationale for the chosen technical approach.
4628	appropriate for assessing model performance (e.g., Pearson correlation for
4629	gene expression prediction, metrics for generalization to unseen
4030	perturbations) would provide clearer guidance for subsequent modeling efforts.
4632	5 OVERALL RECOMMENDATION
4633	This Task Analysis report provides an exceptionally strong foundation for follow-
4634	up modeling efforts. The dataset and task features are comprehensively
4635	insights for model development. The limitations of simpler models are
4636	accurately and thoroughly diagnosed, with specific, well-reasoned
4637	Justifications for employing more complex architectures. This analysis demonstrates a sophisticated understanding of both biological and
4638	computational aspects, aligning exceptionally well with current state-of-the-
4639	art practices in CRISPR-based single-cell perturbation modeling. The report would serve as an outstanding basis for advancing scientific understanding in
4640	genetic interaction networks and functional genomics research. With minor
4641	refinements in the areas noted, this analysis could be considered a model
404 <i>2</i> 4643	
4644	LIMAs Judge Output Cone
4645	LLM As Judge Output-Gene
4646	Method Design-Claude 3.7
4647	1.NUMERICAL SCORES
4648	Scientific Validity: 8/10
4649	Innovation Level: 7/10
4651	Experimental Design: 7/10
4652	
4653	2.DETAILED JUSTIFICATION Scientific Validity (8/10)
4654	The research plan demonstrates strong biological relevance by addressing the
4655	challenge of predicting post-perturbation gene expression profiles in single cells, a critical area in understanding gene regulatory networks and cellular
4656	responses to CRISPR-based interventions. The integration of VAE, GNN, and
4657	Transformer components provides a theoretically sound foundation for handling
4658	scientific understanding in single-cell biology, particularly in leveraging
4039	gene regulatory networks and attention mechanisms for interpretability.
4661	attention weights could be further elaborated to strengthen mechanistic
4662	insight.
4663	Technical Feasibility (8/10)
4664	The proposed pipeline is practical and well-detailed, with clear steps for data
4665	established libraries (scanpy, anndata) and deep learning frameworks (PyTorch
4666	) increases feasibility. The computational resource requirements appear
4667	components may demand substantial GPU memory. The plan addresses scalability
4668	through dimensionality reduction and efficient model components. The data
4669	preprocessing steps are comprehensive, though the dynamic construction of gene interaction graphs requires careful implementation to ensure biological
4670	relevance.
4672	Innovation Level (7/10)
4072	The hybrid model architecture combining VAE, GNN, and Transformer represents a
-1075	
4674	

4675	
4676	novel approach in the field of single-cell perturbation prediction. The
4677	representations offers gradings with gene expression latent
1670	approaches have been explored in other biological contexts, which slightly
4078	reduces the novelty. The plan has potential for new biological insights
4679	through attention mechanisms and graph analysis, but the unique contributions
4680	in methodology could be more clearly articulated compared to existing state-
4681	of-the-art approaches.
4682	Experimental Design (7/10)
/683	The validation methodology includes appropriate loss functions, regularization
4604	techniques, and early stopping criteria, ensuring statistical soundness. The
4084	inclusion of control samples and data augmentation techniques strengthens the
4685	experimental design. However, the plan lacks details on specific evaluation
4686	metrics (beyond loss functions) that would demonstrate biological relevance,
4687	such as correlation with observed gene expression changes or pathway
4688	perturbations for validation and testing, which is appropriate, but the plan
4689	could benefit from more explicit discussion of replicates and robustness
4600	testing.
4090	
4691	Impact Potential (8/10)
4692	The research has significant relevance to advancing single-cell biology by
4693	notential is substantial, particularly for drug discovery and therapeutic
4694	design where understanding cellular responses to genetic perturbations is
4695	critical. The approach is scalable to broader biological questions involving
4606	gene regulatory networks and could inspire follow-up research in multimodal
4090	single-cell analysis. The potential for community adoption is enhanced by the
4697	use of standard frameworks and the provision of biological interpretability
4698	Ieatures.
4699	3. KEY STRENGTHS
4700	Comprehensive integration of VAE, GNN, and Transformer components addressing
4701	different aspects of single-cell data complexity
4702	Well-structured training strategy with regularization techniques to prevent
4702	Overilling
4705	Practical data preprocessing pipeline using established single-cell tools
4704	Incorporation of expert recommendations for addressing class imbalance and
4705	improving generalization
4706	A ADEAS FOR IMPROVEMENT
4707	Enhance discussion of biological validation metrics beyond loss functions
4708	Provide more explicit details on gene regulatory network construction and
4709	updating
4710	Strengthen justification for the specific dimensions chosen for latent spaces and
4710	embeddings
4/11	consider inclusion of additional evaluation strategies such as cross-dataset
4712	Clarify how the model will handle novel cell types or contexts not present in
4713	training data
4714	
4715	5.UVERALL RECOMMENDATION This research plan represents a strong condidate for further development with
4716	high potential for scientific impact in single-cell perturbation prediction
/717	The hybrid model architecture addresses kev challenges in the field while
7/1/	maintaining biological interpretability. The next logical steps would be to
4/18	implement rigorous biological validation using additional metrics and
4719	experimental data, and to compare performance against existing state-of-the-
4720	art methods in perturbation prediction. The plan compares favorably to high-
4721	deep learning approaches and focus on biological relevance
4722	abop rearning approaches and rocas on prorogrear rerevance.
4723	
4724	
4725	
4725	
4/26	
4727	
4728	
4729	
	07
	86

### 4730 I.4. Detailed Results

To comprehensively evaluate the performance of SCAGENTS, 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 SCAGENTS and three
DeepResearch variants (OpenAI, Perplexity, and Gemini) across multiple rounds to ensure statistical robustness. Tables 7
and 8 present the averaged scores from five independent evaluation runs, providing insights into both the consistency and
performance differences across systems.

### 4738 I.4.1. TASK ANALYSIS PHASE EVALUATION

Table 7 reveals several key insights about the Task Analysis capabilities of different systems. SCAGENTS 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 SCAGENTS, achieving comparable scores in certain categories (e.g., 9.0 in drug dataset analysis). However, both Perplexity  $(DR^P)$  and Gemini  $(DR^G)$  variants 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.

Table 7: LLM evaluation of the Task Analysis phase. Three LLM judges evaluated SCAGENTS (scAg) and Deep Research
 (DR) (OpenAI, 2025) pipeline across four key capabilities and three perturbation types. SCAGENTS consistently matched or
 exceeded human expert performance, with particular strength in dataset analysis and identifying baseline model limitations.

