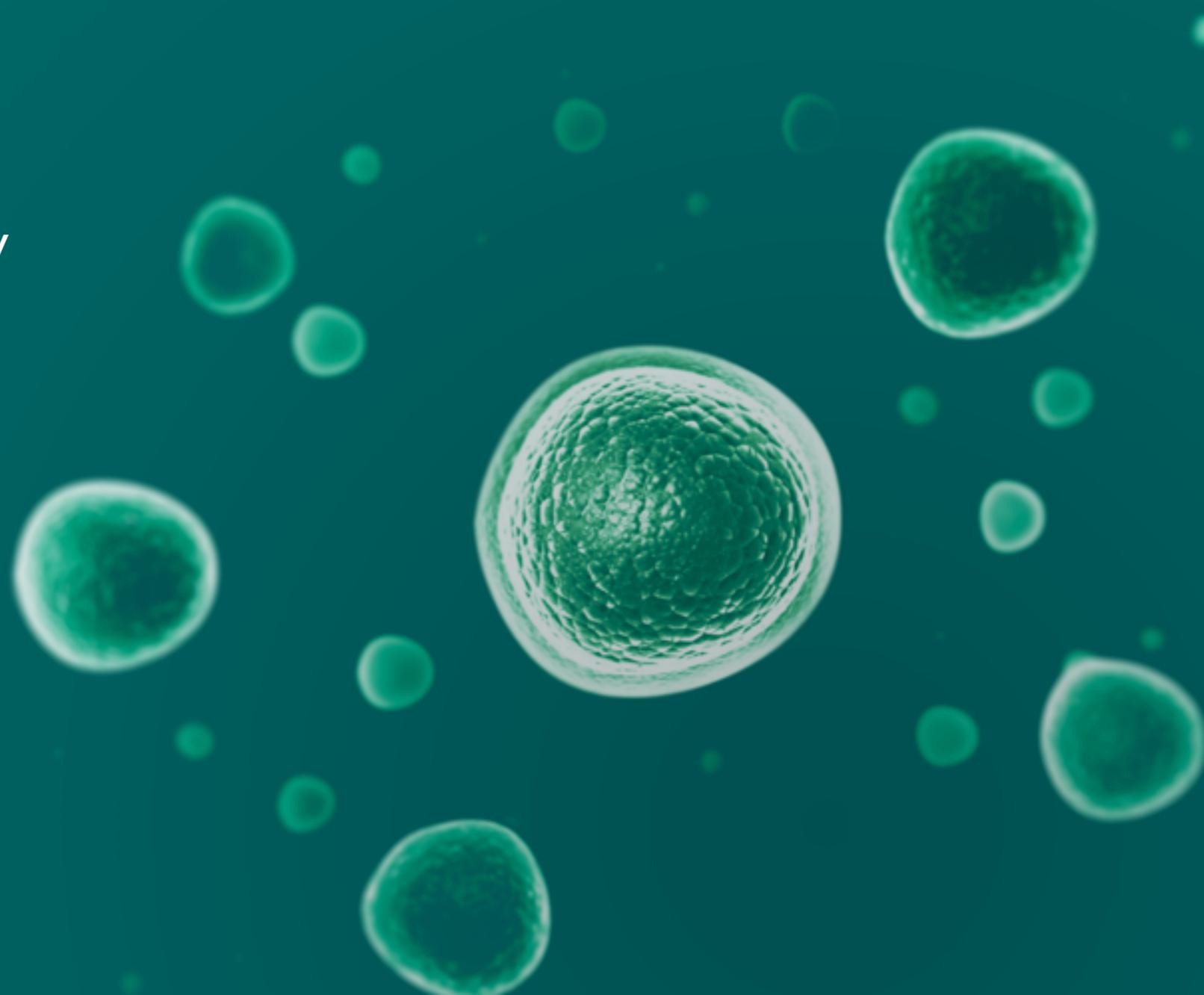


Trajectory or pseudotime analysis of single-cell RNA-seq data

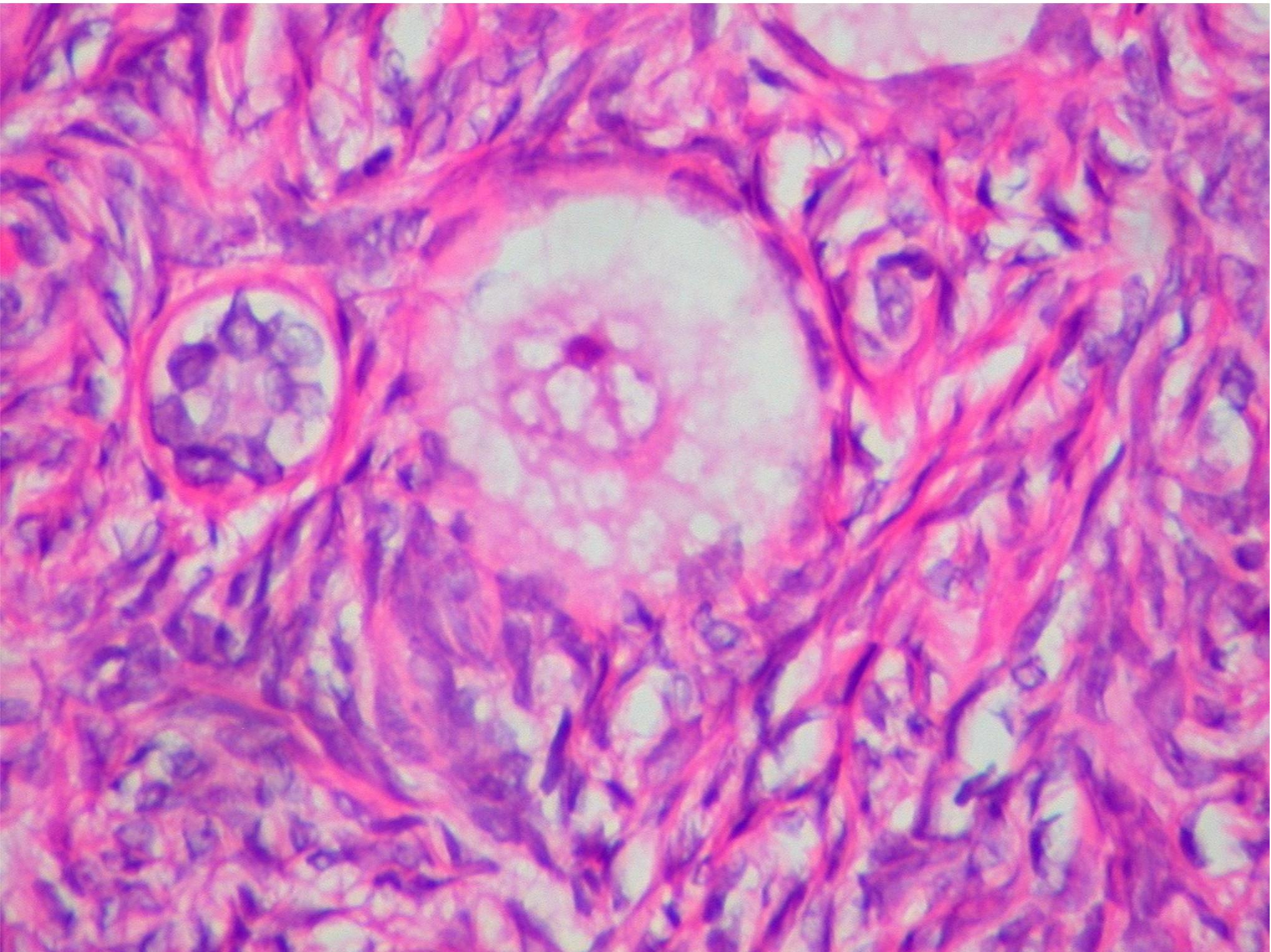
Davis McCarthy
NHMRC Early Career Fellow
Stegle Group

