Single cell RNA-seq analysis Part II: understanding cell types **BMI/CS 776** www.biostat.wisc.edu/bmi776/ Spring 2022 Daifeng Wang daifeng.wang@wisc.edu

Thanks to Ting Jin for slides!

These slides, excluding third-party material, are licensed under CC BY-NC 4.0 by Mark Craven, Colin Dewey, Anthony Gitter and Daifeng Wang

## Outline

- scRNA-seq data analysis
  - Cell type annotation
    - SingleR
  - Cell type markers identification
  - Pseudo timing
    - Monocle
  - Cell-type gene regulatory networks
    - SCENIC
    - BEELINE
  - Single cell deconvolution
    - CIBERSORTx

## Outline

- scRNA-seq data analysis
  - Cell type annotation
    - SingleR
  - Cell type markers identification
  - Pseudo timing
    - Monocle
  - Cell-type gene regulatory networks
    - SCENIC
    - BEELINE
  - Single cell deconvolution
    - CIBERSORTx

# Cell type annotation

- Cell types -> cellular functions
- Assign the cell type for each cell



https://bioconductor.org/books/release/OSCA/cell-type-annotation.html

# Cell type annotation tools

- Supervised methods: a training dataset • labeled with the corresponding cell population is needed to train the classifier
  - SingleR, ACTINN, CaSTle

**Prior-knowledge based methods:** either a marker gene file is required as an input or a pretrained classifier for specific cell populations is provided