Judges		Dr	ug			Gen	e KO			Cyte	okine	
	scAg	$DR^O$	$\mathrm{DR}^P$	$DR^G$	scAg	$DR^O$	$\mathrm{DR}^P$	$DR^G$	scAg	$DR^O$	$DR^P$	$DR^G$
					Analyse	Dataset	* ↑					
Claude3.7	8.8	9.0	6.0	7.0	8.0	7.0	3.0	5.2	7.2	7.0	4.0	5.2
R1	9.0	9.0	6.2	7.0	8.0	7.2	4.0	6.2	7.0	6.2	3.2	5.6
o1	9.0	9.0	6.0	6.8	8.2	7.8	4.4	6.0	7.2	6.0	3.0	5.2
Qwen-plus	8.0	7.2	5.2	6.2	8.0	7.0	4.0	5.4	7.0	6.8	2.8	6.0
Llama 3.1	8.2	7.4	5.2	5.8	8.8	6.8	5.0	6.2	7.8	7.2	4.0	6.0
Average	8.60	8.32	5.72	6.56	8.20	7.16	4.08	5.80	7.24	6.64	3.40	5.60
					Analyse	Task Typ	e †					
Claude3.7	8.0	7.8	6.0	5.0	8.4	8.0	6.0	6.8	6.6	4.0	3.2	4.8
R1	7.0	7.0	4.0	6.0	7.8	7.0	6.0	7.0	6.6	5.8	4.8	5.0
o1	8.6	8.2	6.0	6.6	8.8	8.0	6.0	6.0	7.0	6.2	5.0	5.2
Qwen-plus	8.0	8.0	6.2	6.2	8.2	8.0	7.0	7.4	7.0	6.2	5.8	4.2
Llama 3.1	8.2	7.6	6.2	6.0	8.8	8.0	6.0	7.2	7.2	6.2	5.0	5.0
Average	7.96	7.72	7.80	6.20	6.88	6.88	5.68	4.76	4.84	5.96	6.88	4.84
				Anc	alyse Bas	eline De	fects $\uparrow$					
Claude3.7	6.2	5.0	3.0	4.2	6.8	6.0	3.0	4.2	5.2	4.0	2.0	2.8
R1	7.0	6.0	3.6	4.2	7.0	7.0	4.8	4.2	6.0	5.2	2.0	3.4
o1	6.6	5.8	3.0	4.0	7.2	7.0	5.6	5.0	6.0	5.2	2.0	3.2
Qwen-plus	6.0	5.2	2.2	3.2	7.0	7.0	4.2	5.0	6.0	3.8	1.8	3.0
Llama 3.1	7.0	6.2	4.0	4.8	7.2	7.6	5.0	5.8	6.8	5.2	3.0	3.0
Average	6.56	5.64	3.16	4.08	7.04	6.92	4.52	4.84	6.00	4.68	2.16	3.08

4785 I.4.2. METHOD DESIGN PHASE EVALUATION

Table 8 presents a more nuanced evaluation across five dimensions of research plan quality. The results demonstrate SCAGENTS' 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 SCAGENTS, 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 SCAGENTS
in drug perturbation tasks (7.24 vs. 6.88). This could indicate that our system occasionally proposes more ambitious
but technically challenging solutions. However, SCAGENTS maintains superiority in gene knockout and cytokine tasks,
demonstrating better adaptability to diverse biological contexts.

4800 The most striking performance gap appears in *Impact Potential*, where Perplexity's DeepResearch variant scores as low as 4801 1.8 for drug perturbation tasks. This dramatic difference underscores the importance of our comprehensive approach that 4802 considers not only technical correctness but also the broader scientific implications of proposed methods.

- 4803
- 4804 I.4.3. CROSS-TASK PERFORMANCE ANALYSIS

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

This task-specific performance difference likely reflects the varying complexity of biological mechanisms involved. Gene knockouts typically produce more predictable, direct effects, while cytokine perturbations involve complex signaling cascades and cell-cell communication networks that require more sophisticated modeling approaches. The superior performance of SCAGENTS in these challenging scenarios validates our multi-agent architecture's ability to capture complex biological interactions through collaborative reasoning.

4816

### 4817 I.4.4. INTER-JUDGE AGREEMENT AND RELIABILITY

The consistency of scores across different LLM judges provides confidence in our evaluation methodology. The highest
agreement occurs in the *Innovation Level* dimension, where judges show remarkable consensus (coefficient of variation ;
0.1 for most comparisons). Greater variability appears in *Experimental Design* evaluations, possibly reflecting different
interpretations of what constitutes rigorous experimental validation in computational biology.

These detailed results collectively demonstrate that SCAGENTS not only achieves superior performance but does so consistently across different evaluation criteria, task types, and independent judges. The framework's ability to maintain high standards across all dimensionsfrom technical feasibility to scientific impactunderscores its potential as a comprehensive solution for automated scientific discovery in single-cell biology.

- 4827
- 4828
- 4829
- 4830 4831
- 4832
- 4833
- 4834
- 4835
- 4836
- 4837
- 4838
- 4839

Table 8: LLM evaluation of the Method Design phase. LLM judges assessed the quality of research plans proposed by
 scAGENTS (scAg) and Deep Research (DR) (OpenAI, 2025) pipeline across five dimensions. ScAGENTS consistently
 outperformed on scientific validity, innovation, experimental design, and impact potential.

Judges		Dı	rug			Gen	e KO			Cyte	okine	
	scAg	$DR^O$	$\mathrm{DR}^P$	$DR^G$	scAg	$DR^O$	$\mathbf{DR}^P$	$DR^G$	scAg	$DR^O$	$\mathrm{DR}^P$	DR <sup>G</sup>
					Scientifi	c Validit	y↑					
Claude3.7	7.4	8.0	3.0	4.6	7.8	7.2	3.2	5.0	6.8	6.2	2.8	4.0
R1	8.2	7.4	4.0	4.8	7.8	7.8	4.0	6.2	7.0	6.0	3.6	5.8
o1	7.8	7.0	3.4	6.2	8.2	8.4	3.6	6.2	6.8	7.0	3.0	5.2
Qwen-plus	6.8	6.4	3.0	5.8	7.6	6.8	4.0	5.4	6.6	6.6	2.6	5.0
Llama 3.1	7.0	6.4	5.2	5.8	7.8	6.8	5.0	6.8	7.2	7.0	4.4	6.0
Average	7.44	7.04	3.72	5.44	7.84	7.40	3.96	5.92	6.88	6.56	3.28	5.20
				1	Technical	Feasibi	<i>lity</i> ↑					
Claude3.7	7.0	7.0	2.4	5.2	7.0	5.8	2.6	5.8	6.4	5.8	2.2	5.6
R1	5.8	7.0	4.0	5.2	6.8	6.8	5.0	6.0	6.0	5.6	4.0	5.4
o1	7.4	8.0	4.0	6.6	8.0	7.8	5.0	6.0	7.0	6.8	4.2	5.0
Qwen-plus	7.0	6.8	3.6	5.0	7.4	5.8	3.0	6.0	6.8	6.8	3.0	5.0
Llama 3.1	7.2	7.4	4.0	6.0	8.2	6.8	4.0	5.2	6.2	5.6	5.0	5.4
Average	6.88	7.24	3.60	5.60	7.48	6.60	3.92	5.80	6.48	6.12	3.68	5.28
					Innovat	ion Leve	$l\uparrow$					
Claude3.7	8.0	7.0	5.8	6.8	8.2	7.2	5.0	6.4	7.2	6.2	4.0	6.0
R1	8.2	6.8	4.2	4.8	8.2	7.8	5.0	7.0	8.0	7.4	5.0	6.2
o1	8.0	8.2	6.0	5.0	9.0	9.0	5.2	6.8	7.6	8.0	5.0	6.0
Qwen-plus	7.6	7.4	4.2	5.6	8.0	8.0	5.0	6.6	7.2	7.0	4.0	5.2
Llama 3.1	8.4	7.6	5.0	6.2	8.0	8.0	5.2	7.0	7.2	7.2	5.0	7.0
Average	8.04	7.40	5.04	5.68	8.28	8.00	5.08	6.76	7.44	7.16	4.60	6.08
				Ε	Experimen	ntal Des	ign †					
Claude3.7	7.0	7.2	3.2	4.4	7.0	6.0	2.0	4.4	6.0	5.8	2.2	4.0
R1	8.0	7.2	4.0	4.0	7.2	6.0	2.8	4.2	6.4	6.0	3.0	4.0
01	8.2	8.4	5.0	5.2	7.2	7.0	3.0	4.8	6.8	6.8	3.4	4.0
Qwen-plus	7.2	7.0	4.2	5.8	7.8	5.0	2.0	4.0	6.0	5.0	2.4	4.0
Llama 3.1	7.8	7.2	4.2	5.0	7.2	6.2	4.0	5.0	6.8	6.0	4.0	5.0
Average	7.64	7.40	4.12	4.88	7.28	6.04	2.76	4.48	6.40	5.92	3.00	4.20
					Impact	Potentia	$l\uparrow$					
Claude3.7	6.0	5.0	1.8	3.0	6.8	5.2	2.8	4.0	6.0	5.2	3.8	4.4
	7.0	6.0	2.2	4.0	7.2	6.0	2.2	4.2	6.2	5.0	2.2	4.0
R1	= 2	7.0	3.2	4.8	4.0	7.0	6.0	3.0	6.6	6.0	2.2	4.0
R1 01	7.2							< 0				
R1 o1 Qwen-plus	7.2	6.0	2.0	3.2	6.0	5.6	2.4	6.0	6.2	5.2	2.0	4.0
R1 o1 Qwen-plus Llama 3.1	7.2 7.2 7.4	6.0 7.0	2.0 4.0	3.2 4.8	6.0 7.0	5.6 6.8	2.4 5.8	<b>6.0</b> 4.0	6.2 6.8	5.2 5.8	2.0 4.0	4.0 5.0