Physalia Course, Berlin
08 February 2018

www.ebi.ac.uk



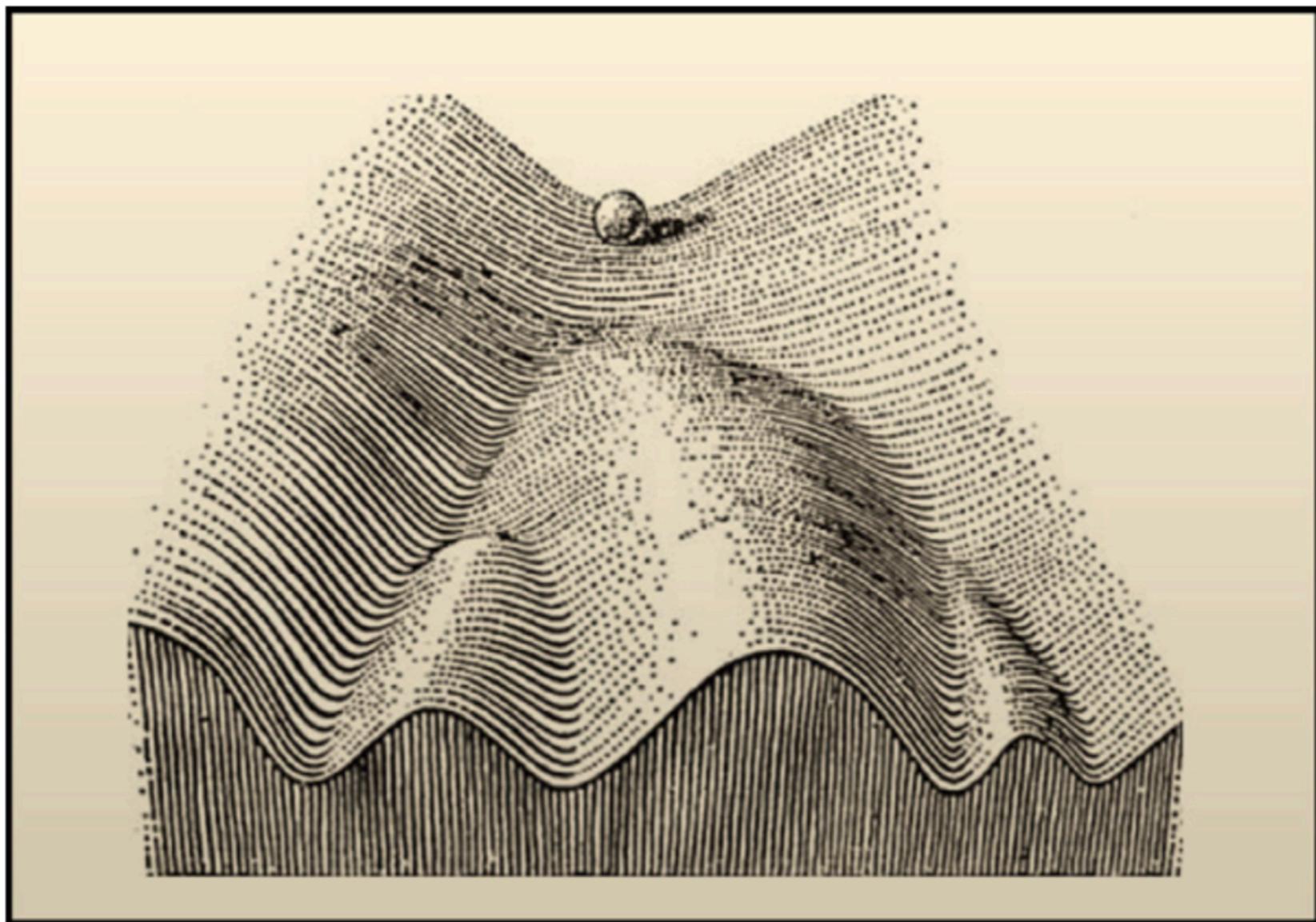
Importance of studying development



Photograph: Ed Uthman, Flickr Creative Commons

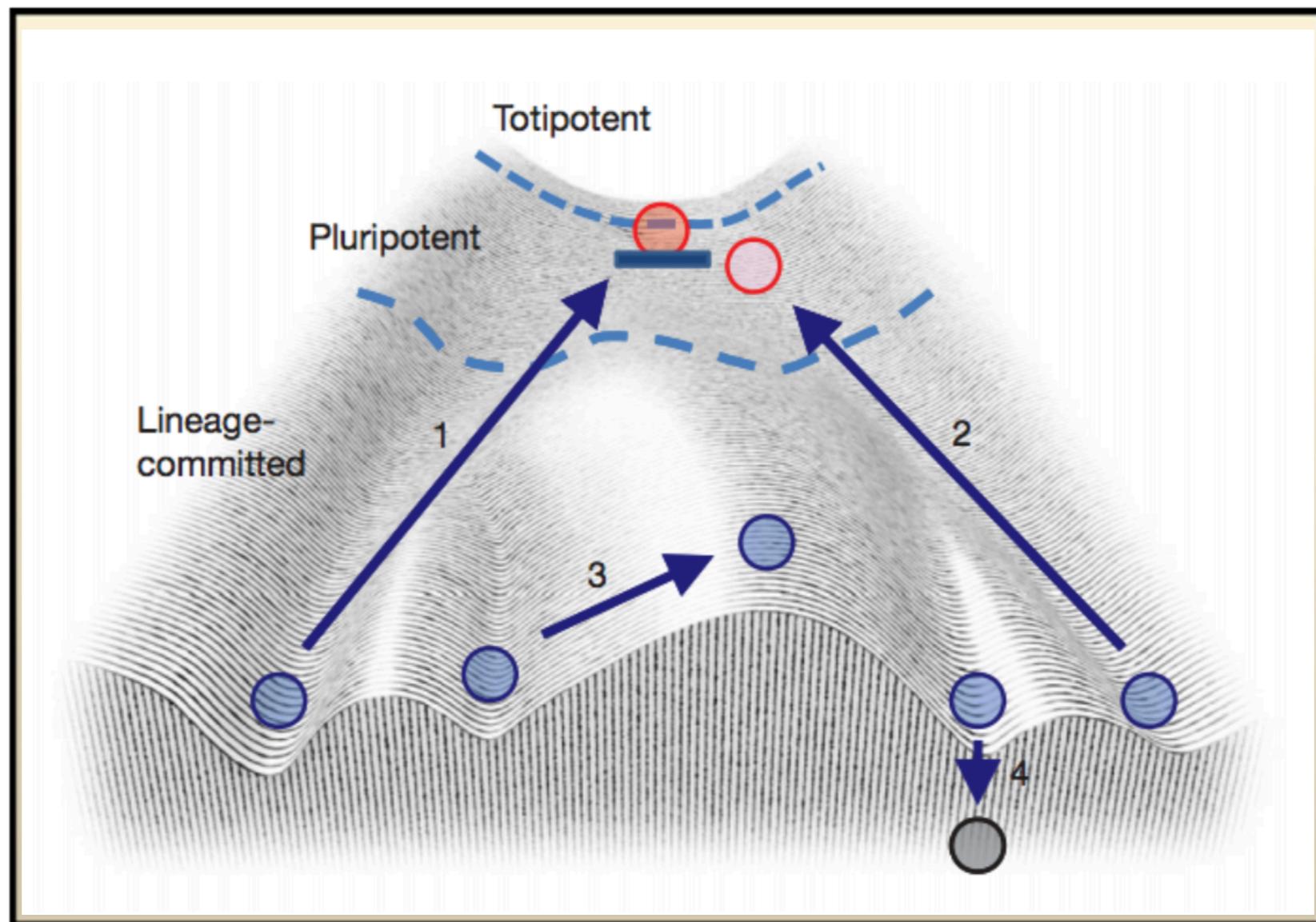
How do cells “make decisions” about their development?

