



Impacts of Dataset Imbalance on Single-Cell Foundation Models

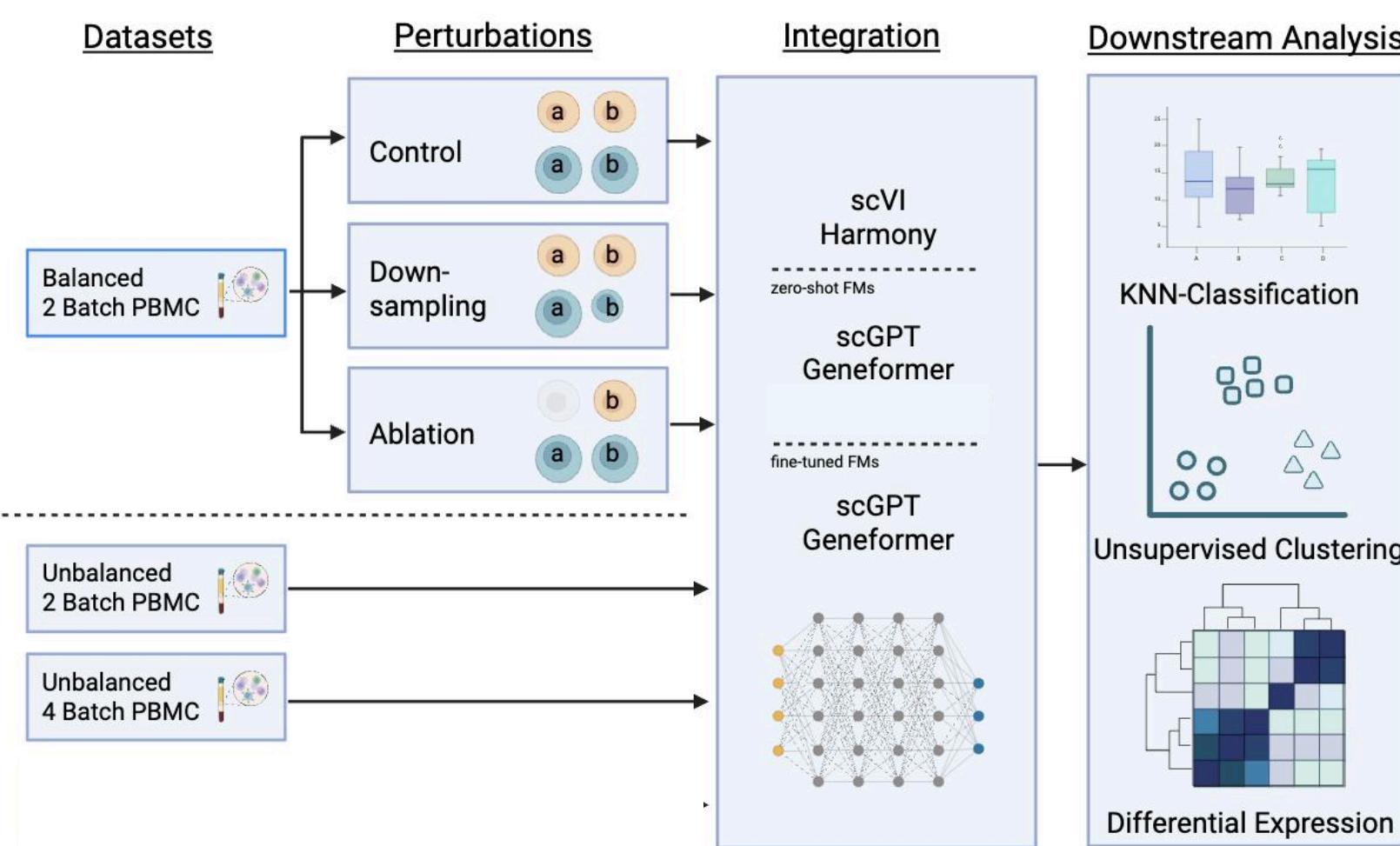
Izumi Ando^{1,2}, Hassaan Maan^{* 3 4 5}, Kieran R. Campbell^{* 1 2 4 6 7 8}

^{*}Jointly supervised this work, **1:** Lunenfeld-Tanenbaum Research Institute at Sinai Health, **2:** Department of Computer Science at the University of Toronto, **3:** Peter Munk Cardiac Centre at University Health Network, **4:** Vector Institute, **5:** Department of Medical Biophysics at the University of Toronto, **6:** Department of Molecular Genetics at the University of Toronto, **7:** Department of Statistical Sciences at the University of Toronto, **8:** Ontario Institute for Cancer Research