# 4895 J. Human scientists' evaluation details

#### 4896 4897 **J.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 Task Analysis Module, Method Design Module and confidence score in Graph-based discussion across cytokine, drug, and gene perturbation tasks.

For each task type(cytokine, drug, gene), experts evaluated 8 outputs: 5 generated by different LLM backends of SCAGENTS (Claude 3.7, o1, DeepSeek R1, Qwen-plus and Llama 3.1) and 3 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 SCAGENTS during graph-based multi-turn
 reasoning. Strong alignment between expert judgments and model confidence was observed, supporting the reliability of
 model self-evaluation.

# 49144915J.2. Detailed Results

Table 9 presents the scores given by human scientists across different tasks (cytokine, drug, and gene perturbation) for outputs from various SCAGENTS (with different LLM backends) and DeepResearch agents.

The results indicate that SCAGENTS generally outperform DeepResearch agents, with versions like scAgents-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. SCAGENTS versions show a more balanced performance compared to the DeepResearch agents, which sometimes score low in certain dimensions.

Table 10 compares the expert ratings with scAgents' confidence scores during graph-based multi-turn reasoning. The
alignment between the confidence scores and expert ratings suggests that scAgents' self-evaluation mechanism is reliable.
This correlation confirms that the confidence scores can serve as a valid indicator of the quality of scAgents' outputs. Overall,
these results highlight scAgents' effectiveness in handling scientific analysis and design tasks and validate the utility of their
confidence assessment.

4927 4928 4929

4955Table 9: SCAGENTS Performance Scores Evaluated by 3 human scientists for 10 hours. (scAg<sup>cld</sup>: scAgents-Claude3.7,4956scAg<sup>o1</sup>: scAgents-o1, scAg<sup>ds</sup>: scAgents-DeepSeek R1, scAg<sup>qw</sup>: scAgents-Qwen-plus, scAg<sup>lm</sup>: scAgents-llama 3.1, DR<sup>O</sup>:4957OpenAI DeepResearch, DR<sup>P</sup>: Perplexity DeepResearch, DR<sup>G</sup>: Gemini DeepResearch)

4959									
4960	Dimension	$\mathbf{scAg}^{cld}$	$scAg^{o_1}$	$\mathbf{scAg}^{ds}$	$\mathbf{scAg}^{qw}$	$\mathbf{scAg}^{lm}$	$\mathbf{DR}^O$	$\mathbf{DR}^P$	$\mathbf{DR}^G$
4961	Α	. Analysis R	eports by	Task Anal	ysis Modu	le			
4962	Cytokines Task Analysis								
4963	Biological Significance	7	6	7	5	6	4	0	4
4964	Gap Analysis Insight	6	2	6	4	5	2	6	2
4965	Data Characterization Accuracy	7	4	7	6	4	3	2	3
4966	Literature Integration Quality	7	4	5	6	4	3	2	2
4967	Drug Perturbation Task							-	. <u> </u>
4068	Biological Significance	7	6	7	6	5	5	3	5
4908	Gap Analysis Insight Task Formulation Clarity	8	5	6 7	6	6	5	3	5
4969	Data Characterization Accuracy	7	8	8	6	4	3	4	3
4970	Literature Integration Quality	8	8	8	6	5	6	4	3
4971	Gene Perturbation Task								
4972	Biological Significance	7	7	6	5	4	5	5	5
4973	Gap Analysis Insight	7	7	5	4	5	5	0	3
4974	Data Characterization Accuracy	8	8	6	7	3	5	5	3
4975	Literature Integration Quality	8	8	5	6	3	5	5	3
4976	B	Hypothesis	Plan by N	lethod De	sign Modu	le			
4977	Cytokines Perturbation Task								
4978	Technical Novelty	6	6	5	4	4	5	2	2
4070	Feasibility	5	6	7	5	5	5	5	5
4979	Clarity and Consistency Biological Plausibility	6	5	6 7	5	5	3	5	63
4980	Mechanism Explanation Ouality	6	5	6	5	4	1	0	2
4981	Pathway Relevance	5	6	6	3	4	3	Õ	3
4982	Cross-Perturbation Generalizability	6	7	4	5	5	5	0	5
4983	Drug Perturbation Task		_	_	_	_			
4984	Technical Novelty Feasibility	8	7	7	5	5	4	3	4
4985	Clarity and Consistency	8	8	7	6	5	4	3	4
4986	Biological Plausibility	8	8	6	5	4	4	2	3
4987	Mechanism Explanation Quality	9	8	5	4	4	4	0	2
4988	Cross-Perturbation Generalizability	9	8 8	5 6	4 5	3 4	4	1	2
4989	Gene Perturbation Task	-	-	-	-				
/000	Technical Novelty	7	7	6	5	4	5	2	4
4001	Feasibility	7	7	5	4	3	5	2	5
4771	Clarity and Consistency	9	8	6	5	4	6	3	5
4992	Mechanism Explanation Quality	7	7	3 4	4	2	4	1	4
4993	Pathway Relevance	7	7	4	3	$\frac{2}{2}$	2	0	2
4994	Cross-Perturbation Generalizability	7	7	5	4	3	2	0	2
4995	Overall Ranking	1	3	2	4	5	6	8	7
1005									

Table 10: Expert Human Scores compare with scAgents' Confidence Scores on graph-based discussions Across Tasks and Rounds

Experts		Cytoki	ne 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	0.82	0.72	0.77	0.84	0.84	0.74	0.76	0.82	0.82
Data Expert	8	9	10	10	8	10	10	10	9	9	10	10
Confidence Score	0.65	0.78	0.81	0.81	0.65	0.78	0.81	0.81	0.65	0.78	0.81	0.81
Training Expert	7	8	9	9	8	9	10	10	9	9	9	9
Confidence Score	0.65	0.80	0.81	0.81	0.77	0.80	0.82	0.82	0.78	0.80	0.81	0.81
Pathway Analyst	9	8	9	9	8	10	10	10	10	10	10	10
Confidence Score	0.77	0.79	0.83	0.83	0.77	0.81	0.81	0.80	0.77	0.85	0.86	0.86
Self Critic	8	8	9	9	8	9	10	10	9	9	9	9
Confidence Score	0.78	0.80	0.84	0.84	0.78	0.80	0.83	0.83	0.78	0.80	0.84	0.84

# 5060 K. Performance Varies Across Different LLMs and AI Coders

To comprehensively evaluate the robustness of our framework, we conducted extensive experiments comparing SCAGENTS with various baseline approaches across six challenging single-cell perturbation datasets. Table 11 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.