Waddington's landscape



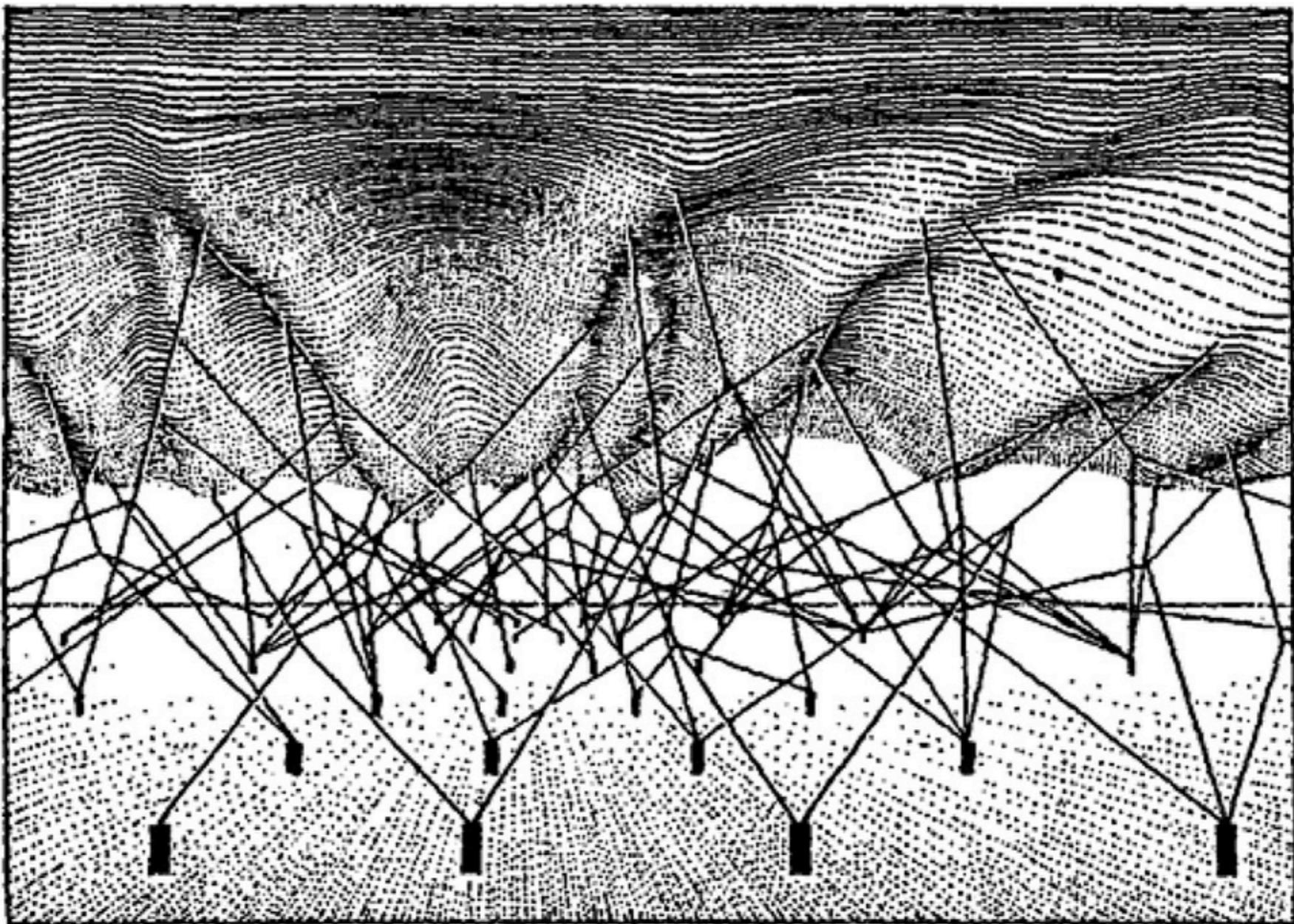
Waddington's **epigenetic landscape** is a metaphor for how **gene regulation** modulates development.

Waddington's landscape



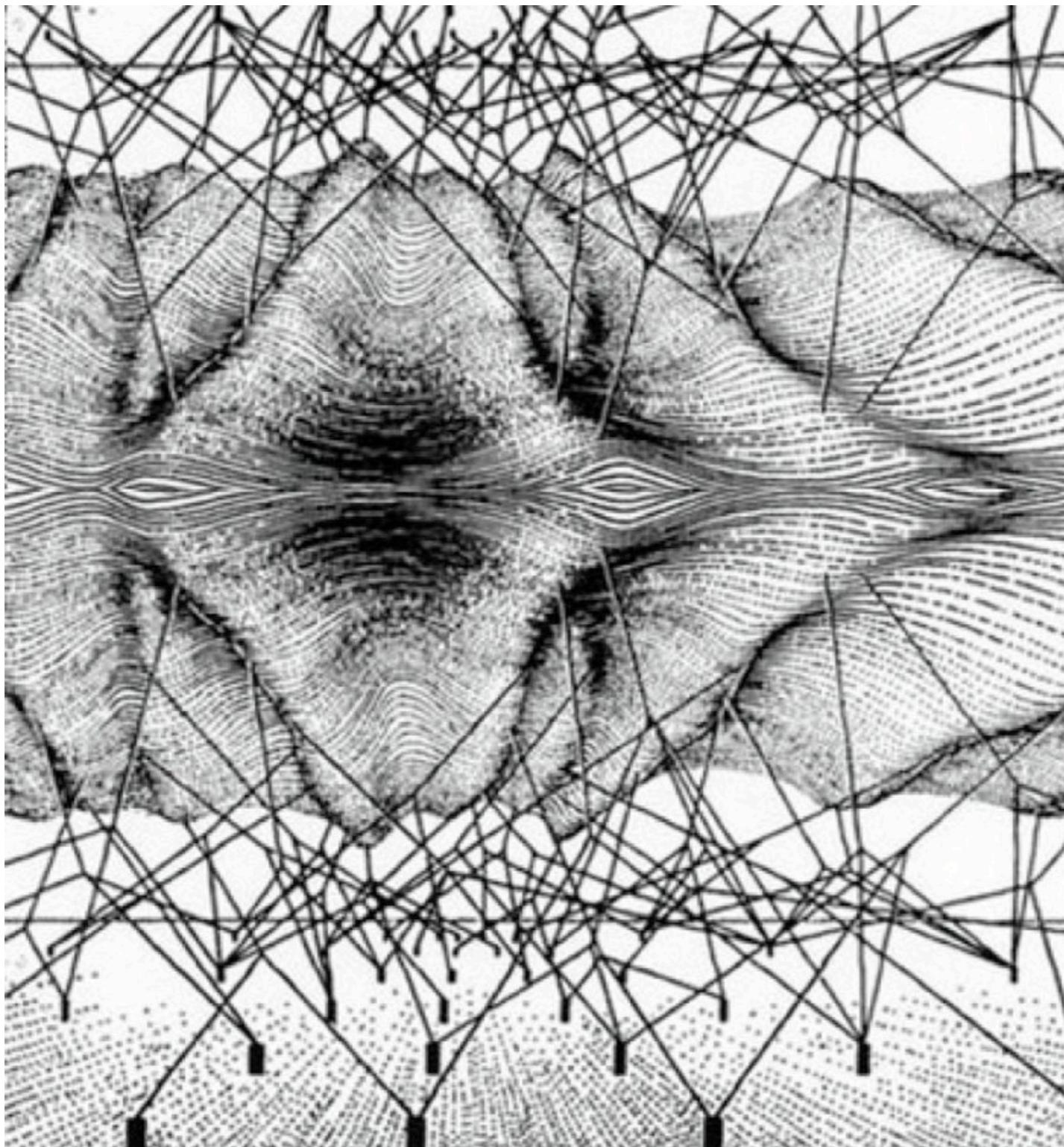
Waddington's **epigenetic landscape** is a metaphor for how **gene regulation** modulates development.

