# Kernel-based perturbation testing for single-cell data

Franck Picard

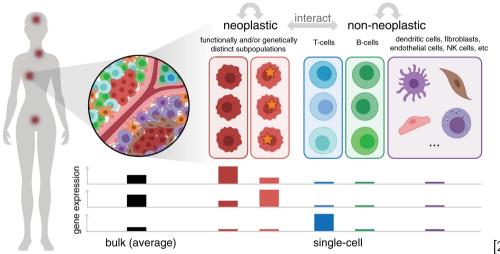
Laboratoire Biologie et Modélisation de la Cellule. CNRS ENS-Lyon



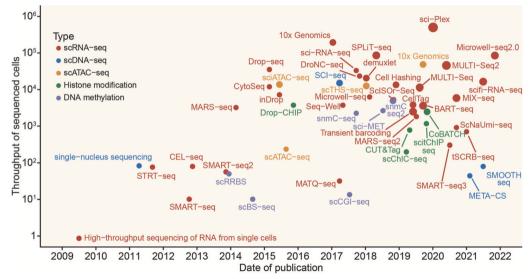
#### **Outline**

- 1. The Single-Cell Revolution
- 2. Comparison of Gene Expression Distributions
- 3. Introduction to kernel testing
- 4. Illustration
- 5. Towards perturbation analysis

#### From bulk to distributions of gene expression

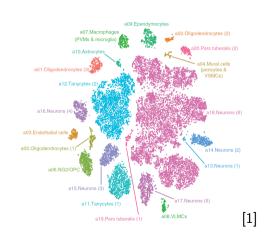


## A timeline: produced data

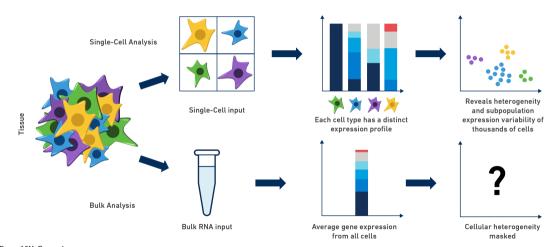


#### **Machine Learning challenges**

Dimension Reduction / Visualization
Clustering cell-type discovery
Datasets alignments
Catch cells-ecosystems behaviors
Simulation of fake data
Data integration
Genes expression comparison



#### Single-Cell from a statistician's perspective



From 10X Genomics

## **Comparing Biological Conditions**

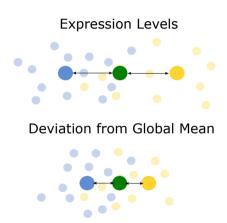
Gene-wise comparison

Statistical Testing

- → Score the difference
- $\rightarrow$  Control type-I errors

Single-cell data  $n \sim 10^6$ 

Try non-parametrics!

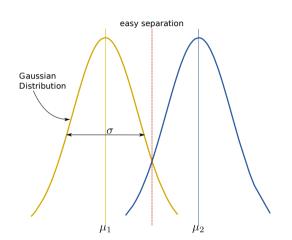


#### **Statistical Setting: two-sample test**

logFC are valid provided  $\mu$  and  $\sigma$  are good summaries of the information

Easy linear separation

Not adapted to single-cell assays



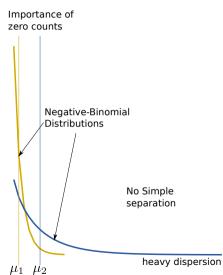
#### sc-RNAseq data are count data

Specificities: discrete, zeros

How to define the signal-to-noise ratio?