#### 5066 5067 K.1. Experimental Setup

5068 We evaluated four distinct categories of code generation approaches: 5069

(1) SCAGENTS with Different LLM Backends: Our full framework integrated with five state-of-the-art LLMs (Claude 3.7, OpenAI o1, DeepSeek R1, Qwen-plus, and Llama 3.1), leveraging the complete multi-agent architecture with collaborative reasoning and iterative refinement.

5073 (2) Single-LLM Direct Generation: Each LLM operating independently without the multi-agent framework, tasked with
 5074 generating the complete solution in a single pass given the same input specifications.

**(3) DeepResearch Variants:** Three commercial implementations of automated research systems, representing the current state-of-the-art in end-to-end scientific code generation.

(4) AI Coding Assistants: Two popular open-source coding frameworks (OpenHands and Aider) integrated with the same five LLMs, representing specialized code generation tools designed for software development tasks.

### 5081 5082 **K.2. Key Findings**

Table 11: Expert Human Scores compare with scAgents' Confidence Scores on graph-based discussions Across Tasks and
 Rounds

5086	Tool	Adamson	Norman	Liscovitch	Papalexi	Srivatsan	Schiebinge
5087		SCAGEN	TS with diffe	erent LLMs inte	egrated		
5088	sa Aganta Clauda 27	4	5	4	4	4	4
5089	scAgents-Claude5.7	4	3 4	4	4	4	4
5090	scAgents-DeepSeek R1	4	4	3	3	3	3
5091	scAgents-Qwen-plus	4	3	4	2	3	3
5092	scAgents-llama 3.1	2	2	1	1	1	1
5093		S	ingle-LLM g	enerated code			
5094	Claude3.7 only	2	2	1	0	1	1
5095	o1 only	1	1	0	1	1	0
5096	DeepSeek R1 only	1 1	1	0	1	1 1	0
5097	Llama 3.1 only	1	1	0	0	0	0
5098		Dee	epResearch s	generated code	?s		
5099	OpenAI DeepResearch	1	2	1	1	1	1
5100	Perplexity DeepResearch	0	0	0	0	0	0
5101	Gemini DeepResearch	0	0	0	0	0	0
5102			AI Co	oders			
5103	OpenHands-Claude3.7	3	4	3	3	2	2
5104	OpenHands-01	3	2	2	2	1	1
5105	OpenHands-DeepSeek R1	2	3	2	2	1	2
5105	OpenHands-Qwen-plus	2	2	1	2	2	2
5106	OpenHands-Llama 3.1	2	1	0	1	1	1
5107	aider-ol	2	2	2	2	1	1
5108	aider-DeepSeek R1	3	2	$\frac{1}{2}$	2	1	1
5109	aider-Qwen-plus	1	1	0	1	0	0
5110	aider-Llama 3.1	1	1	0	0	0	0

5111

5112 The results reveal several critical insights:

5113 5114 **Multi-Agent Architecture Superiority:** SCAGENTS consistently outperforms all baseline approaches, with success rates 5115 ranging from 40-100% depending on the LLM backend and dataset complexity. The multi-agent framework provides an average improvement of 2.3 over single-LLM approaches and 3.5 over AI coding assistants.

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

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

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

#### K.3. Analysis of Failure Modes

The dramatic performance gap between SCAGENTS and other approaches can be attributed to several factors:

(1) Domain Knowledge Integration: Single-LLM approaches often generate syntactically correct but biologically mean-ingless 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 SCAGENTS 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 further optimize scientific AI systems. 

# 5170 L. Designed Models

5171 5172 To understand how SCAGENTS adapts its architectural choices to different biological contexts, we conducted a post-hoc 5173 analysis of the model components selected across various perturbation tasks. Table 12 presents the frequency of different 5174 neural network components appearing in models automatically designed by SCAGENTS across multiple independent runs 5175 for each dataset.

### 5176 5177 **L.1. Methodology**

5178 For each of the six benchmark datasets, we executed SCAGENTS five times with different random seeds and analyzed the 5179 resulting model architectures. Each value in Table 12 represents the number of times a specific component appeared in the 5180 generated models. This analysis reveals how our framework naturally adapts its design choices to the unique characteristics 5181 of each perturbation type and data modality.

Table 12: Type of model components automatically selected by SCAGENTS across different perturbation tasks. Each
value represents the frequency of component usage across five independent runs. The adaptive selection of architectures
demonstrates the framework's ability to tailor solutions to specific biological contexts.

5187	Model Type	Transformer	Attention	GNN	VAE	MLP	RNN	CNN	GAN	XGBoost
5180	Gene Knock Out-RNAseq (Adamson Dataset)	5	4	4	3	4	-	1	2	-
5105	Gene Knock Out-RNAseq (Norman Dataset)	4	6	4	2	5	1	1	2	-
5190	Gene Knock Out-ATACseq (Liscovitch Dataset)	3	2	2	1	2	-	-	2	1
5191	Gene Knock Out-CITEseq (Papalexi Dataset)	6	3	4	2	1	-	-	1	2
5192	Drug Perturbation (Srivatsan Dataset)	3	1	1	1	2	-	-	2	-
5193	Cytokine Perturbation (Schiebinger Dataset)	4	1	1	1	2	-	-	1	1

# 5196 L.2. Architecture Selection Patterns

5197 5198 Several notable patterns emerge from this analysis:

Transformer Dominance in Gene Expression Tasks: Transformer architectures appear most frequently in traditional
 scRNA-seq datasets (Adamson: 5/5, Norman: 4/5), reflecting their effectiveness in capturing long-range gene-gene
 dependencies. The self-attention mechanism enables modeling of complex regulatory relationships without explicit prior
 knowledge of gene interaction networks.

Attention Mechanism Prevalence: Beyond full transformer architectures, standalone attention mechanisms are particularly
 favored for the Norman dataset (6 occurrences), suggesting that the multi-gene perturbation patterns in this dataset benefit
 from selective information aggregation. This aligns with the biological reality that combinatorial perturbations often involve
 non-linear interactions requiring dynamic weighting of different genes' contributions.

Graph Neural Networks for Regulatory Modeling: GNNs appear consistently across gene knockout experiments (4/5 for
 both Adamson and Norman datasets), indicating that SCAGENTS recognizes the value of explicitly modeling gene regulatory
 networks. The slightly lower frequency in drug and cytokine perturbations (1/5) suggests that molecular interaction networks
 may be less directly applicable to these perturbation types.

**Modality-Specific Adaptations:** The scATAC-seq dataset (Liscovitch) shows unique patterns with the inclusion of XGBoost (1/5), the only dataset to utilize this traditional machine learning approach. This likely reflects the extreme sparsity of chromatin accessibility data, where gradient boosting can effectively handle the binary nature of peak calling.

5216

5194 5195

## 5217 L.3. Biological Interpretation of Architectural Choices

The architectural diversity reflects SCAGENTS' understanding of different biological mechanisms:

5220 **Multi-Modal Integration:** The CITE-seq dataset (Papalexi) shows the highest transformer usage (6/5 possible due to 5221 hybrid architectures), consistent with the need to integrate RNA and protein measurements. The lower MLP frequency (1/5) 5222 suggests that simple concatenation approaches are insufficient for this multi-modal challenge.

5223 5224 **Temporal Dynamics:** The appearance of RNNs exclusively in the Norman dataset (1/5) indicates recognition of potential temporal aspects in perturbation responses, even though the data itself is not explicitly time-series. This may reflect the framework's attempt to model the cascade effects of genetic perturbations.