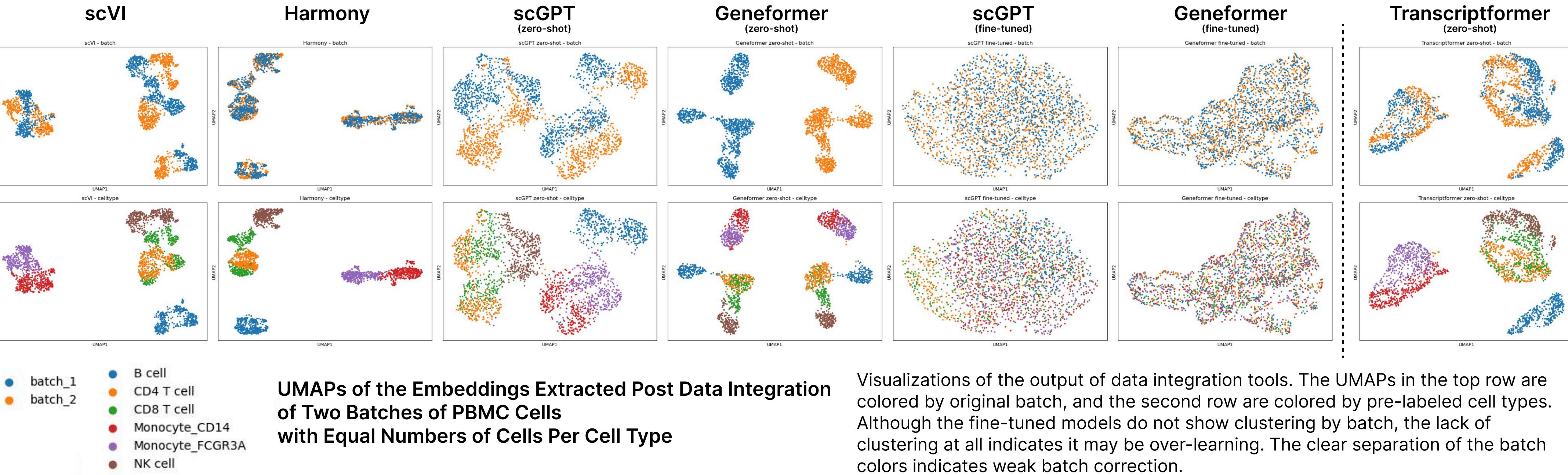
Pipeline

The goal of this project is to evaluate how novel single-cell foundation models perform compared to simpler models in handling imbalances in cell type proportions during data integration. Further, we examine the differences in biological signal that arise when imbalanced datasets are integrated using different tools.

To do so, we used the pipeline illustrated below. We employed perturbation experiments to study the impact of dataset imbalance and used inherently imbalanced datasets to study the potential differences in biological signal that arise from different integration methods. Foundation models were tested in both fine-tuned and zero-shot (non-fine-tuned) settings.