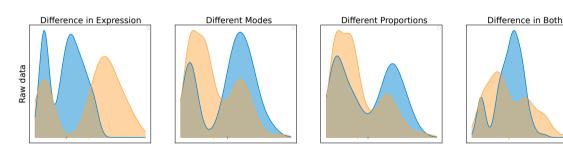
Standard: Negative Binomial distribution

Linear separation with GLM (parametric)



#### sc-RNASeq are complex count distributions

Compare Gene Expression distributions  $\mathbb{P}_1$  vs  $\mathbb{P}_2$ 



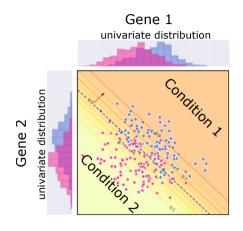
 $\rightarrow$  No simple linear separation

#### **Gene-Wise Strategy: a Good Option?**

Gene Expressions are highly dependent

Multivariate distributions

Calls for non-linear embedding



Joint distribution

[6]

#### **Statistical Challenge**

Statistical testing is based on what is expected under  $\mathcal{H}_0$ 

Control the random fluctuations of the embeddings under the null

Li et al. Genome Biology (2022) 23:79 https://doi.org/10.1186/s13059-022-02648-4 Genome Biology

#### **SHORT REPORT**

**Open Access** 

Exaggerated false positives by popular differential expression methods when analyzing human population samples

Yumei Li<sup>1†</sup>, Xinzhou Ge<sup>2†</sup>, Fanglue Peng<sup>3</sup>, Wei Li<sup>1\*</sup> and Jingyi Jessica Li<sup>2,4,5,6,7\*</sup> o

ightarrow Risk: detect a difference whereas the appropriate model there would not

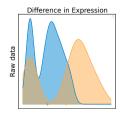
## Take-Home Message Slide (1)

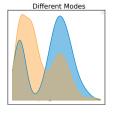
- ✓ Single-cell data are complex distributions
- ✓ the logFC may not be adapted to every situation
- √ pseudo-bulk approaches are possible (GLM)
- ✓ Only based on summary statistics
- √ A dedicated framework is required to perform differential analysis based on distribution

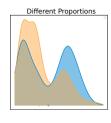
#### **Outline**

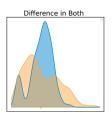
- 1. The Single-Cell Revolution
- 2. Comparison of Gene Expression Distributions
- 3. Introduction to kernel testing
- 4. Illustration
- 5. Towards perturbation analysis

#### **Comparing Gene Expression Distributions**







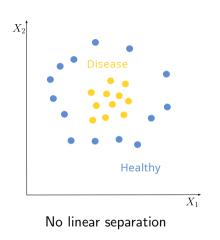


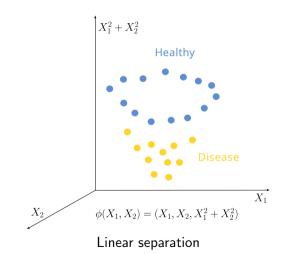
• Single-cell differential expression by distributions comparison :

$$\mathcal{H}_0:\left\{\mathbb{P}_1=\mathbb{P}_2\right\}$$

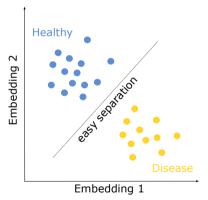
- No simple linear separation: SNR is not relevant anymore
- Idea: transform data into a new space
- Use SNR and linear separation on the transformed data

#### Data transformation for better separation

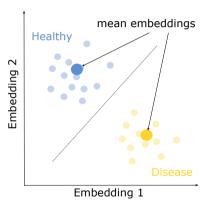




#### Rich Representations of complex data



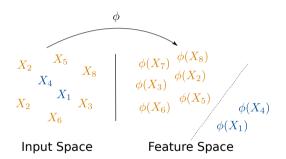
Work on joint transcriptomic embeddings



Mean embeddings by condition

## What is an embedding?

- Transform the input data  $X_i \to \phi(X_i)$
- New representation (UMAP, tSNE)
- Easy separation after transformation ?
- How to choose  $\phi$  ?



## Kernel Methods provide powerful embeddings

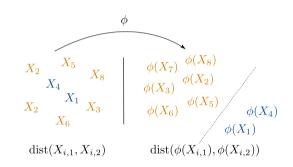
Similarity between data

$$\operatorname{dist}(X_{i,1},X_{i,2})$$

• Similarity between embeddings

$$\operatorname{dist}\!\left(\phi(X_{i,1}),\phi(X_{i,2})\right)$$

This is what does a kernel!