Generative Modeling: GANs appear sporadically across datasets (ranging from 1-2 occurrences), primarily in scenarios where data augmentation might benefit model training. Their higher frequency in scATAC-seq tasks suggests utility in handling the severe class imbalance inherent in chromatin accessibility data. 

#### L.4. Task Complexity and Model Sophistication

A clear correlation exists between task complexity and architectural sophistication:

Simple Perturbations, Simple Models: Drug and cytokine perturbations, which typically involve more direct molecular mechanisms, show lower architectural diversity. The predominant use of transformers (3-4/5) with minimal additional components suggests that these tasks can be adequately modeled with standard sequence-to-sequence approaches. 

Complex Biology, Complex Architecture: Gene knockout experiments, particularly those involving combinatorial perturbations, consistently employ multiple architectural components in tandem. The co-occurrence of transformers, attention mechanisms, and GNNs in these models reflects the need to capture both local (gene-specific) and global (network-wide) effects.

**Sparse Data, Specialized Solutions:** The unique challenges of scATAC-seq dataextreme sparsity (*j*95% zeros) and binary naturelead to distinctive architectural choices. The reduced reliance on VAEs (1/5) and increased use of discriminative models reflects the difficulty of learning meaningful latent representations from such sparse data. 

#### L.5. Implications for Automated Model Design

This analysis demonstrates that SCAGENTS has developed an implicit understanding of the relationship between biological data characteristics and appropriate model architectures. The framework's ability to consistently select transformers for sequence modeling, GNNs for network effects, and specialized components for unique challenges validates our multi-agent approach to scientific discovery. 

Moreover, the architectural diversity observed across different perturbation types underscores the limitations of one-size-fits-all approaches in computational biology. The success of SCAGENTS lies not just in its ability to generate working code, but in its capacity to reason about the fundamental nature of each biological problem and select appropriate computational tools accordingly. 

These findings suggest that effective automation of scientific model design requires not just technical expertise but also deep integration of domain knowledgea capability that emerges naturally from our multi-agent collaborative framework.

# 5280 M. Cost Analysis

5281
 5282
 5283
 5284
 5284
 5285
 5286
 5286
 5287
 5288
 5288
 5284
 5284
 5284
 5284
 5284
 5284
 5285
 5286
 5286
 5286
 5287
 5286
 5288
 5288
 5289
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 5284
 <li

### 5285 5286 **M.1. Training Infrastructure**

All models designed by SCAGENTS, with parameter counts ranging from 10 million to 30 million, were trained and evaluated
 on a consistent computational environment. These relatively lightweight modelscompared to large language models with
 billions of parametersdemonstrate that effective scientific modeling does not necessarily require massive computational
 resources.

### 5292 Hardware Specifications:

5293 5294

5295

5299

5300 5301 • GPU: Two NVIDIA H20-NVLink with 96GB VRAM each (192GB total)

- **CPU:** 16-core AMD EPYC 9K84 @ 2.6GHz
- 5297 5298 • **Memory:** 150GB DDR5 RAM
  - Storage: 2TB NVMe SSD for dataset caching

5302 This setup facilitated stable multi-GPU training and inference processes without encountering memory bottlenecks. The 5303 dual-GPU configuration enabled:

#### 5304 5305

5310

5315

5316 5317

5318

5331

5332 5333

5334

- Parallel training with data parallelism for larger batch sizes 5306
- 5307 Distributed evaluation across multiple perturbation conditions
   5308
- Simultaneous model development for different datasets

Training Time Analysis: On average, models converged within 4-8 hours depending on dataset size and complexity:

- 53135314 Gene knockout (RNA-seq): 4-5 hours
  - Multi-modal (CITE-seq): 6-8 hours
  - Sparse data (ATAC-seq): 3-4 hours

### 5319 5320 **M.2. Token Utilization and Cost Estimation**

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

5325 M.2.1. TOKEN USAGE BREAKDOWN BY PHASE

The specific token usage varies significantly based on task complexity. Our empirical analysis across 50+ experiments
 revealed the following patterns:

5329 For cost estimation purposes, we use the observed average: 5330

- Prompt tokens (input): 60,000
- Completion tokens (output): 300,000

Submission and Forma	atting Instruction	s for ICML	2025 GenBio	Workshop
----------------------	--------------------	------------	-------------	----------

53	340
53	841
53	342
53	343
53	344
53	345
53	346
53	347
53	348
53	349
53	350
53	851
53	352
53	353
53	354
53	855
53	856
53	857
53	358
53	359
53	860
53	861
53	862
53	863
53	864
53	865
53	866
53	867
53	868
53	869
53	370
53	871
53	872
53	373
53	374
53	875
53	876
53	877
53	378
53	379
53	880
53	381

Table 13: Token usage breakdown by framework phase and task complexity

Phase	Simpl	e Tasks	<b>Complex Tasks</b>			
	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		

### 46 M.2.2. PER-REQUEST COST CALCULATION

Given that vendors report token pricing per million tokens (\$/M), the cost per request was computed using:

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

3 Table 14: Per-million-token pricing and per-call cost estimates based on average usage (60K input and 300K output tokens)

Model	Prompt (\$/M)	Completion (\$/M)	Cost per Request (\$)
Claude 3.7 (Anthropic)	3.00	15.00	4.68
OpenAI o1	15.00	60.00	18.90
DeepSeek-R1	0.27	2.19	0.67
Qwen-Plus	0.40	1.20	0.38
LLaMA 3.1	3.50	3.50	1.26
Average	_	-	5.18

### M.3. Total Experimental Costs

To provide practical cost estimates, we analyze the total expense for complete experimental workflows:

Table 15: Total cost estimates for different experimental scenarios

Scenario	Runs	LLM	Cost/Run	Total Cost
Single experiment	1	Claude 3.7	\$4.68	\$4.68
Hyperparameter search	10	DeepSeek-R1	\$0.67	\$6.70
Full benchmark (6 datasets)	30	Qwen-Plus	\$0.38	\$11.40
Production deployment	100	Mixed	\$5.18	\$518.00

#### M.4. Cost-Effectiveness Analysis

79 When compared to traditional approaches, SCAGENTS demonstrates significant cost advantages:

vs. Human Expert Time: A skilled bioinformatician typically requires 40-80 hours to develop and validate a novel
 perturbation prediction model. At standard rates (\$75-150/hour), this translates to \$3,000-12,000 per model. SCAGENTS achieves comparable or superior results for \$5-20.

vs. Cloud Computing: Traditional hyperparameter optimization on cloud GPUs (e.g., AWS p3.2xlarge at \$3.06/hour) for
 100 configurations would cost approximately \$1,200-2,400. Our approach identifies optimal architectures through reasoning
 rather than exhaustive search.

vs. Consulting Services: Commercial bioinformatics consulting for custom model development typically costs \$10,000 50,000 per project. SCAGENTS provides similar capabilities at ¡0.1% of the cost.

### 5390 M.5. Cost Optimization Strategies

Based on our extensive experimentation, we recommend the following strategies to minimize costs while maintaining
 performance:

5394 1. Tiered LLM Usage: Use expensive, high-performance models (e.g., OpenAI o1) only for critical reasoning steps in
 5395 Method Design, while employing cost-effective alternatives (e.g., DeepSeek-R1) for routine code generation.

5396
 5397
 **2. Caching and Reuse:** Implement aggressive caching of Task Analysis results, as dataset characteristics remain constant across experiments. This can reduce costs by 30-40% for repeated experiments.

**3. Early Termination:** Monitor confidence scores during graph-based discussions and terminate when consensus is reached, 5400 potentially saving 20-30% of Method Design tokens.

4. Batch Processing: When analyzing multiple datasets, batch similar tasks together to leverage shared context and reduce
 redundant reasoning.

