# Paper Critique: scGPT: toward building a foundation model for single-cell multi-omics using generative AI

Dan Peng

Department of Computer Science

June 20, 2024

danpeng@unc.edu

## 1. Research Problem

### 1.1. What research problem does the paper address?

This paper introduces scGPT, a foundation model for single-cell biology based on a generative pretrained transformer architecture. It tackles the challenge of effectively leveraging the rapidly expanding volume and diversity of single-cell sequencing data, which has grown to tens of millions of cells across various tissues and conditions.

### 1.2. What is the motivation of the research work?

The explosion of single-cell sequencing data has created a treasure trove of cellular information, yet current computational methods are fragmented into task-specific models with limited datasets. The authors envision a unifying approach—akin to how GPT models revolutionized NLP—that can learn universal representations across diverse cellular contexts. The core motivation is to create a "cellular language model" that can be pretrained once on massive datasets and then fine-tuned for specific downstream tasks, overcoming the scattered nature of current approaches.

## 2. Technical Novelty

### 2.1. What are the key technical challenges identified by the authors?

The authors identify several key challenges:

- Adapting transformer architectures to non-sequential gene expression data, unlike the sequential nature of text in NLP

- Designing an attention masking mechanism for generative modeling of gene expression

- Efficiently handling the high dimensionality of gene-level data (tens of thousands of genes)

- Creating a unified framework that enables transfer learning across diverse biological tasks

- Developing a model that can simultaneously learn both gene and cell embeddings

### 2.2. How significant is the technical contribution of the paper?

The paper makes significant technical contributions, representing a paradigm shift in computational approaches to single-cell biology. Rather than creating separate models for different tasks, scGPT offers a unifying "pretraining-then-fine-tuning" framework that works across diverse tasks. The novel attention masking mechanism for handling non-sequential omics data and the joint modeling of gene-gene and cell-cell relationships are particularly innovative. This approach connects the dots between disparate computational biology methods and establishes a new foundation for the field.

### 2.3. Identify 1-5 main strengths of the proposed approach.

- The model demonstrates remarkable versatility, achieving state-of-the-art performance across five distinct downstream tasks after fine-tuning

- The foundation model shows strong scaling effects, with performance improving as pretraining data size increases from 30,000 to 33 million cells

- scGPT learns biologically meaningful gene embeddings that capture functional relationships, as demonstrated by its ability to separate HLA class I and II genes without explicit supervision

- The attention-based gene network inference provides interpretable insights into gene-gene relationships and cellular regulatory mechanisms

- The approach excels at cross-condition generalization, such as predicting perturbation responses for unseen genetic interventions

## 2.4. Identify 1-5 main weaknesses of the proposed approach.

- The model's current design doesn't inherently address batch effects, requiring additional fine-tuning objectives for integration tasks

- Limited evaluation of the model's ability to generalize beyond human cells to other organisms

- The paper doesn't systematically assess how different model architectures and hyperparameters affect performance

- There's no thorough ablation study showing the individual contribution of each component in the pretraining pipeline

- The computational resources required for pretraining (33M cells) may be prohibitive for many research groups, raising accessibility concerns

## 3. Empirical Results

### 3.1. Identify 1-5 key experimental results, and explain what they signify.

- scGPT significantly outperforms specialized models in cell type annotation, achieving better precision and recall across multiple datasets. This demonstrates the value of transfer learning from large-scale pretraining in cellular biology.

- The model accurately predicts gene expression changes in response to unseen genetic perturbations, outperforming other methods by 5-20%. This showcases the model's capacity to capture complex gene regulatory relationships that generalize beyond its training data.

- When fine-tuned for multi-omic integration, scGPT effectively merges RNA, ATAC-seq, and protein data, revealing more distinct cell populations than alternative methods. This highlights the model's ability to learn complementary information across different molecular modalities.

- The pretrained gene embeddings capture biological knowledge without explicit supervision, shown by the spontaneous clustering of functionally related genes (like HLA class I vs. class II). This suggests the model has learned meaningful "semantics" of cellular biology.

- The model exhibits strong scaling properties, with downstream task performance consistently improving as pretraining data size increases. This mirrors findings in language models and suggests further gains as more cellular data becomes available.

### 3.2. Are there any weaknesses in the experimental section?

Yes, there are several weaknesses in the experimental design:

First, the paper lacks a thorough ablation study of the pretraining components. The authors don't systematically evaluate how different attention masking strategies, input representations, or training objectives contribute to the final performance, making it difficult to determine which innovations are most critical.

Second, the comparisons with baseline methods sometimes use different preprocessing pipelines, potentially confounding the source of performance improvements. A more rigorous approach would standardize preprocessing steps across all compared methods.

Third, the paper doesn't explore failure cases or limitations in depth. Understanding when and why the model underperforms would provide valuable insights for future improvements and appropriate application contexts.

Finally, while the authors show impressive results on human data, there's limited evaluation on cross-species generalization. This leaves open questions about the model's applicability to model organisms like mice, which are widely used in biomedical research.

## 4. Summary

scGPT represents a breakthrough in computational biology by successfully adapting the foundation model paradigm to single-cell genomics. Like GPT did for text, scGPT demonstrates that pretraining on massive cellular datasets can yield representations that transfer remarkably well across diverse downstream tasks. The model's ability to learn the "language of cells" without explicit supervision is particularly intriguing, as shown by how it naturally captures biological relationships between genes.

I am impressed by this model because it elegantly solves multiple challenging problems in single-cell analysis while offering an interpretable window into cellular biology. However, I still hold skepticism about its accessibility, generalizability beyond human data, and the limited ablation studies.

## 5. QA Prompt for a Paper Discussion

### 5.1. Discussion Question

How might foundation models like scGPT change the landscape of drug discovery and personalized medicine in

the coming years?

## 5.2. Your Answer

Foundation models like scGPT could revolutionize drug discovery and personalized medicine by acting as biological "crystal balls" that predict cellular responses to novel compounds or genetic interventions.

Just as ChatGPT can extrapolate from text it's seen to generate new content, scGPT can "imagine" how cells would react to treatments it's never encountered. This could dramatically accelerate the typically decade-long drug development process by rapidly narrowing candidate molecules and predicting off-target effects before a single wet lab experiment.

For personalized medicine, these models could serve as digital twins of a patient's cellular ecosystem, forecasting individual responses to therapies. Rather than the current trial-and-error approach, doctors could virtually test treatments on the patient's digital cellular representation first, identifying the path of least resistance to healing.

The most transformative aspect may be how these models democratize expertise. Just as anyone can now use AI to write code, soon biologists without computational backgrounds could interrogate complex cellular mechanisms through natural language interfaces built on models like scGPT.