#### How to choose the kernel?

• Popular kernel : Gaussian kernel (h hyperparameter)

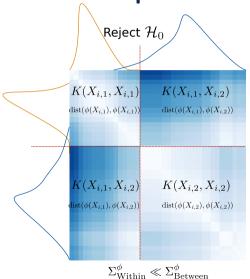
$$K(X_i, X_{i'}) \propto \exp \left\{-\frac{1}{2} \left(\frac{X_i - X_{i'}}{h}\right)^2\right\}$$

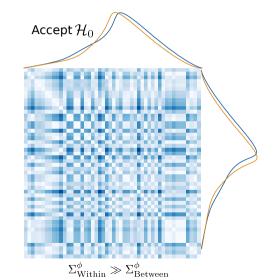
Theory ensures it is a distance between embeddings

$$K(X_i, X_{i'}) = dist(\phi(X_i), \phi(X_i))$$

ullet The embedding  $\phi$  exists but does not need to be defined

#### Kernels to compare distributions





## Take-Home Message Slide (2)

- √ Standard Differential Expression procedures can be applied by averaging data (pseudo bulk)
- √ Propose tests based on distributions comparisons
- √ Work on the embedding of distributions using a kernel
- √ Describe the distributions by the mean and the covariance of the embeddings

#### **Outline**

- 1. The Single-Cell Revolution
- 2. Comparison of Gene Expression Distributions
- 3. Introduction to kernel testing
- 4. Illustration
- 5. Towards perturbation analysis

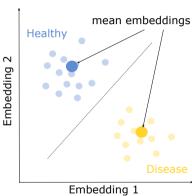
## **Embedding distributions**

• Mean Embedding of  $\mathbb{P}$ :

$$\mu_{\mathbb{P}} = \mathbb{E}_{X \sim \mathbb{P}} \Big( \phi(X) \Big)$$

• Covariance of the embeddings under  $\mathbb{P}$ :

$$\Sigma_{\mathbb{P}} = \mathbb{E}_{X \sim \mathbb{P}} \Big[ (\phi(X) - \mu_{\mathbb{P}})^{\otimes 2} \Big]$$



Mean embeddings by condition

#### Metric between distributions

• Testing  $H_0$  requires a metric between distributions

$$\mathcal{H}_0:\left\{\mathbb{P}_1=\mathbb{P}_2
ight\}$$

Expected property of the metric

$$\mathbb{P}_1 = \mathbb{P}_2 \quad \Leftrightarrow \quad \mu_{\mathbb{P}_1} = \mu_{\mathbb{P}_2}.$$

• The Maximal Mean Discrepancy:

$$\mathsf{MMD}^2(\mathbb{P}_1, \mathbb{P}_2) = \|\mu_1 - \mu_2\|_{\mathcal{H}}^2$$

#### Computing the empirical MMD

Empirical mean embeddings and MMD

$$\widehat{\mu}_1 = \frac{1}{n_1} \sum_{i=1}^{n_1} \phi(X_{i,1}) \quad \widehat{\mu}_2 = \frac{1}{n_2} \sum_{i=1}^{n_2} \phi(X_{i,2})$$

$$\widehat{\mathsf{MMD}}^{2} = \|\widehat{\mu}_{2} - \widehat{\mu}_{1}\|_{\mathcal{H}}^{2} \\
= \frac{1}{n_{1}(n_{1}-1)} \sum_{i \neq i'} k(X_{i,1}, X_{i',1}) + \frac{1}{n_{2}(n_{2}-1)} \sum_{i \neq i'} \sum_{i=1}^{n_{2}} \sum_{i'=1}^{n_{2}} k(X_{i,2}, X_{i',2}) \\
- \frac{2}{n_{1}n_{2}} \sum_{i,i'} k(X_{i,1}, X_{i',2})$$