SVM <sub>rejection</sub> 0.99         0.99         0.89         1         0.98         1         1         0.99         1         1         0.99         0.9         0.9           SCPred         1         0.98         0.98         0.97         1         0.99         1         0.97         0.97         0.98         0.97         1         1         0.94         1         0.99         0.88         0.87         0.99         1         1         0.94         1         0.99         0.88         0.87         0.99         1         1         0.94         1         0.99         0.88         0.87         0.99         1         1         0.94         1         0.99         0.88         0.87         0.88         0.87         0.88         0.88         0.87         0.88         0.87         0.88         0.88         0.88         1         1         0.99         0.88		Pancreas								ТМ	Allen Mouse Brain		PBMC		
ScPred         1         0.98         0.98         1         0.95         1         1         0.97         1         0.98         0.97         0.98		SVM <sub>rejection</sub>				1	0.98	1	1	0.99	1	1	0.98	0.99	0.92
SVM-       0.38       0.39       1       0.99       1       1       0.98       1       0.98       0.89       0.99       1       0.99       0.99       0.99       0.89       0.89       0.89       0.99		scPred-	1	0.98	0.98	1	0.95	1	1	0.97	1	1	0.69	0.96	
singleCellNet         0.97         0.96         0.97         0.99         1         1         1         0.94         1         0.99         0.87         0.88         0.77           ACTINN         0.97         0.98         0.97         1         0.95         1         0.97         1         0.99         0.80         0.88         0.77           CaSTL         0.98         0.98         0.96         0.98         0.80         1         0.97         1         0.98         0.80         0.88         1         0.98         1         0.99         0.99         0.97         0.73         0.88         1         0.99         0.99         0.99         0.99         0.99         1         1         1         1         0.98         0.88         0.88         0.88         0.88         0.88         0.88         0.88         0.89         0.89         0.89         0.89         0.88         <		SVM-	0.98	0.98	0.97	1	0.99	1	1	0.98	1	0.99	0.89	0.95	0.7
ACTINN-       0.97       0.98       0.97       1       0.95       1       1       0.97       1       0.98       0.88       0.88       0.77         CaSTLe       0.93       0.94       0.96       0.98       0.96       1       0.99       0.94       1       0.99       0.79       0.84       0.7         Scmapclub       0.98       0.97       1       0.73       0.84       0.7       0.84       0.7       0.84       0.7       0.84       0.7       0.84       0.7       0.84       0.7       0.84       0.7       0.84       0.7       0.84       0.7       0.84       0.7       0.84       0.7       0.84       0.85       1       1       0.87       1       0.99       0.88       0.88       0.88       0.8       0.88       0.88       0.85       1       1       0.88       1       0.99       0.88       0.88       0.85       0.81       1       0.91       0.99       0.88       0.88       0.85       0.81       1       1       0.88       0.89       0.86       0.85       0.81       0.81       0.97       0.81       0.81       0.81       0.81       0.81       0.81       0.81       0.81 <td></td> <td>singleCellNet-</td> <td>0.97</td> <td>0.96</td> <td>0.97</td> <td>0.99</td> <td>1</td> <td>1</td> <td>1</td> <td>0.94</td> <td>1</td> <td>0.99</td> <td>0.87</td> <td>0.88</td> <td>0.74</td>		singleCellNet-	0.97	0.96	0.97	0.99	1	1	1	0.94	1	0.99	0.87	0.88	0.74
CaSTLe       0.93       0.94       0.96       0.98       0.99       1       0.99       0.94       1       0.99       0.79       0.64       0.7         Scmapcell       0.88       0.97       0.90       0.99       0.99       1       1       0.98       1       0.99       0.84       0.73       0.68         LDA<		ACTINN-	0.97	0.98	0.97	1	0.95	1	1	0.97	1	0.99	0.86	0.88	0.74
scmapcell       0.98       0.97       1       0.73       1       1       0.98       1       1       0.91       0.73       0.8         LDA       0.94       0.97       0.96       0.99       0.89       1       1       0.95       1       0.95       0.88       0.88       0.88       0.85       1       1       0.95       1       0.96       0.88       0.86       1       1       0.95       0.88       0.88       0.85       1       1       0.95       0.88       0.86       0.85       1       1       0.86       1       0.99       0.86       0.88       0.85       1       1       0.88       1       0.99       0.86       0.86       0.86       0.85       1       1       0.88       1       0.99       0.86       0.86       0.86       0.85       1       1       0.88       0.81       0.86       0.86       0.86       0.89		CaSTLe-	0.93	0.94	0.96	0.98	0.96	1	0.99	0.94	1	0.99	0.79	0.84	0.79
LDA       0.94       0.97       0.96       0.99       0.89       1       1       0.95       1       0.99       0.88       0.63       0.63         Scmapcluster       0.99       0.95       0.97       1       1       1       0.87       1       0.98       0.88       0.73       0.41         KFF       0.94       0.96       0.98       0.88       0.85       1       1       0.91       1       0.99       0.73       0.81       0.81       0.81       0.91       0.99       0.88       0.83       0.81       0.91       0.91       0.99       0.83       0.81       0.83       0.91       0.81       0.91       0.99       0.31       0.81       0.93       0.91       0.99       0.99       0.91       0.99       0		scmapcell-	0.98	0.98	0.97	1	0.73	1	1	0.98	1	1	0.91	0.73	0.64
scmapcluster       0.99       0.95       0.97       1       1       1       1       0.87       1       0.98       0.88       0.73       0.44         RF       0.94       0.96       0.97       0.95       0.97       0.99       1       1       0.91       1.0       0.99       0.73       0.81       0.66       0.33         LAmbDA       0.92       0.97       0.95       0.97       1       1       0.82       1       0.99       0.84       0.99       0.99       1       1       0.88       1       0.99       0.84       0.99		LDA-	0.94	0.97	0.96	0.99	0.89	1	1	0.95	1	0.99	0.88	0.63	0.66