### **M.6. Environmental Considerations**

5406 While focusing on monetary costs, we acknowledge the environmental impact of extensive LLM usage. Based on published 5407 estimates, our average experiment generates approximately 0.5-2.0 kg CO equivalent, depending on the data center's energy 5408 source. This is comparable to a 5-20 km car journey, but significantly less than training a large neural network from scratch 5409 (100-1000 kg CO).

# 54115412M.7. Return on Investment

5413 Despite the per-experiment costs, SCAGENTS offers compelling returns:

- Time Savings: 40-80 hours of expert time reduced to 4-8 hours of computation
- Quality Improvements: 15-49% better prediction accuracy than baseline methods
- Reproducibility: 100% reproducible results with complete audit trails
- Accessibility: Enables small research groups without ML expertise to conduct cutting-edge analyses

This analysis demonstrates that while SCAGENTS incurs non-trivial API costs, it remains highly cost-effective compared
 to traditional approaches. The framework democratizes access to sophisticated computational biology methods, enabling
 broader participation in scientific discovery at a fraction of traditional costs.

### 5445 **N. Failure Case Analyses**

5452

5453

In this section, we analyze common failure modes of SCAGENTS across various single-cell perturbation tasks. We manually reviewed 20 randomly selected failed experiment cases generated by SCAGENTS across cytokine, drug, and gene response prediction scenarios. Based on qualitative inspection of agent behaviors and outputs, we identified 7 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.

Table 16: Failure Types of SCAGENTS code generation

Failure Type	Definition & Examples
Model Configuration Error	The agent misconfigures the model architecture or fails to define required hyper- parameters. This includes mismatched layer dimensions, invalid GNN configura- tions, 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 that are 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, interpre- tations) that are not grounded in the available data or context. This includes fabricating values, inventing data structures or statistical conclusions, or reason- ing 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 9, 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



5530 specific modeling or reasoning tasks. Error Recovery Failure comprises 16% of the failures, where the agent fails to adapt to 5531 execution failures. Hallucination makes up 4% of the errors, where outputs are not grounded in available data or context.

5532 Notably, we found that implementing code to print array or matrix shapes can aid SCAGENTS in subsequently reading 5533 the command region's output for modification, thereby enhancing their ability to identify and resolve shape-related issues 5534 during tensor operations. This approach proved particularly effective given the complexity of data processing workflows 5535 in scAgents. Even though SCAGENTS utilize a data parser to obtain the original dataset dimensions and incorporate data 5536 experts during the graph-based discussion phase, the subsequent data splitting and complex model processing steps often 5537 introduce intricate dimension transformations. These transformations can lead to matrix dimension mismatches, especially 5538 when handling dynamic data structures or applying multi-layered model architectures. According to the chart, 48% of the 5539 errors in the Computation Execution Error category have been mitigated by allowing the agent to read the printed array or 5540 matrix shapes from the command output and adjust accordingly. This self-debugging capability significantly enhances the 5541 agent's ability to resolve shape-related issues during tensor operations, improving overall system robustness. Figure 10 5542 provides an example of the printed data shapes during tensor operations, which SCAGENTS can utilize to dynamically adjust 5543 and correct dimension mismatches. 5544

```
5545
       Using device: cuda
       Loading data...
5546
       Number of perturbation types in training set: 189
5547
       Number of perturbation types in test set: 48
5548
5549
       [I 2025-05-12 23:07:44,902] A new study created in memory with name: no-name-93169e84-3cf2-44c8-963d-c5a725c457ca
5550
       Training data shape: (85536, 840)
5551
       Test data shape: (18930, 840)
5552
                                Figure 10: A probable example of printing array or matrix shapes.
5553
5554
```

# **O. UMAP of Perturbation Results**

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

#### **O.1. Visualization Methodology**

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

- <sup>5</sup> 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)

## O.2. Results and Interpretation

Figure 11 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 (Norman et al., 2019)), drug perturbation
(Srivatsan et al. Dataset (Srivatsan et al., 2020)), and cytokine perturbation (Schiebinger et al. Dataset (Schiebinger et al., 2019)).



Figure 11: **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.

## O.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 SCAGENTS successfully learned the gene regulatory logic.

5610 O.2.2. DRUG PERTURBATION ANALYSIS 5611 The drug perturbation results (middle panel) reveal a more complex landscape: 5612 5613 • Global Structure: The overall "comet-like" shape is well-preserved, indicating successful capture of major drug 5614 5615 response trajectories 5616 • Increased Dispersion: Predicted cells show slightly more spread than ground truth, particularly in transition regions 5617 5618 5619 • Gradient Effects: The model captures the continuous nature of dose-response relationships, visible as smooth transitions rather than discrete clusters 5620 5621 5622 The increased variability in drug responsesdue to factors like off-target effects and cell-specific metabolismpresents a greater 5623 challenge that our model handles reasonably well. 5624 5625 **O.2.3.** CYTOKINE PERTURBATION ANALYSIS 5626 The cytokine perturbation visualization (right panel) shows remarkably tight correspondence despite the inherent complexity 5627 of immune signaling: 5628 5629 5630 • Circular Organization: Both predicted and ground truth cells form a characteristic circular pattern, likely representing 5631 cell cycle or differentiation trajectories 5632 5633 • Uniform Coverage: The model achieves uniform coverage across the entire manifold without gaps or over-densification 5634 5635 • Fine Structure: Subtle substructures within the main circular pattern are preserved, indicating capture of nuanced 5636 biological states 5637 5638 **O.3.** Quantitative Assessment 5639 5640 To complement the visual analysis, we computed several quantitative metrics on the UMAP embeddings: 5641 5642 Table 17: Quantitative metrics for UMAP embedding similarity 5643 5644 Metric Gene Drug Cytokine 5645 **Procrustes Distance** 0.12 0.18 0.14 5646 Centroid Distance 0.08 0.15 0.10 5647 KL Divergence 0.09 0.16 0.11 5648 Silhouette Score (Overlap) 0.92 0.84 0.89 5649 5650 5651 These metrics confirm the visual observations: gene perturbations show the highest fidelity (lowest distances), while drug 5652 perturbations exhibit more variability. All values indicate strong overall correspondence between predicted and ground truth 5653 distributions. 5654 5655 **O.4. Biological Significance** 5656 The UMAP visualizations reveal that SCAGENTS captures not just individual gene expression values but also: 5657

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

### **O.5. Limitations and Considerations**

5667 While these visualizations provide compelling evidence of model quality, several caveats should be noted:

- UMAP Parameters: Different parameter choices can affect the visual appearance while preserving the same underlying relationships
- Projection Artifacts: Some apparent differences may be artifacts of the 2D projection rather than true prediction errors
- Sampling Effects: For visualization clarity, we show a random subset of 5,000 cells per condition

### **O.6. Implications for Model Development**

5677 The UMAP analysis provides several insights for future model improvements:

5678
 5679
 5679
 5679 benefit from specialized model components

- 5681 2. Uncertainty Quantification: Regions with lower overlap could guide uncertainty estimation mechanisms
- 3. Biological Constraints: Incorporating known constraints (e.g., cell cycle boundaries) could improve predictions in ambiguous regions
- 5685 These visualizations ultimately demonstrate that SCAGENTS successfully generates models that capture both fine-grained 5686 expression patterns and global transcriptional landscapes across diverse perturbation types, validating our approach for 5687 automated scientific discovery in single-cell biology.

# 5720 P. Additional Visualizations

### P.1. Comparative Performance Analysis

To provide a comprehensive visual assessment of SCAGENTS' performance advantages, Figure 12 presents comparative bar charts across three evaluation dimensions for different perturbation types. These visualizations offer complementary insights to the numerical results in Tables 7 and 8.