The MMD is a testing framework based on kernelized distances

#### Between/Within kernel trade-off

Intra-condition distances

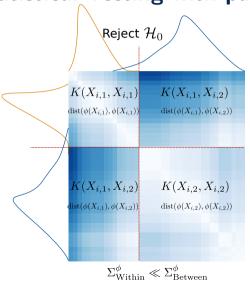
$$\frac{1}{n_1^2} \sum_{i=1}^{n_1} \sum_{i'=1}^{n_1} k(X_{i,1}, X_{i',1}) \quad \text{and} \quad \frac{1}{n_2^2} \sum_{i=1}^{n_2} \sum_{i'=1}^{n_2} k(X_{i,2}, X_{i',2})$$

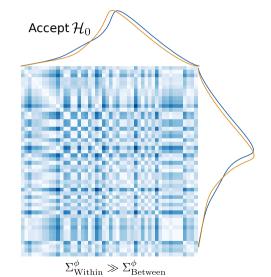
- → If small, conditions are homogeneous
- Inter-condition distance

$$\frac{1}{n_1} \frac{1}{n_2} \sum_{i=1}^{n_1} \sum_{i'=1}^{n_2} k(X_{i,1}, X_{i',2})$$

 $\rightarrow$  If high, conditions are well separated

Statistical Testing with pair-wise distances





#### **Considering Variances**

Separated Conditions:

$$\Sigma_{Within} \ll \Sigma_{Between}$$

Similar conditions :

$$\Sigma_{Within} \sim \Sigma_{Between}$$

• Construct the discriminant ratio between conditions:

$$\mathsf{D}^2(\mathbb{P}_1,\mathbb{P}_2) = \Sigma_{\mathsf{Within}}^{-1} \Sigma_{\mathsf{Between}}$$

• Similar to a 1-way-ANOVA:

$$\phi(X) \sim ext{Condition}$$

## **Definition of Intra/Inter Variance of embeddings**

• The MMD is linked to the between-group covariance

$$\widehat{\Sigma}_B = \frac{n_1 n_2}{n^2} \Big( \widehat{\mu}_2 - \widehat{\mu}_1 \Big)^{\otimes 2}$$

• Define the within-group covariances  $\widehat{\Sigma}_1$  and  $\widehat{\Sigma}_2$ 

$$\widehat{\Sigma}_1 = \frac{1}{n_1} \sum_{i=1}^{n_1} \left( \phi(X_{1,i}) - \widehat{\mu}_1 \right)^{\otimes 2}, \quad \widehat{\Sigma}_2 = \frac{1}{n_2} \sum_{i=1}^{n_2} \left( \phi(X_{2,i}) - \widehat{\mu}_2 \right)^{\otimes 2}$$

$$\Sigma_W = \frac{n_1}{n} \Sigma_1 + \frac{n_2}{n} \Sigma_2$$

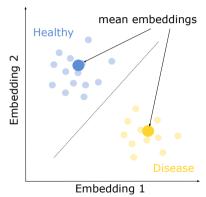
#### The Normalized MMD

The normalized MMD statistics is

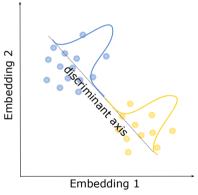
$$\mathsf{D}^2(\mathbb{P}_1, \mathbb{P}_2) = \frac{n_1 n_2}{n} \left\| \Sigma_W^{-\frac{1}{2}} (\mu_2 - \mu_1) \right\|_{\mathcal{H}}^2$$
$$\sim \frac{1}{n} \operatorname{Tr} \left( \Sigma_W^{-1} \Sigma_B \right)$$

- ullet It is a kernelized discriminant ratio with  $\chi^2$  distribution under  $\mathcal{H}_0$
- The method is based on Kernel Fisher Discriminant Analysis