RF       0.94       0.94       0.96       0.98       0.85       1       1       0.91       1       0.99       0.73       0.81       0.85         SingleR       0.96       0.97       0.95       0.97       0.99       1       1       0.88       1       0.97       0.86       0.97       0.99       1       1       0.88       1       0.99       0.86       0.97       0.86       0.97       1       1       0.82       1       0.99       0.84       0.97       0.96       0.97       0.99       0.99       0.99       0.99       0.99       0.99       0.99       0.97       0.81       0.71       0.55         CHETAH       0.91       0.94       0.96       0.97       0.96       1       1       0.83       1       0.97       0.81       0.71       0.55         ScU1       0.95       0.56       0.97       0.96       1       1       0.90       1       0.97       0.68       0.71       0.55       0.55       0.55       0.55       0.55       0.55       0.55       0.55       0.55       0.55       0.55       0.55       0.55       0.55       0.55       0.55       0.55       0.55		scmapcluster-	0.99	0.95	0.97	1	1	1	1	0.87	1	0.98	0.88	0.73	0.44
SingleR       0.96       0.97       0.95       0.97       0.99       1       1       0.88       1       0.97       0.86       0.66       0.33         LAmbDA<		RF-	0.94	0.94	0.96	0.98	0.85	1	1	0.91	1	0.99	0.73	0.81	0.66
LAmbDA- 0.92 0.8 0.95 0.96 0.97 11 1 1 0.62 11 0.99 0.84 0.84 0.9 NMC- 0.92 0.91 0.84 0.93 0.99 0.92 0.9 0.69 0.99 0.99 0.99 0.99 0.99 0.99		SingleR-	0.96	0.97	0.95	0.97	0.99	1	1	0.88	1	0.97	0.86	0.66	0.32
NMC       0.92       0.91       0.84       0.93       0.99       0.90       0.97       0.81       0.71       0.81         CHETAH       0.91       0.94       0.96       0.97       0.96       1       1       0.83       1       0.96       0.97		LAmbDA-	0.92	0.8	0.95	0.96	0.97	1	1	0.62	1	0.99	0.84		0.4
CHETAH- 0.91       0.94       0.96       0.97       0.96       1       1       0.83       1       0.96       0.81       0.65       0.1         ScVI- 0.98       0.56       0.97       0.99       1       1       0       1       0.97       0.9       0.97       0.9       0.91       0.91       0.97 </td <td></td> <td>NMC-</td> <td>0.92</td> <td>0.91</td> <td>0.84</td> <td>0.93</td> <td>0.99</td> <td>0.92</td> <td>0.9</td> <td>0.69</td> <td>0.99</td> <td>0.97</td> <td>0.81</td> <td>0.71</td> <td>0.55</td>		NMC-	0.92	0.91	0.84	0.93	0.99	0.92	0.9	0.69	0.99	0.97	0.81	0.71	0.55
scVI-       0.98       0.56       0.97       0.99       1       1       0       1       0.97       0.97       0.67         scID-       0.75       0.59       0.99       0.85       0.88       1       1       0.42       11       0.99       0.63       0.61       0.4         Cell_BLAST-       0.11       0.89       0.79       0.08       0.63       1       0.99       0.97       1       0.99       0.63       0.61       0.4         kNN-       0.91       0.95       0.95       0.85       0.03       1       0.99       0.97       1       0.99       0.63       0.91       0.77         kNN-       0.91       0.95       0.95       0.85       0.03       1       0.99       0.97       1       0.99       0.76       0.91       0.77         kNN-       0.91       0.95       0.95       0.85       0.03       1       0.98       0.92       1       0.64       0.91       0.55         SCINA-       SCINA-       I       I       I       I       0.99       I       I       0.99       0.7       0.91       0.7       0.91       0.7       0.91       0.7		CHETAH-	0.91	0.94	0.96	0.97	0.96	1	1	0.83	1	0.96	0.81	0.65	0.11
sclD       0.75       0.59       0.95       0.85       0.8       1       1       0.42       1       0.95       0.63       0.4       0.4         Cell_BLAST       0.11       0.89       0.79       0.08       0.63       1       0.99       0.97       1       0.99       0.76       0.91       0.7         kNN       0.91       0.95       0.95       0.85       0.03       1       0.99       0.97       1       0.99       0.76       0.91       0.7         kNN       0.91       0.95       0.95       0.85       0.03       1       0.99       0.97       1       0.99       0.76       0.91       0.55         SCINA       0.91       0.95       0.95       0.85       0.03       1       0.98       0.92       1       0.64       0.4       0.45       0.5         SCINA       0.91       0.95       0.75       0.91       0.95       0.75       0.91       0.91       0.7       0.91       0.7       0.91       0.7       0.91       0.7       0.91       0.7       0.91       0.7       0.91       0.7       0.91       0.7       0.91       0.7       0.91       0.7       0.91		scVI-	0.98	0.56	0.97	0.99	1	1	1	0	1	0.97	0	0.97	0.64
Cell_BLAST       0.11       0.89       0.79       0.08       0.63       1       0.99       0.97       1       0.99       0.76       0.91       0.7         kNN<		scID-	0.75	0.59	0.95	0.85	0.8	1	1	0.42	1	0.95	0.63	0.61	0.42
kNN-       0.91       0.95       0.95       0.85       0.03       1       0.98       0.92       1       0.64       0.13       0.45       0.5         SCINA-       C <thc< th="">       C       C       <thc< <="" td=""><td></td><td>Cell_BLAST-</td><td>0.11</td><td>0.89</td><td>0.79</td><td>0.08</td><td>0.63</td><td>1</td><td>0.99</td><td>0.97</td><td>1</td><td>0.99</td><td>0.76</td><td>0.91</td><td>0.74</td></thc<></thc<>		Cell_BLAST-	0.11	0.89	0.79	0.08	0.63	1	0.99	0.97	1	0.99	0.76	0.91	0.74
SCINA-       Image: Colore of the colore of th	kNN			0.95	0.95	0.85	0.03	1	0.98	0.92	1	0.64	0.13	0.45	0.54
DigitalCellSorter- Garnett <sub>CV</sub> - Moana-       Image: Comparison of the comparison													1*	1*	
Garnett <sub>CV</sub> -       Garnett <sub>pretrained</sub> -       Image: Comparison of the		DigitalCellSorter-												0.99*	0.78*
Garnettpretrained-       Image: Construction of the second s		Garnett <sub>CV</sub> -												0.94*	0.6*
Moana-		Garnettpretrained -												0.98*	0.54*
Median F1-score     Garnett <sub>DE</sub> 0     0.25       0.75       0       0.25       0.75       0       0.75       0       0.75       0       0.75       0       0.75       0.75       0       0.75       0       0.75       0       0.75       0       0.75       0       0.75       0       0.75       0       0.75       0       0.75       0       0.75        0.75 <td< td=""><td></td><td>Moana-</td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td>0.93*</td><td>0.5*</td></td<>		Moana-												0.93*	0.5*
0 0.25 0.5 0.75 Digital Call Sector	Median F1-score	Garnett <sub>DE</sub> -												0.65	0.37
		SCINA <sub>DE</sub> -												0.38	0.47
	0 0.20 0.0 0.75	DigitalCellSorter <sub>DE</sub> -												0	0
https://biocellgen-public.svi.edu.au/mig 2019 scrnaseq- mdble/ublai/s/biocellgen-public.svi.edu.au/mig 2019 scrnaseq- bittps://biocellgen-public.svi.edu.au/mig 2019 scrnaseq- mdble/ublai/s/biocellgen-public.svi.edu.au/mig 2019 scrnaseq- bittps://biocellgen-public.svi.edu.au/mig	https://btep.ccr.cancer.gov/wp- content/uploads/Celltype Annota https://biocellgen-public.svi.edu.	tion final.pdf au/mig 2019 scrnaseq-	aron Mouse-	aron Human-	Muraro-	Segerstolpe-	Xin-	10X-	CEL-Seq2-	-MT	AMB3-	AMB16-	AMB92-	heng sorted-	Zheng 68K-