Waddington's landscape



Waddington's **epigenetic landscape** is a metaphor for how **gene regulation** modulates development.

Waddington's landscape



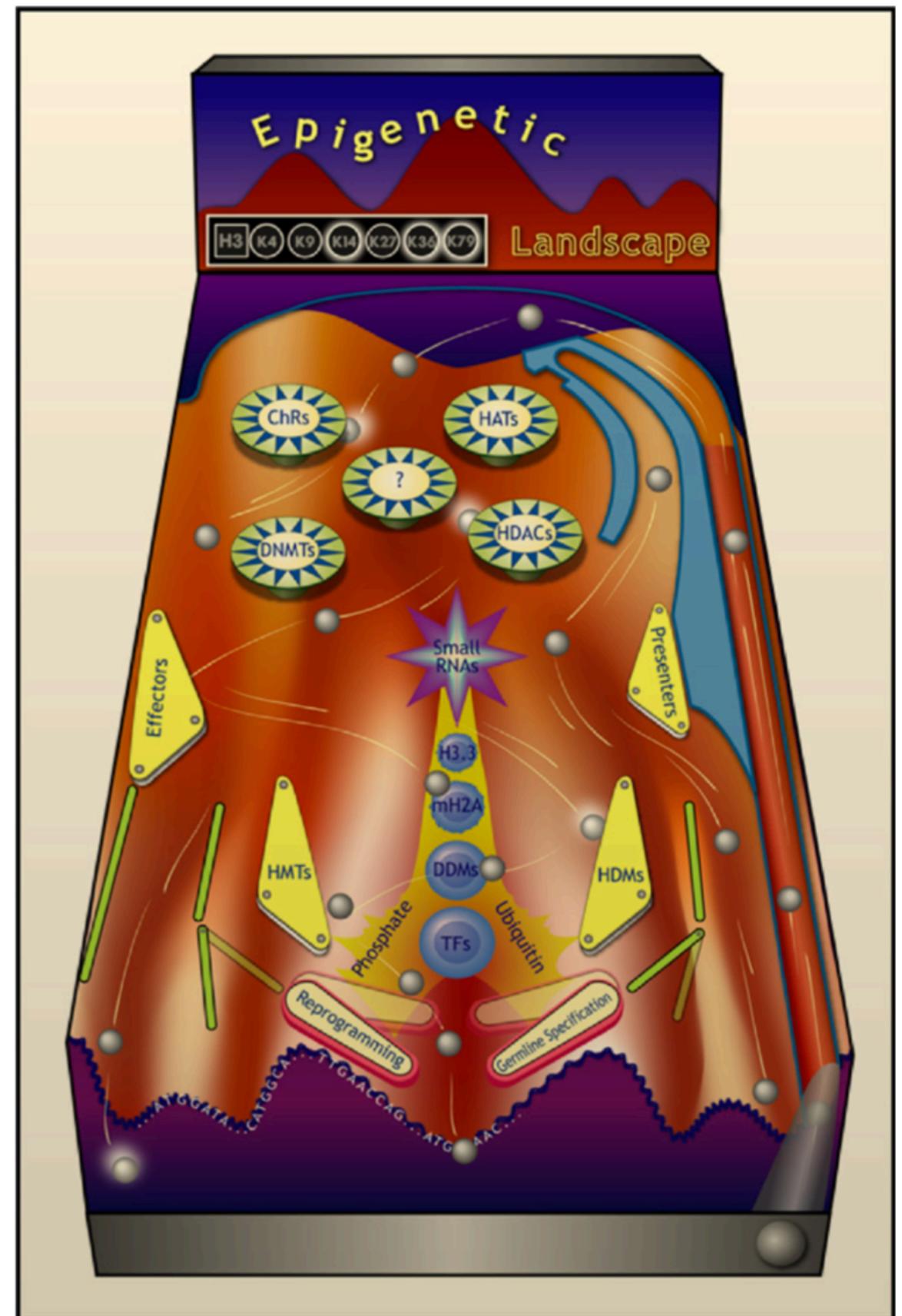
Influence of
environment

Developmental
landscape

Functional networks

Genes

Waddington's pinball landscape





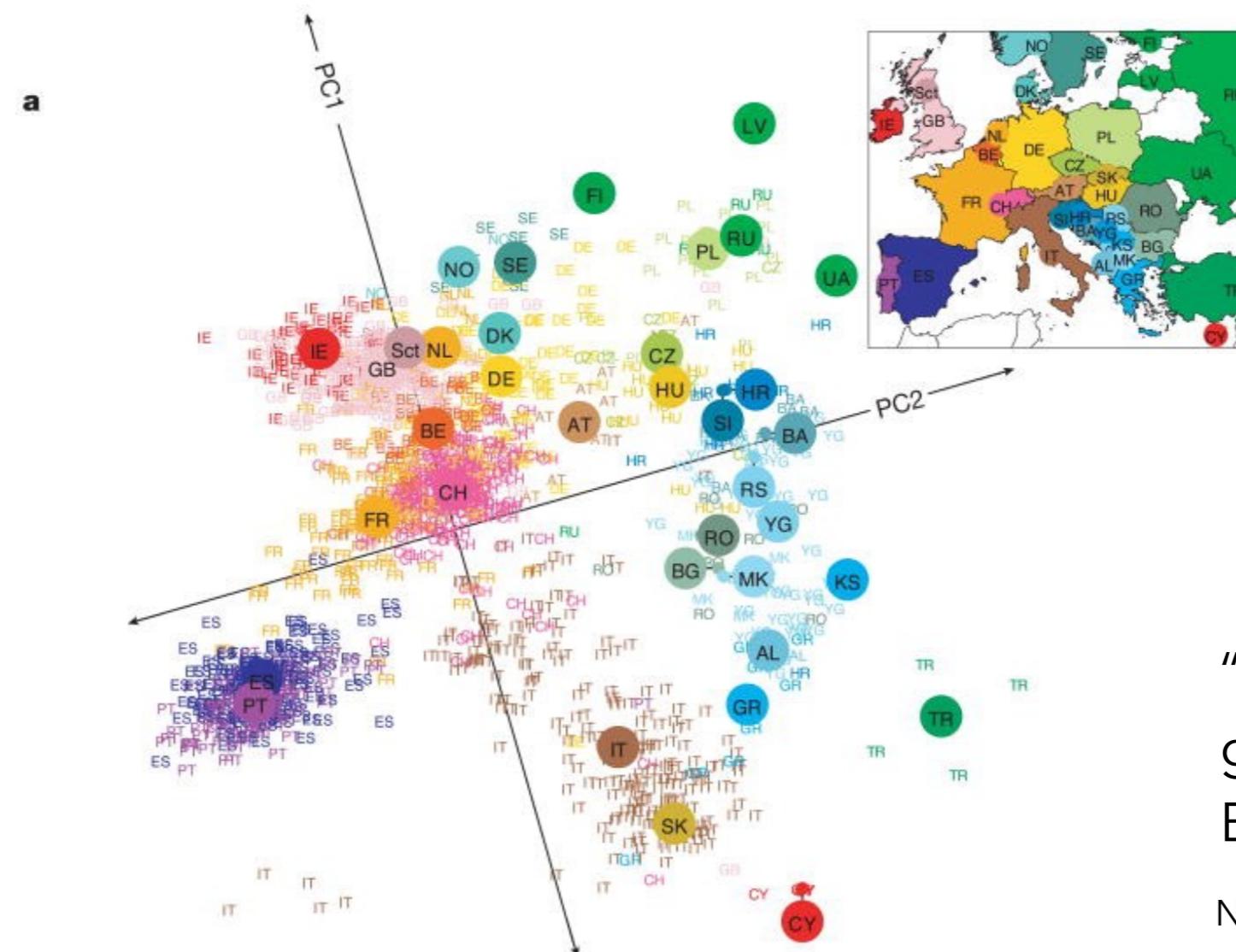


How do we represent extremely high dimensional data in a way that we can interpret?