#### From Separation to Discrimination

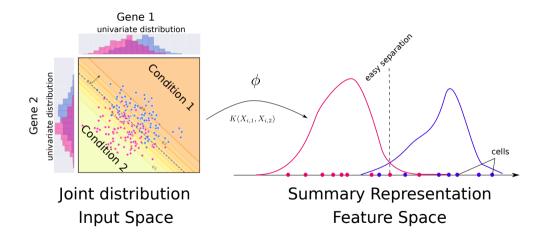


Mean embeddings by condition



Discrimination of cell populations

#### **Towards Classification**



## Take-Home Message Slide (3)

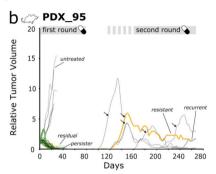
- ✓ Kernel methods can be used to define discrepancies between distributions
- ✓ Kernel tests are based on pair-wise distances between embeddings
- ✓ These distances can be normalized by embeddings variability
- √ pvalues can be obtained (approximations)
- √ The method is based on a classifier

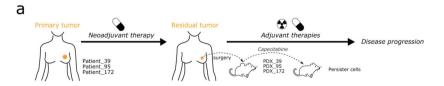
#### **Outline**

- 1. The Single-Cell Revolution
- 2. Comparison of Gene Expression Distributions
- 3. Introduction to kernel testing
- 4. Illustration
- 5. Towards perturbation analysis

#### ChemoResistance in Triple Negative Breast Cancer

- Emergence of resistant phenotypes is a multi-step process
- After drug insult only a pool of drug-tolerant persister cells manage to tolerate the treatment and survive.
- Reservoir from which drug-resistant cells can ultimately emerge.

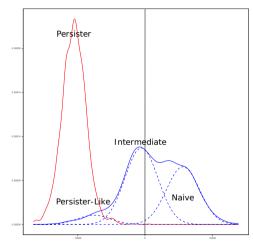




[5]

#### Kernel testing on Persister vs. Naive cells

- Persister cells survived the first treatment
- Reservoir for resistant cells
- Epigenomic data: 6376 features
- Compare untreated ( $\sim$  3000 cells) vs. persister ( $\sim$  2000 cells)
- Did we identify the reservoir of persister cells based on their epigenomic signatures
   ?



Summary of Whole Epigenome differences

#### METHOD Open Access

# Kernel-based testing for single-cell differential analysis



A. Ozier-Lafontaine<sup>1\*</sup>, C. Fourneaux<sup>2</sup>, G. Durif<sup>2</sup>, P. Arsenteva<sup>1</sup>, C. Vallot<sup>3,4</sup>, O. Gandrillon<sup>2</sup>, S. Gonin-Giraud<sup>2</sup>, B. Michel<sup>1\*†</sup> and F. Picard<sup>2\*†</sup>

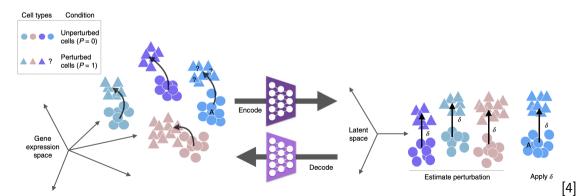
https://github.com/LMJL-Alea/ktest

#### **Outline**

- 1. The Single-Cell Revolution
- 2. Comparison of Gene Expression Distributions
- 3. Introduction to kernel testing
- 4. Illustration
- 5. Towards perturbation analysis

## **Context in Single-Cell Transcriptomics**

- Modeling/predicting the effects of perturbations is a key task in systems biology.
- Capture the heterogeneity of cell populations using Supervised Dimension Reduction
- Latent space captures the signal to predict a cell's response to perturbation.