Madian Ed acare

– DigitalCelllSorter, Moana

https://bioconductor.org/books/release/OSCA/cell-type-annotation.html

#### SingleR: Reference-based annotation of scRNA-seq

- SingleR pipeline is based on correlating reference bulk transcriptomic data sets of pure cell types with single-cell gene expression.
- Reference set: a comprehensive transcriptomic dataset (microarray or RNA-seq) of pure cell types
- Human
  - Human Primary Cell Atlas (HPCA) : 38 main cell types, 169 subtypes, 713 samples
  - Blueprint+Encode: 43 cell types, 259 bulk RNAseq samples
- Mouse
  - Immunological Genome Project (ImmGen) : 20 main cell types, 830 microarray samples
  - mouse RNA-seq samples (brain-specific) : 28 cell types, 358 RNA-seq samples



#### SingleR: Reference-based annotation of scRNA-seq



# Step 1: Identifying variable genes among cell types in the reference set

- For each cell type, identify the top *N* variable genes that have a higher median expression in that cell type than in every other cell type
- Take the 'red' cell type as an example
  - For every gene, median expression values grouped by cell type were obtained.
  - Differential expression between each other cell type and the 'red' cell type was calculated and all genes with positive differential expression values were selected.
  - All selected genes were sorted by differential expression values, and then the top N genes were selected as variable genes for the **'red'** cell type.



https://biocellgen-public.svi.edu.au/mig\_2019\_scrnaseq-workshop/public/clustering-and-cell-annotation.html

Aran, D., Looney, A.P., Liu, L. et al. Reference-based analysis of lung single-cell sequencing reveals a transitional profibrotic macrophage. Nat Immunol 20, 163–172 (2019). https://doi.org/10.1038/s41590-018-0276-y

# Step 2: Correlating each single-cell transcriptome with each sample in the reference set

- Spearman coefficient is calculated for single cell expression with each of the samples in the reference dataset.
- The correlation analysis is performed only on variable genes in the reference dataset.



# Step 3: Iterative fine-tuning - reducing the reference to only top cell types

- For a single cell and each cell type, multiple Spearman correlation coefficients are aggregated into a "cell-type score"
  - The SingleR score for each cell type is the 80 percentile in each of the boxplots.
- Cell types with the lowest score or a score below will be removed
- Repeat from step 1 until only one cell type remained



#### SingleR: Reference-based annotation of scRNA-seq



Aran, D., Looney, A.P., Liu, L. *et al.* Reference-based analysis of lung single-cell sequencing reveals a transitional profibrotic macrophage. *Nat Immunol* **20**, 163–172 (2019). <u>https://doi.org/10.1038/s41590-018-0276-y</u>

## Outline

- scRNA-seq data analysis
  - Cell type annotation
    - SingleR
  - Cell type markers identification
  - Pseudo timing
    - Monocle
  - Cell-type gene regulatory networks
    - SCENIC
    - BEELINE
  - Single cell deconvolution
    - CIBERSORTx

## Cell type markers identification Differential expression analysis

- Non-parametric tests
  - Wilcoxon rank sum test
  - Student's t-test
- Methods specific for scRNA-seq
  - MAST : GLM-framework that treats cellular detection rate as a covariate (Finak et al, Genome Biology, 2015)
- Methods for bulk RNA-seq
  - DESeq2 : DE based on a model using the negative binomial distribution (Love et al, Genome Biology 2014)

• https://satijalab.org/seurat/archive/v3.1/immune\_alignment.html

Finak, G., McDavid, A., Yajima, M. et al. MAST: a flexible statistical framework for assessing transcriptional changes and characterizing heterogeneity in single-cell RNA sequencing data. Genome Biol 16, 278 (2015). <a href="https://doi.org/10.1186/s13059-015-0844-5">https://doi.org/10.1186/s13059-015-0844-5</a>

<sup>•</sup> Love MI, Huber W, Anders S (2014). "Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2." Genome Biology, 15, 550. doi: 10.1186/s13059-014-0550-8.

## Cell type markers identification Differential testing and visualization in Scanpy



sc.tl.rank\_genes\_groups(adata, 'leiden', method='wilcoxon')
sc.tl.rank\_genes\_groups(adata, 'leiden', method='t-test')

sc.pl.dotplot(adata, marker\_genes, groupby='leiden')



 Finak, G., McDavid, A., Yajima, M. *et al.* MAST: a flexible statistical framework for assessing transcriptional changes and characterizing heterogeneity in single-cell RNA sequencing data. *Genome Biol* 16, 278 (2015). https://doi.org/10.1186/s13059-015-0844-5

- Love MI, Huber W, Anders S (2014). "Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2." *Genome Biology*, 15, 550. doi: <u>10.1186/s13059-014-0550-8</u>.
- <u>https://satijalab.org/seurat/archive/v3.1/immune\_alignment.html</u>
- https://zenodo.org/record/4317764#.YII7gdPMKCg

## Outline

- scRNA-seq data analysis
  - Cell type annotation
    - SingleR
  - Cell type markers identification
  - Pseudo timing
    - Monocle
  - Cell-type gene regulatory networks
    - SCENIC
    - BEELINE
  - Single cell deconvolution
    - CIBERSORTx

## **Pseudo timing**

- Many cell differentiation processes take place during development
- We order the cells along one or more trajectories representing the underlying developmental processes
- This ordering is called 'pseudotime'
- Trajectory inference (TI) aims to reconstruct a cellular dynamic process



http://cole-trapnell-lab.github.io/monocle-releas https://scrnaseq-course.cog.sanger.ac.uk/website