The “crowding problem”. What do we lose? What can go wrong?

PCA: linear combinations of gene expression values
to maximise variance between cells

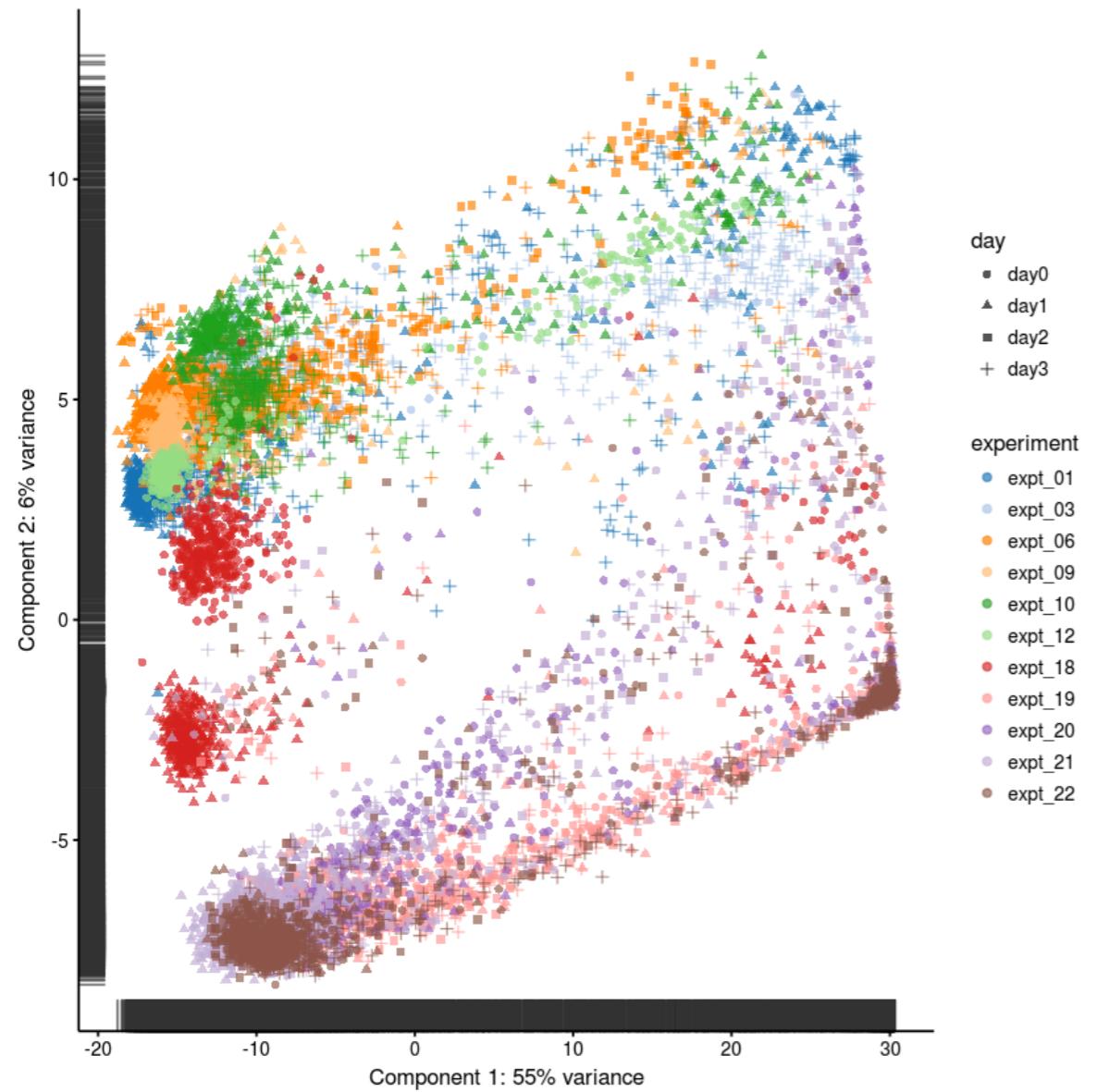
PCA: linear combinations of gene expression values
to maximise variance between cells



“Genes mirror
geography within
Europe”

Novembre et al, *Nature*, 2008

PCA on single-cell data typically captures technical effects, most often number of genes detected