#### BRIEF COMMUNICATION

https://doi.org/10.1038/s41587-020-0605-1



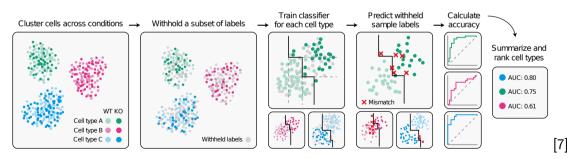


# Cell type prioritization in single-cell data

Michael A. Skinnider <sup>1,2,9</sup> <sup>1,2</sup>, Jordan W. Squair <sup>1,3,4,9</sup> <sup>1,3</sup>, Claudia Kathe <sup>1,3</sup>, Mark A. Anderson <sup>1,3</sup>, Matthieu Gautier <sup>1,3</sup>, Kaya J. E. Matson <sup>5</sup>, Marco Milano <sup>1,3</sup>, Thomas H. Hutson <sup>1,3</sup>, Quentin Barraud <sup>1,3</sup>, Aaron A. Phillips <sup>6</sup>, Leonard J. Foster <sup>2,7</sup>, Gioele La Manno <sup>1</sup>, Ariel J. Levine <sup>5</sup> and Grégoire Courtine <sup>1,3,8</sup> <sup>1,3,8</sup>

#### **New Paradigm**

- Identify cell-types more responsive to biological perturbations
- Hypothesis: responsive cell-types should be more separable
- Cell-types are prioritized based on the area under the curve AUC



# The two sides of Supervised Learning

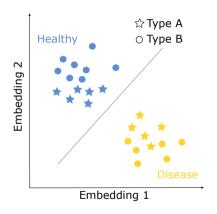
#### Classification

- Observe  $(y_1, x_1), \ldots, (y_n, x_n)$
- Construct a predictor  $f: \mathcal{X} \to \mathcal{Y}$
- Predict a new y
- $\mathcal{H}_1$  is not properly defined
- Pros: differences are not well defined
- Cons: x2 data to reach the same power

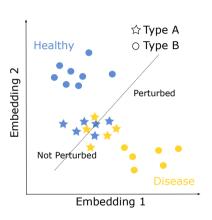
#### **Statistical Testing**

- Observe  $(y_1, x_1), \ldots, (y_n, x_n)$
- Define a test statistic
- Control Type-I error
- Ensure power
- Pros: powerful if  $\mathcal{H}_1$  can be defined
- Cons: when  $\mathcal{H}_1$  is ill-defined

## From Differential Analysis to Perturbation Analysis

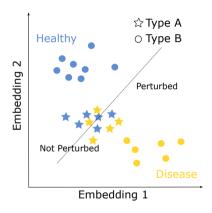


All Cell-types perturbed

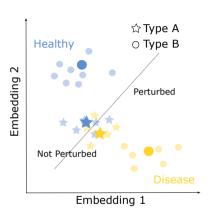


Differential Perturbation

# Perturbed Mean Embeddings



Differential Perturbation



Interaction Treatment × Cell-types

# **ANOVA** for non-linear Embeddings

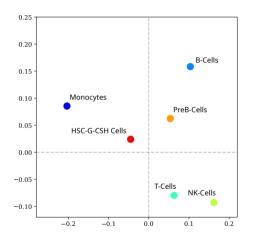
• Complex design : treatment, cell types factors

$$\phi(\mathsf{Expression}) \sim \mathsf{treatment} + \mathsf{celltype} + \mathsf{treatment} \times \mathsf{celltype}$$

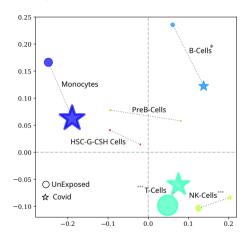
• Identify Perturbed cell types using contrasts

$$\mathcal{H}_0^{\bigstar}: \Big\{ \mathtt{Healthy} \times \bigstar = \mathtt{Disease} \times \bigstar \Big\}$$
 $\mathcal{H}_0^{\mathsf{O}}: \Big\{ \mathtt{Healthy} \times \mathsf{O} = \mathtt{Disease} \times \mathsf{O} \Big\}$ 

