Overcoming Local Optima in Waddington Landscape: A Unified Machine Learning Framework for Single-Cell Perturbation Prediction and Antigen Presentation

Francis Boabang*

INTRODUCTION

The ability of cells to respond to genetic, chemical, and antigenic perturbations is fundamental to progress in precision medicine, drug discovery, and immunotherapy [1]–[3]. However, despite significant advances in data-driven biology, existing computational models often fall short in capturing the intricate, nonlinear dynamics of cell fate decisions and immune recognition. These limitations are frequently due to suboptimal model initialization, convergence to local optima, and a lack of biological interpretability.

This project aims to develop a next-generation artificial intelligence framework that combines reinforcement learning, transformer-based architectures, and biological multiomic data integration to model cellular and immune system responses with unprecedented accuracy and interpretability. The research will focus on two synergistic objectives.

First, we will design a multistage reinforcement learning algorithm that models single-cell responses to genetic and chemical perturbations using data from Single Cell RNA Sequencing (scRNA-seq), Cellular Indexing of Transcriptomes and Epitopes by sequencing (scCITE-seq) and Single Cell ATAC Sequencing (scATAC-seq). By integrating natural gradient trust-region optimization with proximal policy refinement, this approach will overcome local minima and better capture true differentiation trajectories within the cell's Waddington landscape.

Second, we will develop a machine learning framework for Class II antigen recognition, enabling accurate prediction of peptide–Major Histocompatibility Complex interactions that underpin immune activation. Using bidirectional neural encoding, transformer attention mechanisms, and a reinforcement-based feedback loop, the model will iteratively improve its predictions based on biological relevance, enhancing interpretability and reliability for vaccine and immunotherapy applications.

Together, these innovations will yield a unified, data-driven framework that not only predicts how cells and immune systems respond to perturbations but also provides deeper mechanistic insight into underlying biological processes. Expected outcomes include improved accuracy and generalization in single-cell perturbation prediction, enhanced understanding of antigen presentation and immune activation, and a scalable computational platform for translational research in drug design and precision medicine.

The proposed framework can offer personalized therapeutic prescription for patients with degenerative conditions by integrating cell regeneration strategies such as induced pluripotent stem cell reprogramming. Through this integration, the model can help identify optimal regenerative interventions tailored to patient-specific molecular and cellular profiles.

RESEARCH OBJECTIVES

The project aims to develop a scalable, interpretable, and biologically grounded computational model through two synergistic objectives:

1) Modeling Single-Cell Perturbations.

Design a multistage reinforcement learning algorithm that integrates scRNA-seq, CITE-seq, and scATAC-seq. The method combines **Trust Region Policy Optimization (TRPO)** for safe, curvature-aware exploration with Proximal Policy Optimization (PPO) for efficient fine-tuning. This integration mitigates convergence to local optima and captures true differentiation trajectories in the Waddington landscape.

2) Predicting Class II Antigen recognition and Reponse.

Develop a transformer-based neural framework for predicting peptide—Major Histocompatibility Complex (MHC-II) interactions, which underpin immune activation. A bidirectional encoder with attention-based cross-modality learning and RL-based feedback will iteratively refine predictions according to biological relevance, enhancing interpretability for vaccine and immunotherapy applications.

^{*}Francis Boabang is with the Concordia Institute for Information and Systems Engineering (CIISE), Concordia University, Montréal, QC, Canada., boabangf@yahoo.com)

METHODOLOGY

Overview

The unified framework consists of two interacting modules:

- Graph Attention Model with Reinforcement Learning Perturbation Module: This module leverages advanced reinforcement learning (RL) algorithms to model the dynamic and nonlinear behavior of single cells under genetic, chemical, or environmental perturbations. By framing cell state transitions as sequential decision-making processes, the RL policy learns optimal perturbation strategies that can predict downstream effects on gene expression and cell fate. The module is capable of capturing complex interactions between multiple genes and regulatory pathways, enabling more accurate simulations of cellular responses and the identification of critical nodes for targeted interventions.
- Graph Attention Model with Reinforcement Learning Antigen Recognition and Response Module: This module employs transformer architectures to model the sequence-to-structure relationships in peptide–MHC-II binding. By integrating both primary peptide sequences and structural features of the MHC-II molecules, the transformer can effectively predict binding affinities and presentation likelihood. Importantly, the module provides feedback to the RL-based perturbation module, allowing the learning policy to consider immune recognition outcomes when simulating perturbation effects. This integration ensures that the framework not only predicts cellular responses but also evaluates immunogenic consequences, which is critical for applications in immunotherapy and vaccine design.