Figure 12: **Comparative evaluation of SCAGENTS 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. SCAGENTS (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.** SCAGENTS achieves consistently high scores (7.2-8.6) across all perturbation types, 5745 demonstrating robust capability in extracting and interpreting complex biological data characteristics. The most significant 5746 advantage appears in drug perturbation analysis (8.6), where our multi-agent approach effectively handles the complexity of 5747 chemical-biological interactions.

**Task Type Understanding.** While baseline methods show variable performance (3.2-8.0), SCAGENTS 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. SCAGENTS 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, SCAGENTS 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 SCAGENTS' superiority extends beyond marginal improvements represents a funda mental 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.

#### **O. Related Works** 5775

5801

5776 This section provides a comprehensive overview of existing approaches in automated scientific discovery and single-cell 5777 perturbation analysis, positioning SCAGENTS within the broader landscape of AI-driven research systems. 5778

#### 5779 Q.1. Agent Systems for Scientific Discovery 5780

5781 The landscape of agent systems for scientific discovery has evolved rapidly with the advancement of LLMs, creating a 5782 continuum of capabilities that span different aspects of the scientific workflow. As summarized in Table 18, existing systems 5783 can be categorized into several distinct classes based on their primary functionalities. 5784

Literature Analysis Systems. Task analysis systems form the foundation of this ecosystem, with tools like PaperQA2 5785 (Skarlinski et al., 2024) and CHIME (Hsu et al., 2024) demonstrating sophisticated literature retrieval and organization 5786 capabilities. These systems excel at extracting insights from diverse sources while maintaining accuracy and contextualizing 5787 information within broader scientific frameworks. Evaluation frameworks such as ScienceAgentBench (Chen et al., 2024) 5788 and DSBench (Jing et al., 2024) have emerged to systematically assess these capabilities, providing standardized benchmarks 5789 for scientific task decomposition and comprehension abilities. However, these systems remain limited to analysis without 5790 the ability to generate novel hypotheses or implement solutions. 5791

5792 Hypothesis Generation Systems. Building on these foundational capabilities, hypothesis generation systems represent a 5793 more creative dimension of scientific discovery. ResearchAgent (Baek et al., 2024) implements an iterative approach to 5794 research ideation by progressively refining hypotheses through literature analysis. Complementary research has investigated 5795 whether LLMs possess inherent scientific reasoning abilities, with Qi et al. (Qi et al., 2023) demonstrating that models can 5796 generate plausible hypotheses in zero-shot contexts. The quality of these AI-generated research proposals has been critically 5797 assessed through human evaluation studies (Si et al., 2024), revealing that while LLM-generated ideas often meet basic 5798

5799 Table 18: Comparison of Agent Systems for Scientific Discovery. The table highlights the comprehensive capabilities of 5800 SCAGENTS across all evaluated dimensions versus existing methods.

Method	Capability								
	Task Analysis	Hypothesis Gen.	Code Gen.	End-to-End Automation	Multi-Agent System	Domain Knowledge	Evaluation		
		Literature Analysis	Systems						
PaperQA2 (Skarlinski et al., 2024)	1	×	×	×	×	1	×		
CHIME (Hsu et al., 2024)	✓	×	×	×	×	1	×		
PaperBench(Starace et al., 2025)	✓	×	×	×	×	1	1		
	Н	ypothesis Generatio	n Systems						
ResearchAgent (Baek et al., 2024)	<ul> <li>Image: A second s</li></ul>	1	×	×	1	1	×		
VirSci (Su et al., 2024)	×	1	×	×	1	1	×		
CoI Agent (Li et al., 2024a)	✓	1	×	×	1	1	×		
OpenD5 (Zhong et al., 2023)	✓	1	×	×	×	1	1		
DeepResearch (OpenAI, 2025)	✓	1	<ul> <li>Image: A set of the set of the</li></ul>	<ul> <li>Image: A set of the set of the</li></ul>	×	1	×		
		Code Generation S	Systems						
SciCode (Tian et al., 2024)	1	×	1	×	×	1	X		
DA-Code (Huang et al., 2024b)	✓	×	1	×	×	1	×		
MLAgentBench (Huang et al., 2024a)	✓	×	1	×	×	1	×		
DiscoveryBench (Majumder et al., 2024)	✓	1	1	×	×	1	1		
BLADE (Gu et al., 2024)	✓	1	<ul> <li>Image: A second s</li></ul>	×	×	1	<ul> <li>Image: A second s</li></ul>		
	1	End-to-End Scientifi	c Systems						
AI Scientist (Lu et al., 2024)	<ul> <li>Image: A second s</li></ul>	1	1	1	1	1	×		
MLR-Copilot (Li et al., 2024b)	✓	1	1	1	1	×	×		
Future House(House, 2024)	✓	1	<ul> <li>Image: A second s</li></ul>	✓	1	1	✓		
		Domain-Specific S	ystems						
BioDiscoveryAgent (Roohani et al., 2024b)	<ul> <li>Image: A second s</li></ul>	1	1	×	×	1	1		
Coscientist (Boiko et al., 2023)	1	1	1	1	×	1	×		
AtomAgents (Ghafarollahi & Buehler, 2024)	1	1	1	1	1	1	×		
TAIS (Liu et al., 2024a)	$\checkmark$	1	1	1	1	1	1		
MedAgents (Tang et al., 2024)	<ul> <li>Image: A second s</li></ul>	×	×	×	1	<ul> <li>Image: A second s</li></ul>	×		
scAgents					1	1	1		

- 5830 quality thresholds, they frequently lack the creative insights characteristic of high-impact human research. Commercial
- 5831 platforms like DeepResearch (OpenAI, 2025) have operationalized these capabilities, offering integrated solutions for 5832 literature analysis and hypothesis development across disciplines
- 5832 literature analysis and hypothesis development across disciplines.

Code Generation and Validation Systems. The implementation gap between conceptual models and executable experiments has been addressed through specialized code generation frameworks. SciCode (Tian et al., 2024) and DA-Code (Huang et al., 2024b) establish benchmarks for evaluating code quality across scientific domains, while MLAgentBench (Huang et al., 2024a) focuses specifically on machine learning experimentation workflows. More advanced systems like DiscoveryBench (Majumder et al., 2024) and BLADE (Gu et al., 2024) integrate hypothesis generation with code implementation, providing a more comprehensive evaluation of scientific reasoning. These tools ensure that implementation aligns to be the provided and the problem to be provi

with scientific best practices, bridging the gap between theoretical models and practical applications.

5841 The refinement and validation of hypotheses represent a crucial step where several innovative approaches have emerged.

5842 Chain of Ideas (Li et al., 2024a) implements a sequential framework that builds upon previous insights while exploring novel

5843 directions through structured reasoning. Goal-Driven Discovery (Zhong et al., 2023) demonstrates how language-guided

5844 exploration can enhance hypothesis development by highlighting statistically significant relationships in complex datasets.

5845 These developments address a critical gap identified by Honovich et al. (Honovich et al., 2022), who highlighted the

5846 limitations of current prompting methods for complex scientific tasks.

5847

### 5848 Q.2. End-to-End Scientific Discovery Systems

Recent developments have focused on integrating these capabilities into comprehensive systems that manage the entire scientific workflow. The AI Scientist (Lu et al., 2024) represents one of the most ambitious attempts at fully automated scientific discovery, integrating literature analysis, hypothesis generation, experimentation, and result interpretation into a cohesive framework. Similar end-to-end approaches include MLR-Copilot (Li et al., 2024b) for machine learning research and Agent Laboratory (Schmidgall et al., 2025), which employs LLM agents as research assistants throughout the scientific process.