## Non-Linear perturbations following Covid Exposure

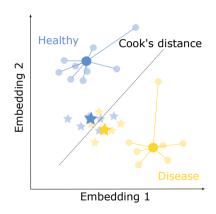


Cell-Type Effect\*\*\*



Interaction Cell-Type × Disease\*\*\*

#### Soon Included

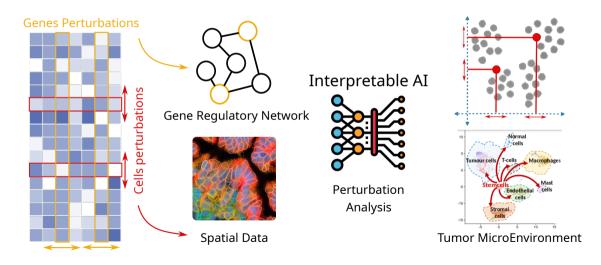


Atypical Cells identification



Multi-patients Designs

# **Perspectives**



## **Perspectives**

ANOVA approach:

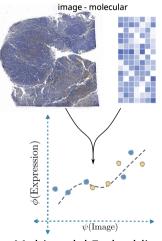
$$\phi$$
(Expression) =  $\mathbf{X}\boldsymbol{\beta} + \boldsymbol{E}$ 

- Test Contrasts:  $\mathbf{C}\beta = 0$
- Regression Approach

$$\phi(\mathsf{Expression}) = \psi(\mathsf{Covariates})\beta + E$$

- Test conditional independence:  $\beta = 0$
- Covariates can be space! (SPARK-X!)

#### Data integration



Multimodal Embeddings

## **Acknowledgments**

- Anthony Ozier-Lafontaine, Bertrand Michel, Perrine Lacroix, Nantes University
- Polina Arsenteva, Ghislain Durif, Lucy Attwood, ENS Lyon
- Vincent Rivoirard, Dauphine University
- Philippe Bertolino, CRCL, Lyon
- PEPR Digital Health (Al4scMed), ANR

#### References

- J. N. Campbell, E. Z. Macosko, H. Fenselau, T. H. Pers, A. Lyubetskaya, D. Tenen, M. Goldman, A. M. Verstegen, J. M. Resch, S. A. McCarroll, E. D. Rosen, B. Lowell, and L. T. Tsai. A molecular census of arcuate hypothalamus and median eminence cell types. *Nat. Neurosci.*, 20(3):484–496, Mar 2017.
- [2] J. Fan, K. Slowikowski, and F. Zhang. Single-cell transcriptomics in cancer: computational challenges and opportunities. Exp Mol Med, 52(9):1452–1465, Sep 2020.
- [3] Q. Jia, H. Chu, Z. Jin, H. Long, and B. Zhu. High-throughput single-ell sequencing in cancer research. Signal Transduction and Targeted Therapy, 7(1), May 2022.
- [4] M. Lotfollahi, F. A. Wolf, and F. J. Theis. scgen predicts single-cell perturbation responses. Nature Methods, 16(8):715-721, July 2019.
- [5] J. Marsolier, P. Prompsy, A. Durand, A.-M. Lyne, C. Landragin, A. Trouchet, S. T. Bento, A. Eisele, S. Foulon, L. Baudre, K. Grosselin, M. Bohec, S. Baulande, A. Dahmani, L. Sourd, E. Letouzé, A.-V. Salomon, E. Marangoni, L. Perié, and C. Vallot. H3k27me3 conditions chemotolerance in triple-negative breast cancer. *Nature Genetics*, 54(4):459–468, Apr. 2022.
- [6] V. Ntranos, L. Yi, P. Melsted, and L. Pachter. A discriminative learning approach to differential expression analysis for single-cell rna-seq. Nature Methods, 16(2):163–166, Jan. 2019.
- [7] M. A. Skinnider, J. W. Squair, C. Kathe, M. A. Anderson, M. Gautier, K. J. E. Matson, M. Milano, T. H. Hutson, Q. Barraud, A. A. Phillips, L. J. Foster, G. La Manno, A. J. Levine, and G. Courtine. Cell type prioritization in single-cell data. Nature Biotechnology, 39(1):30–34, July 2020.