Multi-Stage TRPO_PPO Reinforcement Learning

The TRPO_PPO algorithm [4] is designed to escape local minima in nonconvex biological landscapes. Let $\pi_{\theta}(a|s)$ denote the stochastic policy for selecting perturbation actions a given cell states s, with parameters θ .

The objective is to maximize expected reward:

$$J(\theta) = \mathbb{E}_{\pi_{\theta}} \left[\sum_{t=0}^{T} \gamma^{t} r_{t} \right], \tag{1}$$

where r_t measures similarity between predicted and observed post-perturbation profiles.

The natural gradient step is computed using Fisher-vector products and conjugate gradient:

$$Hv = \nabla_{\theta} \left(\nabla_{\theta} \bar{D}_{KL}(\pi_{\theta_{old}} || \pi_{\theta}) \cdot v \right) + \lambda_{damp} v, \tag{2}$$

and scaled to satisfy a KL-divergence trust region constraint δ :

$$\Delta\theta_{nat} = \alpha d, \quad \alpha = \sqrt{\frac{2\delta}{g^T d + \epsilon}}.$$
 (3)

Starting from $\theta' = \theta_{old} + \Delta \theta_{nat}$, the algorithm switches to PPO-style clipped fine-tuning:

$$L_{CLIP}(\theta) = \mathbb{E}_t \left[\min(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t) \right], \tag{4}$$

where \hat{A}_t is the advantage estimate and r_t the importance ratio.

This two-stage pipeline balances **exploration** (TRPO) and **sample-efficient exploitation** (PPO), enabling robust convergence in rugged loss surfaces typical of cellular differentiation.

Algorithmic Framework

Algorithm: Multi-Stage TRPO-PPO for Single-Cell Perturbation Prediction

- 1. Initialize policy π_{θ} , value network V_{ϕ} , and trust region δ .
- 2. Collect on-policy trajectories $\{s_t, a_t, r_t\}$.
- 3. Compute advantages \hat{A}_t and returns \hat{R}_t .
- 4. Perform TRPO update using conjugate gradient to compute $\Delta\theta_{nat}$.
- 5. Update parameters $\theta' = \theta + \Delta \theta_{nat}$.
- 6. Fine-tune with PPO for N_{ppo} epochs using clipped objective.
- 7. Repeat until convergence or early stopping criterion satisfied.

3.4 Biological Data Integration

Multiomic data from scRNA-seq, scATAC-seq, and CITE-seq will be integrated using contrastive learning (Seurat v5, Harmony). The RL state vector encodes latent gene expression, chromatin accessibility, and surface protein abundance, providing a unified representation for policy learning.

Antigen Presentation Module

We design a transformer encoder-decoder for peptide-MHC-II binding:

- Encoder: Embeds HLA and peptide sequences using self-attention.
- Decoder: Predicts binding affinity via cross-attention with structural embeddings.
- Feedback Loop: RL-based critic adjusts weights using immunogenicity metrics.

This system iteratively refines peptide–MHC predictions, improving interpretability for vaccine development.

EVALUATION FRAMEWORK

Evaluation metrics include both classification and regression criteria:

- Classification: Accuracy, Precision, Recall, F1, and AUPRC.
- **Regression:** MSE, RMSE, MAE, R^2 , and Pearson correlation.

Preliminary Experimental Results Summary

Empirical studies on scRNA-seq, ATAC, and CITE-seq datasets demonstrate substantial improvement of the TRPO_PPO model over PPO or TRPO alone.

Table I Performance comparison of different algorithms on testing Joint dataset(RNA and ATAC single cell data) [4].

Algorithm	Accuracy	Precision	Recall	F1	AUPRC	MSE	RMSE	MAE	R^2	Pearson Corr
PPO	0.6281	0.6577	0.6281	0.6283	0.6095	0.6044	0.7755	0.7004	0.1741	0.6388
TRPO_PPO	0.6661	0.6940	0.6661	0.6675	0.6428	0.5948	0.7672	0.6935	0.1948	0.6495

The multistage framework achieves high correlation with the experimental ground truth while reducing MSE by an order of magnitude compared to baselines, indicating improved landscape navigation and generalization.