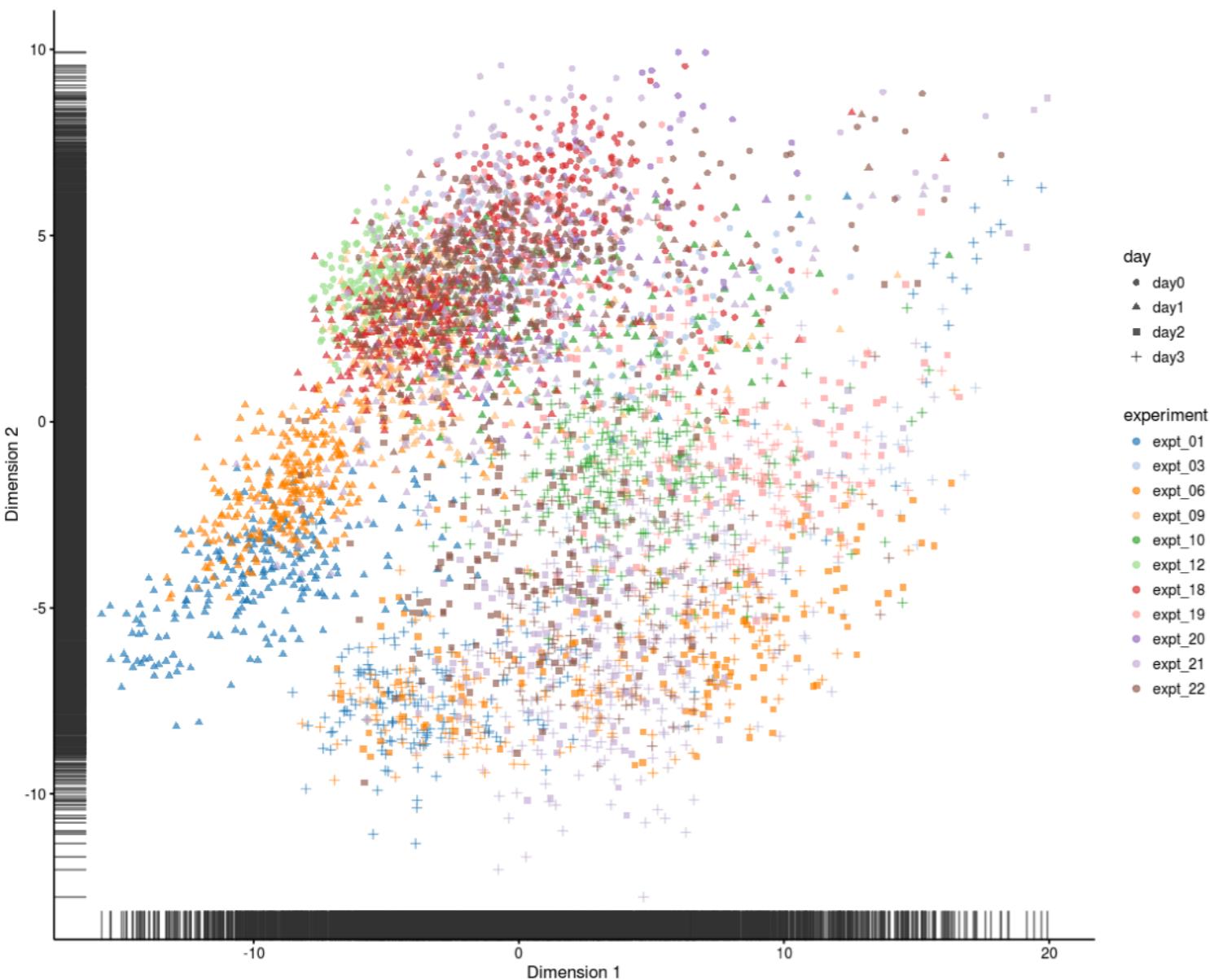


First principal component explains 55% of variance, very strongly correlated with number of genes detected per cell.

PCA extremely useful for QC.

Produced with the R/Bioconductor package scater (McCarthy et al, *Bioinformatics*, 2017)

PCA on single-cell data typically captures technical effects, most often number of genes detected



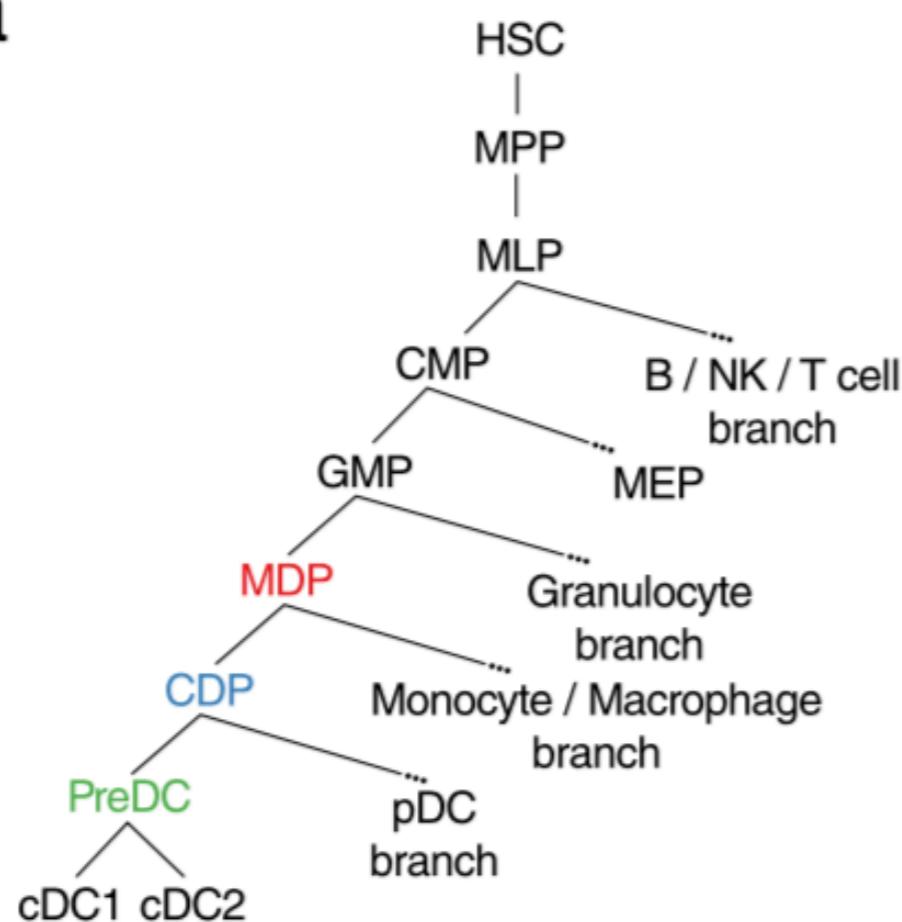
First principal component explains 55% of variance, very strongly correlated with number of genes detected per cell.

PCA extremely useful for QC.

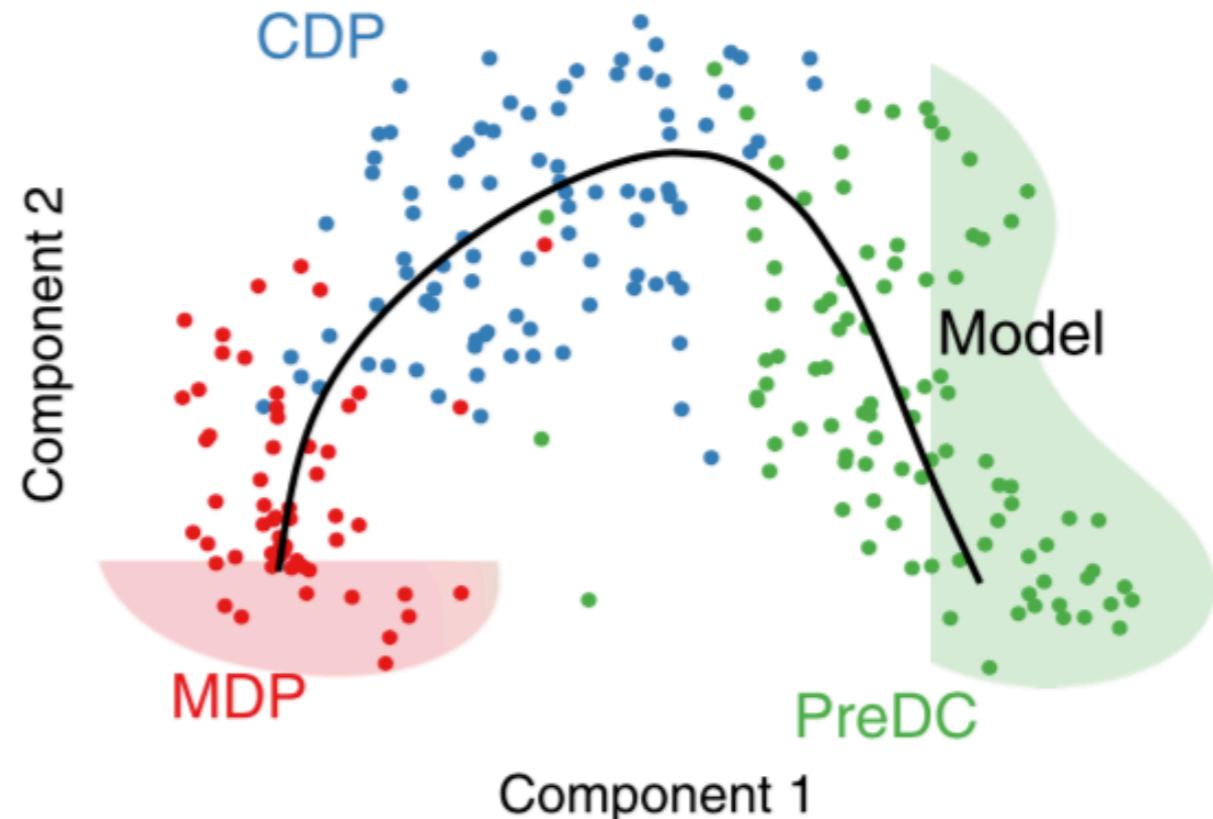
Produced with the R/Bioconductor package scater (McCarthy et al, *Bioinformatics*, 2017)