Saelens, W., Cannoodt, R., Todorov, H. *et al.* A comparison of single-cell trajectory inference methods. *Nat Biotechnol* **37**, 547–554 (2019). <u>https://doi.org/10.1038/s41587-019-0071-9</u> Trapnell, C., Cacchiarelli, D., Grimsby, J. *et al.* The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. *Nat Biotechnol* **32**, 381–386 (2014). https://doi.org/10.1038/nbt.2859

## **Pseudo timing**

 Using single-cell-omics data, many trajectory inference (TI) methods could computationally order cells along trajectories, allowing the unbiased study of cellular dynamic processes



http://cole-trapnell-lab.github.io/monocle-release/docs/#constructing-single-cell-trajectories

https://scrnaseq-course.cog.sanger.ac.uk/website/biological-analysis.html#pseudotime-analysis

Saelens, W., Cannoodt, R., Todorov, H. et al. A comparison of single-cell trajectory inference methods. Nat Biotechnol 37, 547-554 (2019). https://doi.org/10.1038/s41587-019-0071-9

## Monocle

#### Constructing single cell trajectories

Monocle, an unsupervised algorithm to build single-cell trajectories, and find cell fate decisions and dynamically regulated genes.

- Step 1: Choose genes that define progress
- Step 2: Reduce data dimensionality
  - independent component analysis (ICA)
- Step 3: Construct minimum spanning tree (MST) on the cells
- Step 4: Find the longest path through the MST
- Step 5: Order cells along the trajectory

http://cole-trapnell-lab.github.io/monocle-release/docs/#constructing-single-cell-trajectories

Trapnell, C., Cacchiarelli, D., Grimsby, J. *et al.* The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. *Nat Biotechnol* **32**, 381–386 (2014). <u>https://doi.org/10.1038/nbt.2859</u>

## Step 1: Choose genes that define progress

• Represent the expression profile of each cell as a point in a highdimensional Euclidean space, with one dimension for each gene



## Step 2: Reduce data dimensionality

- Reduce dimensionality using independent component analysis (ICA)
- Transform the cell data from a high-dimensional space into a low-dimensional one that preserves essential relationships between cell populations



https://github.com/NBISweden/excelerate-scRNAseq/blob/master/session-trajectories/trajectory\_inference\_analysis.pdf

https://www.cs.cmu.edu/~tom/10701 sp11/recitations/Recitation 11.pdf

Trapnell, C., Cacchiarelli, D., Grimsby, J. *et al.* The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. *Nat Biotechnol* **32**, 381–386 (2014). <u>https://doi.org/10.1038/nbt.2859</u>

#### ICA

- Assumption: the mixed sources signals are independent of each other
- Goal: find linear mapping W which maximize independence and unmix sources signal s
   Mixed variables



https://github.com/NBISweden/excelerate-scRNAseq/blob/master/session-trajectories/trajectory\_inference\_analysis.pdf https://www.slideserve.com/vladimir-kirkland/ica-and-isa-using-schweizer-wolff-measure-of-dependence

## ICA vs PCA

- PCA : Find the directions of maximal variance
- ICA : Find the directions of maximal independence
  - The values in each source have non-Gaussian distributions



https://github.com/NBISweden/excelerate-scRNAseq/blob/master/session-trajectories/trajectory\_inference\_analysis.pdf https://scikit-learn.org/stable/modules/decomposition.html#independent-component-analysis-ica

### Why ICA



# Step 3: Construct minimum spanning tree (MST) on the cells

- Minimum spanning tree (MST)
  - The undirected graph connecting all vertices with the smallest sum of all distances
     Vertex : Cell



Trapnell, C., Cacchiarelli, D., Grimsby, J. *et al.* The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. Nat Biotechnol 32, 381–386 (2014). <u>https://doi.org/10.1038/nbt.2859</u> https://en.wikipedia.org/wiki/Minimum\_spanning\_tree

# Step 4: Find the longest path through the MST

Correspond to the longest sequence of similar cells (e.g., gene expression)



## Step 5: Order cells along the trajectory

 Produce a 'trajectory' of an individual cell's progress through differentiation



#### Developmental trajectory of olfactory neurons in mice

- Each point is a cell, which is connected to an MST
- The pseudotime value of each cell is measured as the distance along the trajectory from its position back to the beginning



## Pseudo timing

• The performance of TI methods mostly depend on the topology of the trajectory in the single-cell data.

Method	SCUBA pseudotime	Wanderlust	Wishbone	SLICER	SCOUP	Waterfall	Mpath	TSCAN	Monocle	SCUBA	
Visual abstract		10000 g		-			••••	K	AL AL	T0 T1 T2 T3	
Structure	Linear	Linear	Single bifurcation	Branching	Branching	Linear	Branching	Linear	Branching	Branching	
Robustness strategy	Principal curves	Ensemble, starting cell	Ensemble, starting cell	Starting cell	Starting population	Clustering of cells	Clustering of cells using external labelling	Clustering of cells	Differential expression	Simple model	
Extra input requirements	None	Starting cell	Starting cell	Starting cell	Starting population	None	Time points	None	Time points	Time points	
Unbiased	+	±	±	±	±	+	-	+	-	-	
Scalability w.r.t. cells	-	-	±	±	-	±	+	+	-	±	
Scalability w.r.t. genes	+	+	+	+	-	+	±	±	±	+	
Code and documentation	-	±	+	±	+	±	+	+	+	±	
Parameter ease-of-use	+	+	+	+	-	±	-	+	+	+	

http://cole-trapnell-lab.github.io/monocle-release/docs/#constructing-single-cell-trajectories

https://indico.math.cnrs.fr/event/3780/contributions/3242/attachments/2195/2550/Slides-maugis-181018.pdf

https://scrnaseq-course.cog.sanger.ac.uk/website/biological-analysis.html#pseudotime-analysis