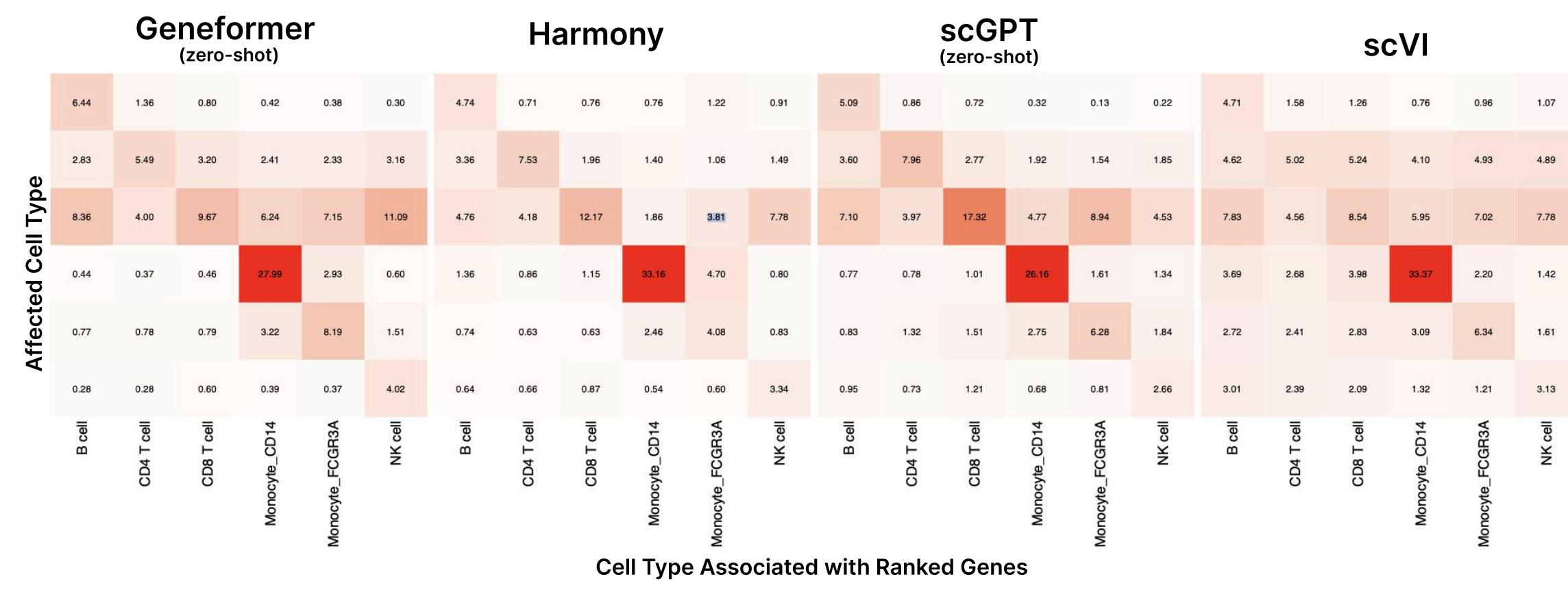


Visualizations of Integrated Embeddings



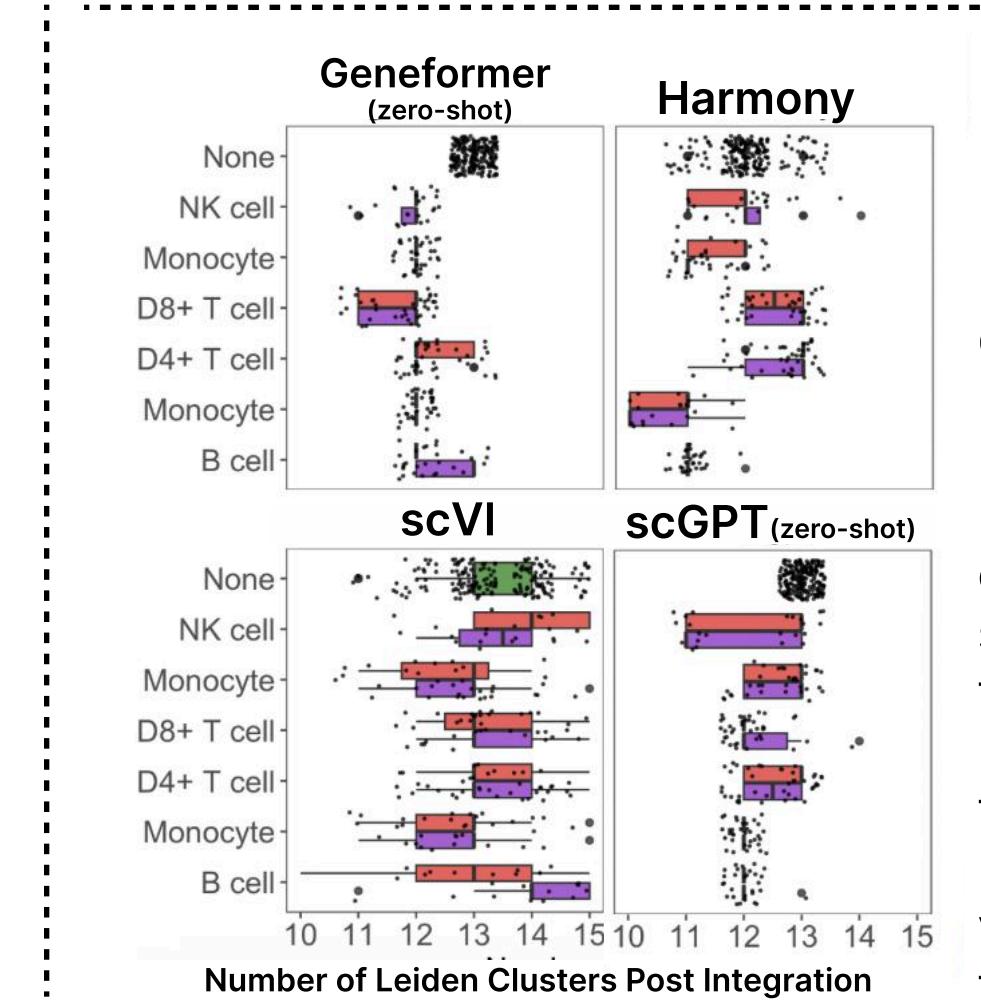
Visualizations of the output of data integration tools. The UMAPs in the top row are colored by original batch, and the second row are colored by pre-labeled cell types. Although the fine-tuned models do not show clustering by batch, the lack of clustering at all indicates it may be over-learning. The clear separation of the batch colors indicates weak batch correction.