Trajectory analysis

a



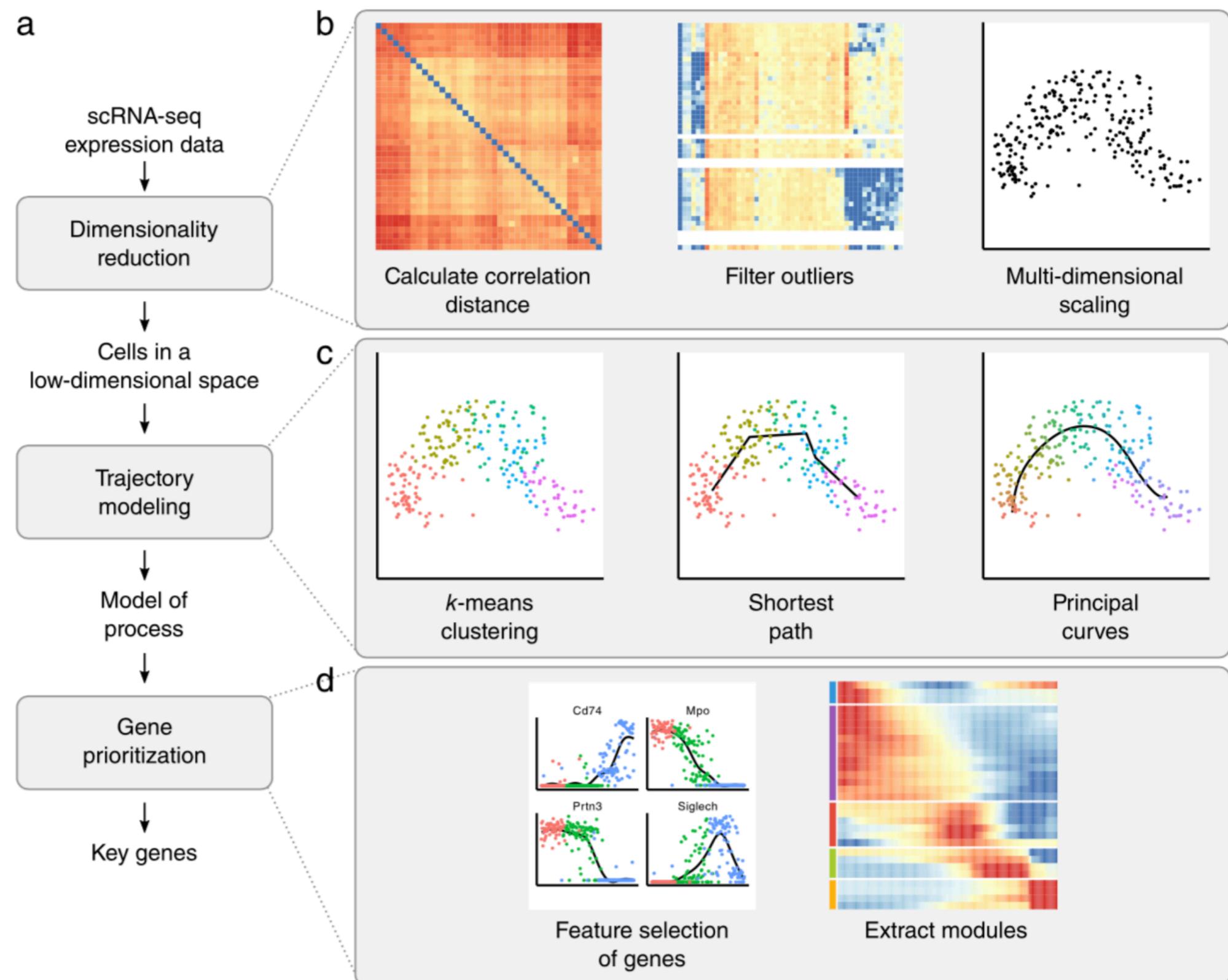
b



Tools:

monocle, SCORPIUS, TSCAN, SLICER, ouija, PhenoPath, GrandPrix (GPLVM)
+ many more