Saelens, W., Cannoodt, R., Todorov, H. et al. A comparison of single-cell trajectory inference methods. Nat Biotechnol 37, 547-554 (2019). https://doi.org/10.1038/s41587-019-0071-9

## Outline

- scRNA-seq data analysis
  - Cell type annotation
    - SingleR
  - Cell type markers identification
  - Pseudo timing
    - Monocle
  - Cell-type gene regulatory networks
    - SCENIC
    - BEELINE
  - Single cell deconvolution
    - CIBERSORTx

## **Gene regulation**



Gene regulation is the process of controlling which genes in a cell's DNA are expressed (used to make a functional product such as a protein).

 $https://www.cs.purdue.edu/homes/ayg/TALKS/STC\_CHICAGO10/Introductory\_material/regulatory\_networks.ppt$ 

# Gene regulatory network

- Gene regulatory networks (GRNs) like on-off switches of a cell operating at the gene level
- Two genes are connected if the expression of one gene modulates expression of another one by either activation or inhibition
- GRN can be inferred from correlations in gene expression data, time-series gene expression data, and/or gene knock-out experiments





 $https://www.cs.purdue.edu/homes/ayg/TALKS/STC\_CHICAGO10/Introductory\_material/regulatory\_networks.ppt$ 

# Cell-type gene regulatory networks

 Cell-type-specific GRNs would be key tools for the study of cellular heterogeneity



# SCENIC

#### single-cell regulatory network inference and clustering

- Simultaneously reconstruct gene regulatory networks and identify stable cell states from single-cell RNA-seq data, based on three tools
  - GENIE3 or GRNboost
  - RcisTarget
  - AUCell
- The gene regulatory network is inferred based on co-expression and DNA motif analysis, and then the network activity is analyzed in each cell to identify the recurrent cellular states.



## Step 1: TF-based co-expression network



# Step 2: Identification of transcription factor binding motifs



### Regulon

• Regulon: a group of genes that are regulated as a unit





## **AUCell score**

- AUCell uses the "Area Under the Curve" (AUC) to calculate whether a critical subset of the input gene set is enriched within the expressed genes for each cell.
- AUCell score: measure how active a regulon is in a cell
  - Step 1: For each cell, build gene-expression ranking
  - Step 2: Calculate enrichment for the gene signatures (AUC)
  - Step 3: Determine the cells with given regulon



## Step 3: Regulon activities in each cell



Regulons

TF1 regulon

TF2 regulon

.

# Top regulons on the Mouse brain



# Microglia GRN on the Mouse brain

- The regulons associated to microglia can be summarized based on the binding motif of the associated TF.
- The predicted network for microglia contains many well-known regulators of microglial fate and/or microglial activation, including PU.1, Nfkb, Irf, and AP-1/Maf.



# **Evaluate GRN inference methods**

- Ground truth of GRNs is usually unknown.
- How do we evaluate the performance of existing GRN inference methods from scRNA-seq data?



## BEELINE

# Benchmarking gene regulatory network inference from single-cell transcriptomic data

- BEELINE is an evaluation framework incorporating 12 diverse GRN inference algorithms to assess the accuracy, robustness, and efficiency of GRN inference techniques for single-cell gene expression data.
- Step 1: preprocessing
  - Input type 1: Simulated data from synthetic networks
  - Input type 2: Simulated data from curated models
  - Input type 3: Experimental single-cell RNA-seq datasets
- Step 2: Run GRN inference algorithms
- Step 3: postprocessing and evaluation



## BoolODE for simulation Generate inputs for BEELINE

- Convert a Boolean model into a system of stochastic differential equations
- Read a model definition file, and outputs the simulated expression data for a given model.



#### Input type 1: Simulated datasets from synthetic networks

Synthetic network types:

- Linear
- Linear long
- Cycle
- Bifurcating
- Bifurcating converging
- Trifurcating

Simulated data from synthetic networks



44



### Input type 1:

#### Simulated datasets from synthetic networks

**Bifurcating trajectory** 

- The color of each 'cell' is determined by the timepoint at which it was sampled in the simulation
- Darker colors indicate earlier time points



## Input type 2:

#### Simulated datasets from curated models (Boolean)

Four published Boolean models:

- Mammalian Cortical Area Development (mCAD)<sup>1</sup>
- Ventral Spinal Cord Development (VSC)<sup>2</sup>
- Hematopoietic Stem Cell Differentiation (HSC)<sup>3</sup>
- Gonadal Sex Determination (GSD)<sup>4</sup>

mCAD





- Giacomantonio CE & Goodhill GJ A boolean model of the gene regulatory network underlying mammalian cortical area development. PLoS Comput. Biol. 6, e1000936 (2010).
- Lovrics A et al. Boolean Modelling reveals new regulatory connections between transcription factors orchestrating the development of the ventral spinal cord. PLoS One 9, e111430 (2014).
- Krumsiek J, Marr C, Schroeder T & Theis FJ Hierarchical Differentiation of Myeloid Progenitors Is Encoded in the Transcription Tactor Network. PLoS One 6, e22649 (2011).
- Ríos O et al. A Boolean network model of human gonadal sex determination. Theor. Biol. Med. Model. 12, 26 (2015).