Preliminary Results : How does imbalance in cell type proportions affect integration performance and subsequent analysis?



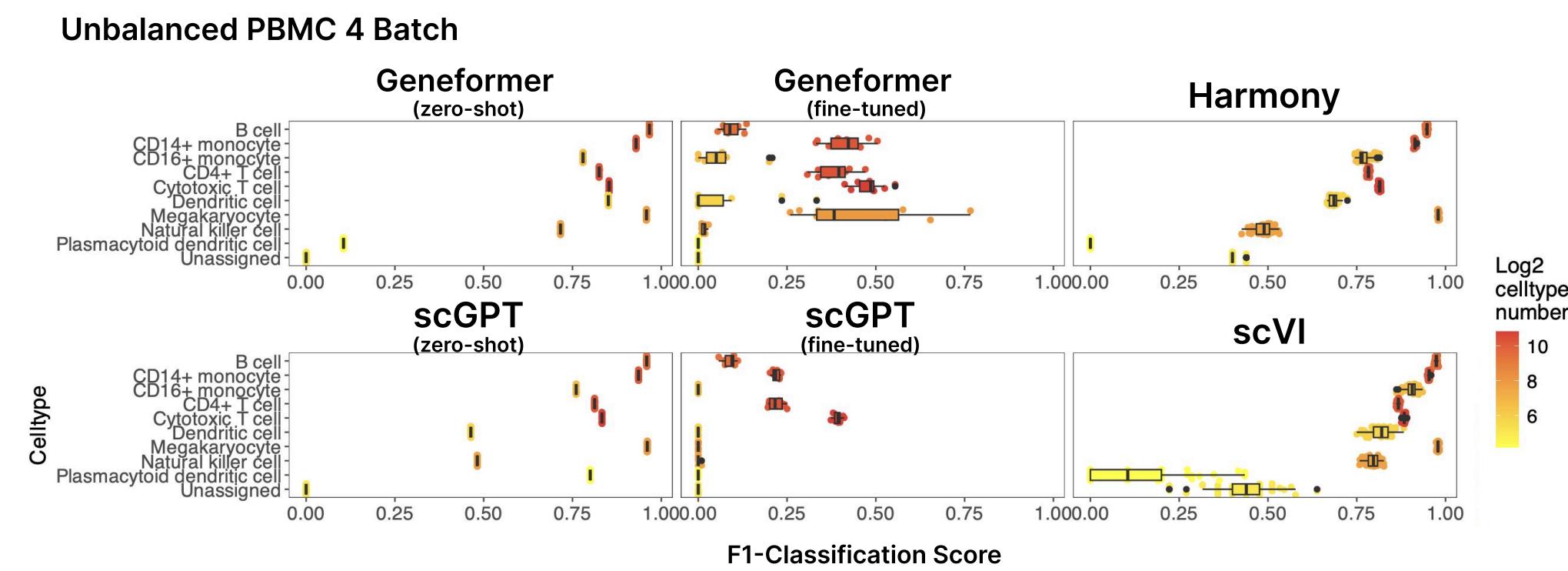
Differential Expression Analysis: Variance of Gene Rankings

It is important to understand whether imbalance in cell type proportions in integrated data sets affects biological signal in differential expression analysis. The top 10 marker genes for each cell type in the datasets being integrated were selected and the variance of those genes' rankings were computed after integration. This heatmap shows that when CD14 monocytes are downsampled or ablated, their ranking can significantly vary regardless of the integration method.