A particularly innovative development is Agentrxiv (Schmidgall & Moor, 2025), which creates a collaborative framework for autonomous research through a shared preprint server, demonstrating that knowledge sharing among agents significantly improves performance on benchmark tasks. These systems represent significant progress toward full automation but often lack domain-specific knowledge integration, limiting their effectiveness in specialized fields like single-cell biology.

5860 Domain-Specific Implementations. Domain-specific implementations have addressed the unique challenges of particular 5861 scientific fields. Autonomous Chemical Research (Boiko et al., 2023) integrates chemical knowledge with laboratory 5862 automation to accelerate compound discovery, while BioDiscoveryAgent (Roohani et al., 2024b) focuses on designing genetic 5863 perturbation experiments. Systems like AtomAgents (Ghafarollahi & Buehler, 2024) for materials science, MedAgents 5864 (Tang et al., 2024) for medical reasoning, and DORA AI Scientist (Naumov et al., 2025) for general scientific exploration 5865 demonstrate how domain knowledge can be effectively incorporated into agent-based systems. Perhaps most impressively, 5866 Sparks (Ghafarollahi & Buehler, 2025) has discovered previously unknown phenomena in protein mechanics without human 5867 intervention, highlighting the potential for truly autonomous scientific discovery. 5868

**Specialized Methodologies and Enhancements.** Specialized methodologies have emerged to enhance specific aspects of the scientific discovery process. NOVA (Hu et al., 2024) employs iterative planning and search to improve the novelty and diversity of generated ideas, while Scideator (Radensky et al., 2024) grounds idea generation in research paper facet recombination. "Literature Meets Data" (Liu et al., 2024b) demonstrates the advantages of combining literature-based insights with empirical data, outperforming approaches that rely exclusively on either source. CODESCIENTIST (Jansen et al., 2025) utilizes genetic search over combinations of research text and executable code for scientific ideation and experiment, enabling large-scale automated discovery with rigorous evaluation across code, review, and replication.

The quality and organization of literature synthesis have been enhanced through systems like SurveyForge (Yan et al., 2025), which employs outline heuristics and memory-driven generation for automated survey writing. Novel evaluation frameworks such as AI Idea Bench (Qiu et al., 2025) and LiveIdeaBench (Ruan et al., 2024) provide standardized metrics for assessing LLMs' scientific creativity and idea generation capabilities, while RAGBench (Friel et al., 2025) offers an explainable benchmark for retrieval-augmented generation systems specifically designed for scientific literature.

Maintaining methodological rigor in automated experimentation has been addressed by Curie (Kon et al., 2025), which
 embeds rigor through modules enhancing reliability, methodical control, and interpretability. The ethical dimensions of AI

5885 in scientific research have been explored by Lin et al. (Lin, 2024), who proposes practical strategies for ethical AI use that 5886 move beyond abstract principles to address concrete implementation challenges.

5887

### 5888 Q.3. Comparison of SCAGENTS with Existing Approaches 5889

Table 19: **Comparison of computational approaches for single-cell analysis.** The table highlights the comprehensive capabilities of SCAGENTS across all evaluated dimensions versus existing methods grouped by categories.

	Dimension								
Method	Automation	Agentic	Domain Knowl.	Cross Modal.	Interpret.	Code Gen			
	Ti	raditional .	ML						
Linear Regression (Dixit et al., 2016)	×	×	×	×	1	×			
Random Forest (Skinnider et al., 2021)	×	×	×	×	1	×			
	Deep (	Generative	Models						
scGen (Lotfollahi et al., 2019)	×	×	<ul> <li>Image: A set of the set of the</li></ul>	×	×	×			
Perturb-CGAN (Hetzel et al., 2022)	×	×	1	×	×	×			
	Netwo	rk-Based I	Methods						
Dynamo (Qiu et al., 2022)	×	×	1	<ul> <li>✓</li> </ul>	1	×			
AttentionPert (Bai et al., 2024)	×	×	1	×	✓	×			
GRNBoost2 (Moerman et al., 2019)	×	×	<ul> <li>Image: A set of the set of the</li></ul>	×	<ul> <li>Image: A second s</li></ul>	×			
	Transfo	ormer Arch	vitectures						
scGPT (Cui et al., 2024)	×	X	<ul> <li>Image: A set of the set of the</li></ul>	<b>√</b>	X	×			
Geneformer (Theodoris et al., 2023)	×	×	1	$\checkmark$	×	×			
scBERT (Yang et al., 2022)	×	×	$\checkmark$	✓	×	×			
	Existi	ng Agent S	Systems						
BioDiscoveryAgent (Roohani et al., 2024b)	×	1	<ul> <li>✓</li> </ul>	×	×	1			
DiscoveryBench (Majumder et al., 2024)	×	1	1	×	×	1			
ScienceAgentBench (Chen et al., 2024)	<i>✓</i>	1	$\checkmark$	✓	×	1			
			1			/			

5912

5931

5932

5934

5913 While general scientific discovery systems have made significant progress, single-cell perturbation analysis presents unique 5914 challenges that require specialized approaches. Table 19 compares SCAGENTS with existing methods specifically designed 5915 for single-cell analysis across six critical dimensions. 5916

**Traditional and Deep Learning Approaches.** Traditional machine learning methods like Linear Regression (Dixit et al., 2016) and Random Forest (Skinnider et al., 2021) offer interpretability but lack the sophistication to capture complex gene regulatory relationships. Deep generative models such as scGen (Lotfollahi et al., 2019) and Perturb-CGAN (Hetzel et al., 2022) improve prediction accuracy but operate as black boxes without domain knowledge integration or automation capabilities.

Network-Based and Transformer Methods. Network-based approaches including Dynamo (Qiu et al., 2022), AttentionPert
(Bai et al., 2024), and GRNBoost2 (Moerman et al., 2019) incorporate biological knowledge through gene regulatory
networks, enhancing interpretability. However, they require manual implementation and lack cross-modal capabilities.
Recent transformer architectures like scGPT (Cui et al., 2024), Geneformer (Theodoris et al., 2023), and scBERT (Yang
et al., 2022) demonstrate impressive cross-modal learning but still require expert knowledge for deployment and lack
automated discovery capabilities.

**The SCAGENTS Advantage.** As shown in both tables, SCAGENTS is the only system that achieves comprehensive coverage across all evaluated dimensions. This unique position stems from our integration of:

- End-to-end automation that eliminates the need for manual intervention
- Multi-agent collaboration that leverages specialized expertise
- **Domain knowledge integration** through agentic retrieval and expert reasoning
- 59365937 Cross-modal capabilities for handling diverse single-cell data types
- **Built-in interpretability** through transparent decision-making processes
## Submission and Formatting Instructions for ICML 2025 GenBio Workshop

5940 5041	• Automatic code generation that produces executable, optimized implementations
5942	This comprehensive approach enables SCACENTS to address the full complexity of single call parturbation analysis while
5943	maintaining the flexibility to adapt to new biological contexts and data modalities. The framework's success in outperforming
5944	specialized models like scGPT while providing complete automation demonstrates the power of combining domain-specific
5945	specialized models like score reasoning complete automation demonstrates the power of comonning domain-specific knowledge with general-purpose reasoning canabilities in a unified system
5946	knowledge with general purpose reasoning capabilities in a diffied system.
5947	
5948	
5949	
5950	
5951	
5952	
5953	
5954	
5955	
5956	
5957	
5958	
5959	
5960	
5062	
5063	
5967	
5965	
5966	
5967	
5968	
5969	
5970	
5971	
5972	
5973	
5974	
5975	
5976	
5977	
5978	
5979	
5980	
5082	
5002	
508/	
5085	
5986	
5987	
5988	
5989	
5990	
5991	
5992	
5993	
5994	