### Input type 2:

Simulated datasets from curated models (Boolean)

mCAD

- Visualizations of t-SNE from the BoolODE output
- The color of each point indicates the corresponding simulation time.



Giacomantonio CE & Goodhill GJ A boolean model of the gene regulatory network underlying mammalian cortical area development. *PLoS Comput. Biol.* 6, e1000936 (2010). Lovrics A et al. Boolean Modelling reveals new regulatory connections between transcription factors orchestrating the development of the ventral spinal cord. *PLoS One* 9, e111430 (2014). Krumsiek J, Marr C, Schroeder T & Theis FJ Hierarchical Differentiation of Myeloid Progenitors Is Encoded in the Transcription Tactor Network. *PLoS One* 6, e22649 (2011). Rios O et al. A Boolean network model of human gonadal sex determination. *Theor. Biol. Med. Model.* 12, 26 (2015).

#### Input type 3: Experimental scRNA-seq datasets (Mouse)

Single cell RNAseq data



#### Embryonic stem cells



: hematopoietic stem and progenitor cell differentiation. Blood 128, 20-31 (2016).

• Hayasni, 1. et al. Single-cell full-length total KINA sequencing uncovers dynamics of recursive splicing and enhancer RNAs. Nat. Commun. 9, 619 (2018).

• Shalek, A. K. et al. Single-cell RNA-seq reveals dynamic paracrine control of cellular variation. Nature 510, 363–369 (2014).

Camp, J. G. et al. Multilineage communication regulates human liver bud development from pluripotency. Nature 546, 533–538 (2017).

• Chu, L. F. et al. Single-cell RNA-seq reveals novel regulators of human embryonic stem cell differentiation to definitive endoderm. Genome Biol. 17, 173 (2016).

#### Input type 3: Experimental scRNA-seq datasets (Human)



Nestorowa, S. et al. A single-cell resolution map of mouse hematopoietic stem and progenitor cell differentiation. Blood 128, 20–31 (2016).

- Hayashi, T. et al. Single-cell full-length total RNA sequencing uncovers dynamics of recursive splicing and enhancer RNAs. Nat. Commun. 9, 619 (2018).
- Shalek, A. K. et al. Single-cell RNA-seq reveals dynamic paracrine control of cellular variation. Nature 510, 363–369 (2014).
- Camp, J. G. et al. Multilineage communication regulates human liver bud development from pluripotency. Nature 546, 533–538 (2017).
- Chu, L. F. et al. Single-cell RNA-seq reveals novel regulators of human embryonic stem cell differentiation to definitive endoderm. Genome Biol. 17, 173 (2016).

## Run GRN inference algorithms

• Incorporation of 12 diverse GRN inference algorithms





## **Evaluation**

- Accuracy (AUPRC and early precision)
- Stability of results (across simulations, in the presence of dropouts and across algorithms)
- Analysis of network motifs
- Scalability

#### **BEELINE** summary

	Properties						Accuracy Stability						Scalability (genes)							
	Calegory	Addi. Inputs	Timeor	dered?	ed? ned?	writhe	uc unate	CRIMAS	Datas	Runs	Dropo	Pseud	otime 100	Tir	ne	21	100	Men	nory	
PIDC	MI	-	X	×	X							_	15	1m	5m	30m	0.1G	0.1G	0.5G	2K
GENIE3	RF	-	×	1	×							-	5m	1h	Зh	12h	1G	2G	2G	2G
GRNBOOST2	RF	-	×	1	X							-	1m	10m	30m	1h	0.1G	0.1G	0.5G	1G
SCODE	ODE+Reg	ODE parameters	1	1	1								1m	5m	5m	30m	1M	0.1G	0.1G	0.5G
PPCOR	Corr	-	X	X	1							-	15	1s	1s	1s	1M	0.1G	0.1G	0.1G
SINCERITIES	Reg	-	1	1	1								15	1m	5m	10m	0.1G	0.1G	0.1G	0.5G
SCRIBE	MI	Type of RDI	1	1	X			-					5m	2h	6h	-	0.1G	0.1G	0.1G	-
SINGE	GC	Regression parameters	1	1	×			-					Зh	>1d	>1d	-	0.5G	0.5G	1G	-
LEAP	Corr	Lag	1	1	X			-					15	1s	1m	5m	1M	0.1G	0.1G	0.5G
GRISLI	ODE+Reg	Regression parameters	1	1	×			-					5m	1h	Зh	-	0.5G	>4G	>4G	-
GRNVBEM	Reg	-	1	1	1			-					1m	>1d	-	-	0.1G	2G	-	-
SCNS	Bool	Boolean model parameters	1	1	1			-				-	-	-	-	-	-	-	-	-
		Low/	Poor		High	/God	bd		L	ow/Po	oor		High	/Goo	d					