THEORETICAL INTERPRETATION OF WADDINGTON LANDSCAPE

Cell reprogramming can be viewed as an optimization trajectory escaping a local optimum, driven by external perturbations (i.e., chemical, genetic, or environmental) that reshape the underlying potential landscape of gene regulation. Within the framework of the Waddington epigenetic landscape, each cell fate represents a local attractor basin stabilized by transcriptional and epigenetic feedback loops. Transitions between these basins require mechanisms that promote exploration and controlled deviation from stability, analogous to optimization strategies in machine learning that prevent premature convergence.

EXPECTED OUTCOMES

- 1) A unified framework that predicts both cellular perturbations and immune responses.
- 2) Mechanical interpretability linking RL states to biological processes.
- 3) Generalization to unseen perturbations and cell types.
- 4) Foundation for personalized regenerative and immunotherapeutic design.

COMPUTATIONAL AND ECONOMIC CONSIDERATIONS

The proposed framework relies on large-scale machine learning models, reinforcement learning algorithms, and multi-agent simulations, which require substantial computational resources. Each experimental run involves multiple API calls to large language models and transformer-based architectures. Based on preliminary benchmarking, the average cost per experiment ranges from 0.38to18.90, depending on the LLM backend and task complexity. Multi-agent simulations can generate up to 400,000 output tokens per task, further increasing compute time and latency.

We initially applied for Google Colab GCP credits to offset computational costs but were unsuccessful in securing this support. As a result, all planned experiments must rely on alternative cloud compute services or dedicated GPU/TPU resources, which increases the financial burden. For 100–200 experiments per month, we estimate a monthly cost of 500–4,000, and a cumulative annual cost of 6,000–48,000, depending on model scale, hyperparameter sweeps, and multi-agent simulations.

Funding support is therefore essential to ensure uninterrupted experimentation, timely progress, and scalability of the framework, particularly for tasks requiring high token throughput and repeated model fine-tuning

Optimization efficiency from multistage RL reduces computational overhead by up to 40% relative to single-stage training.

BROADER IMPACT AND FUTURE WORK

The proposed framework advances explainable AI in precision medicine, enabling:

- Multi-agent RL for modeling intercellular communication;
- Integration with spatial transcriptomics for tissue-level inference
- multi-agent reinforcement learning systems with improved activation functions, optimizers, and loss function formulations to better escape local optima in the Waddington landscape.

This unified platform connects cellular and immune intelligence, bridging single-cell perturbation dynamics with therapeutic discovery.

REFERENCES

- M. Lotfollahi, F. A. Wolf, and F. J. Theis, "scgen predicts single-cell perturbation responses," *Nature Methods*, vol. 16, no. 8, pp. 715–721, 2019.
 L. Hetzel, S. Böhm, N. Kilbertus, S. Günnemann, M. Lotfollahi, and F. Theis, "Predicting cellular responses to novel drug perturbations at a single-cell resolution," *arXiv preprint*, 2022, arXiv:2204.13545.
 H. Cui, C. Wang, H. Maan, K. Pang, F. Luo, N. Duan, and B. Wang, "scgpt: Toward building a foundation model for single-cell multi-omics using a constant of a constant of the con
- generative ai," *Nature Methods*, 2024.

 [4] F. Boabang and S. A. Gyamerah, "Escaping local optima in the waddington landscape: A multi-stage trpo-ppo approach for single-cell perturbation analysis," 2025, preprint available at https://github.com/boabangf/GNN_RL_gene_trajectory_perturbation/blob/main/Multi-Step%20Differentiation%2025
 2025, preprint available at <a href="https://github.com/boabangf/GNN_RL_gene_trajectory_perturbation_formalized-trajectory_perturbation_formalized-trajectory_perturbation_formalized-trajectory_perturbation_formalized-trajectory_perturbation_formalized-trajectory_perturbation_formalized-trajectory_perturbation_formalized-trajectory_perturbation_formalized-trajectory_perturbation_formalized-trajectory_perturbation_formalized-trajectory_perturbation_formalized-trajectory_perturbation_formalized-trajectory_perturbation_formalized-trajectory_perturbation_ pdf.