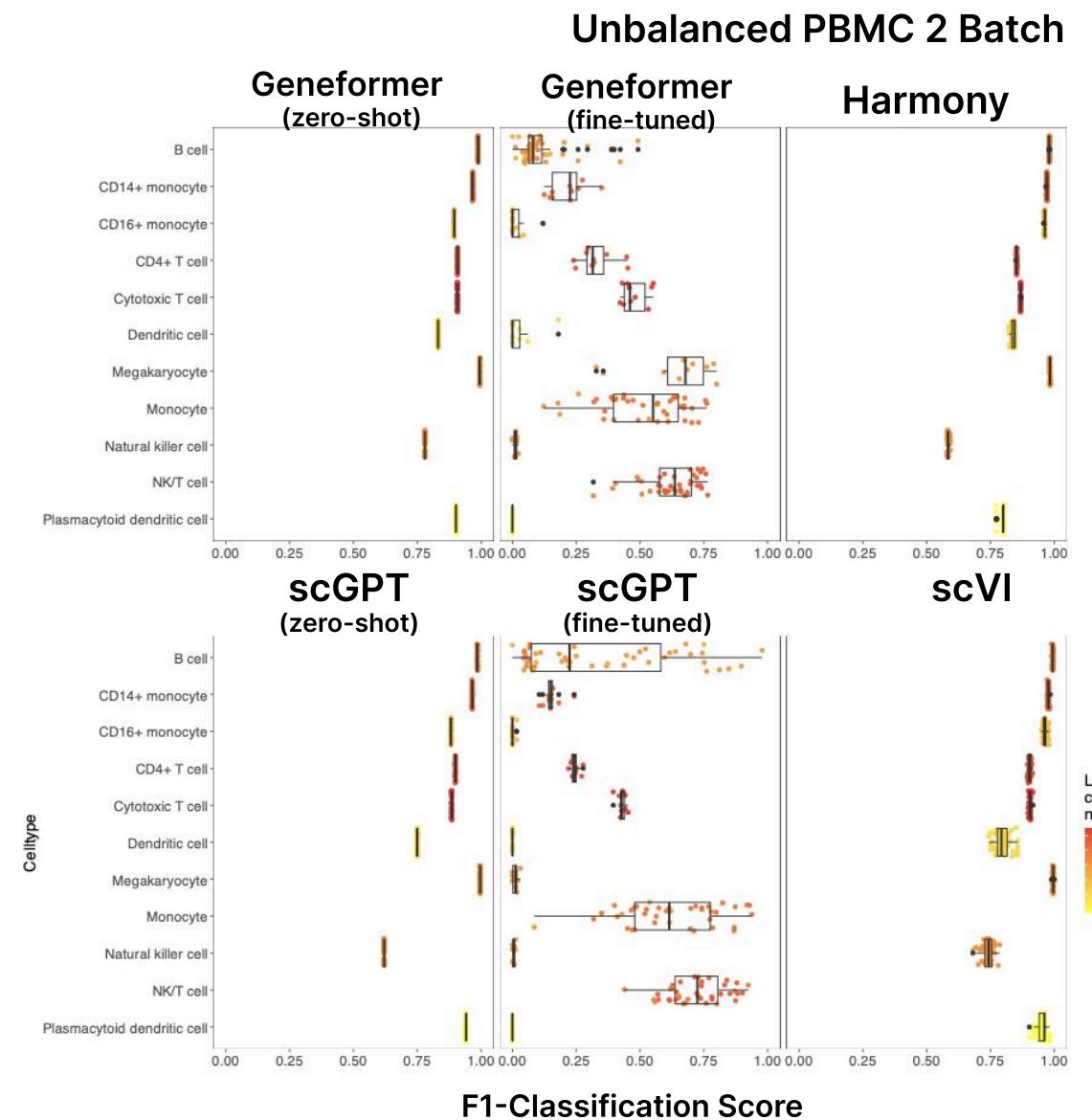
Unsupervised Clustering: Changes in Clustering When Proportions of Cell Types Shift

Unsupervised clustering is a common downstream task in scRNA-seq analysis. For most cell types, downsampling and ablation (in other words, an imbalance in cell type proportions) leads to varying results. This could potentially lead to varying biological interpretations of the data.



KNN Classification: Integration Performance

F1 scores from KNN Classification indicate integration performance. The two sets of boxplots above are from repeated integration attempts of the 2-batch unbalanced PBMC dataset and the 4-batch unbalanced PBMC dataset. The fine-tuned models show poor performance, while the other models display similar performance. In the 2-batch PBMC datasets, all models perform worse for classifying natural killer cells, while in the 4-batch PBMC dataset, plasmacytoid dendritic cells also show low scores agnostic to the integration method.



Takeaways

- Imbalance in the number of cells per cell type in a sc-RNA-seq dataset DOES affect the performance of data integration across all methods.
- Certain cell types affect integration performance, agnostic to integration method.
- Challenges in benchmarking sc-RNA-seq foundation models include the fact that most are not designed to be fine-tuned for data integration, and they each have very different interfaces.

Next Steps