Pratapa A, Jalihal AP, Law JN, Bharadwaj A, Murali TM. Benchmarking algorithms for gene regulatory network inference from single-cell transcriptomic data. *Nat Methods*. 2020;17(2):147-154. doi:10.1038/s41592-019-0690-6

## Network inference algorithms summary



Todorov H., Cannoodt R., Saelens W., Saeys Y. (2019) Network Inference from Single-Cell Transcriptomic Data. In: Sanguinetti G., Huynh-Thu V. (eds) Gene Regulatory Networks. Methods in Molecular Biology, vol 1883. Humana Press, New York, NY. https://doi.org/10.1007/978-1-4939-8882-2\_10

## Outline

- scRNA-seq data analysis
  - Cell type annotation
    - SingleR
  - Cell type markers identification
  - Pseudo timing
    - Monocle
  - Cell-type gene regulatory networks
    - SCENIC
    - BEELINE
  - Single cell deconvolution
    - CIBERSORTx

# Single cell deconvolution

#### **Pros for Bulk-seq**

- Can assay entire sample at once
- Can help identify transcription changes in individual cell types
- Huge amount of data out there already
- Cheap

#### Cons

Lose single cell information

#### Bulk

#### \$200/sample (Novogene)



Bulk transcriptomic analyses lose single cell information

How to computationally figure out what went into the mixture?

#### Single cell \$4000 ~ 10000/sample



Single cell transcriptomic analyses retain single cell

https://projects.iq.harvar d.edu/files/chanbioinfor matics/files/cell\_type\_de convolution.pdf ASHG 2019 scRNAseq HiPlex oral presentation https://www.youtube.com/ watch?v=YlRemO\_TE3Y

55

# Single cell deconvolution

• Qualify cell types in mixture



## CIBERSORTx

cell-type identification by estimating relative subsets of RNA transcripts

Infer cell-type-specific gene expression profiles without physical cell isolation



## CIBERSORTx

cell-type identification by estimating relative subsets of RNA transcripts



## MuSiC

#### Multi-Subject Single Cell deconvolution

• A method for characterizing the cell type composition of large amounts of RNA sequencing data in complex tissues using cell type-specific gene expression in single-cell RNA sequencing (RNA-seq) data



## **BSEQ-sc**

Deconvolution of Bulk Sequencing Experiments using Single Cell Data

 Leverage single-cell sequencing data to estimate cell type proportion and cell type-specific gene expression differences from RNA-seq data from bulk tissue samples



A single-cell transcriptomic map of the human and mouse pancreas reveals inter- and intra-cell population structure M. Baron, A. Veres, S.L. Wolock, A. L. Faust, R. Gaujoux, A. Vetere, J. Hyoje Ryu, B. K. Wagner, S. Shen-Orr, A. M. Klein, D. A. Melton, I. Yanai Cell Systems. 2016 Oct 26 <u>10.1016/j.cels.2016.08.011</u>

## Single cell RNA-seq data for human brain

- 8 excitatory and 8 inhibitory adult neuronal subtypes (i.e., cell expression clusters) ٠
- Major adult non-neuronal types: astrocytes, endothelial, microglia, ٠ oligodendrocytes, and oligodendrocyte progenitor (OPC), pericyte

Ex1 Ex2 Ex3 Ex4 Ex5 Exe Ex7 Ex8

Developmental neuronal and non-neuronal types

#### Read-count based; e.g., Transcripts Per Kilobase Million (TPM)



 $\sim 400$  cells (Darmanis et al., PNAS, 2015)

**Neuronal Subpopulations** Dim In7

> -20 20 Dim1 ~3000 cells (Lake et al., Science, 2016)



Molecular-count based; e.g., Unique molecular identifiers (UMI)



~10319 cells (Lake et al., Nature Biotech, 2018)



 $\sim 17,093$  cells (PsychENCODE)

## Single cell deconvolution Step 1: unsupervised learning to see brain cell types

Non-negative matrix factorization (NMF)

Single cell signatures

PNAS, 2015)

Science, 2016&2018)



## Single cell deconvolution Step 2: supervised learning to estimate cell fractions







#### Comparison with existing methods



Wang, et al., Science, 2018

# Neuronal and glial cell fraction changes in gender and disorders



Excitatory to Inhibitory imbalance at neuronal subtype level for ASD\*

#### Astrocyte and Microglia increase in ASD\*\*



\* Rubenstein et al., Model of autism: increased ratio of excitation/inhibition in key neural systems, Genes Brain Behav. 2003

\*\* Gandal et al., Shared molecular neuropathology across major psychiatric disorders parallels polygenic overlap, Science 2018



## anges in Age

#### natostatin (SST)



Astrocyte



Oligodendrocy





# Cell-type fraction changes in human brain development



## Resources

#### Tutorial

- <u>https://github.com/hbctraining/scRNA-seq</u>
- <u>https://bioconductor.org/books/release/OSCA/</u>
- <u>http://data-science-sequencing.github.io/</u>
- <u>https://broadinstitute.github.io/2019\_scWorkshop/</u>
- <u>https://biocellgen-public.svi.edu.au/mig\_2019\_scrnaseq-</u> workshop/public/index.html

#### Tools

• <u>https://github.com/seandavi/awesome-single-cell</u>