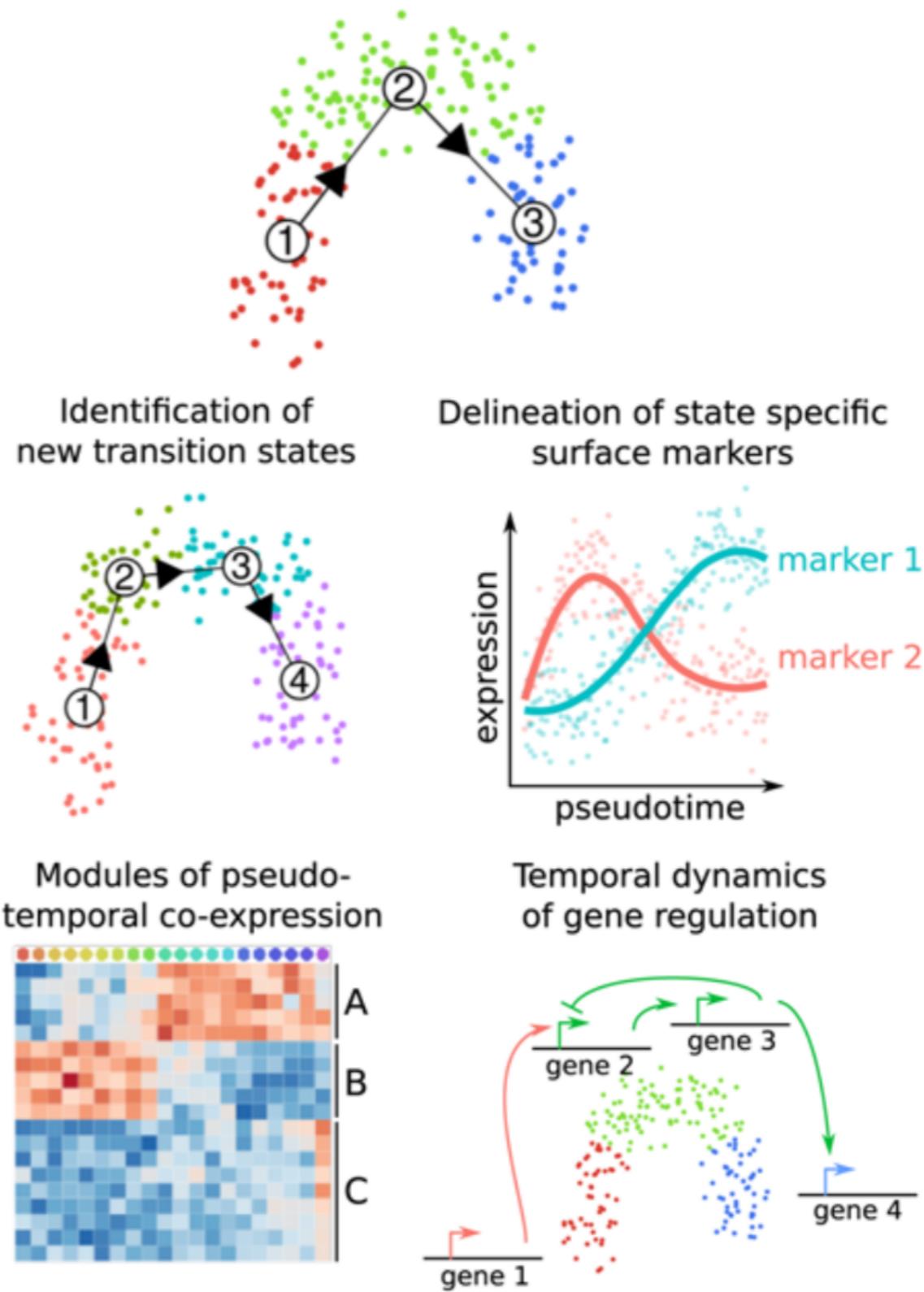
Trajectory analysis



Linear trajectories

B

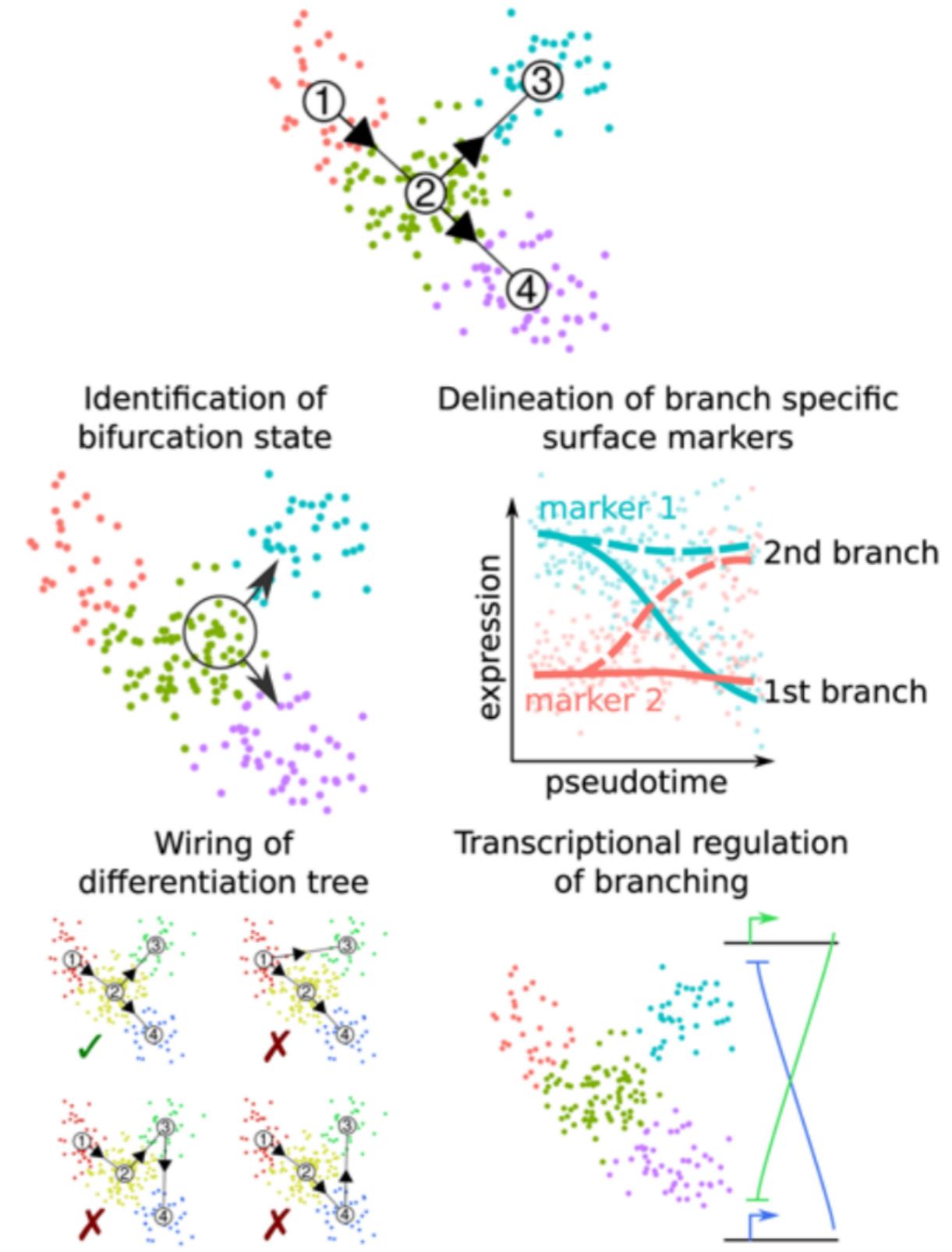
Linear trajectories



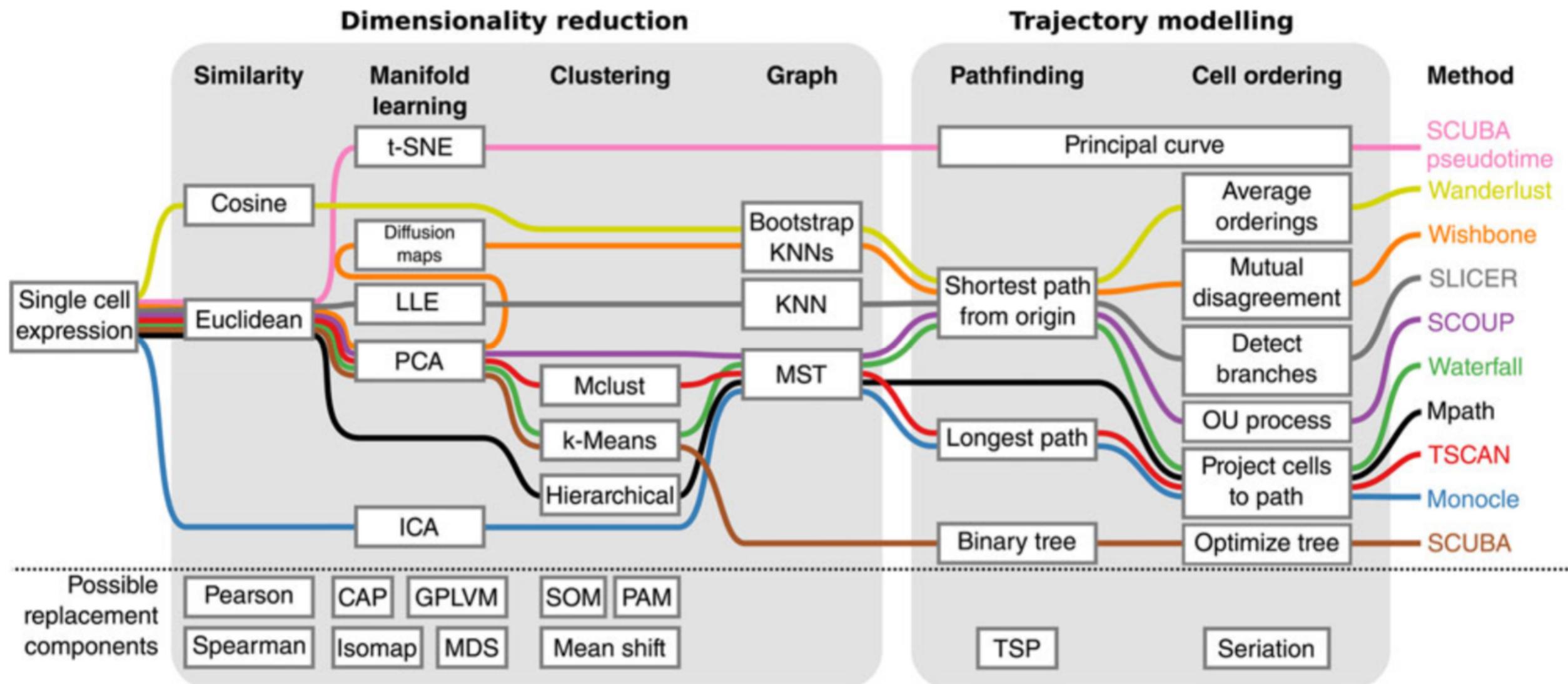
Branched trajectories

C

Branched trajectories



Generalised trajectory inference



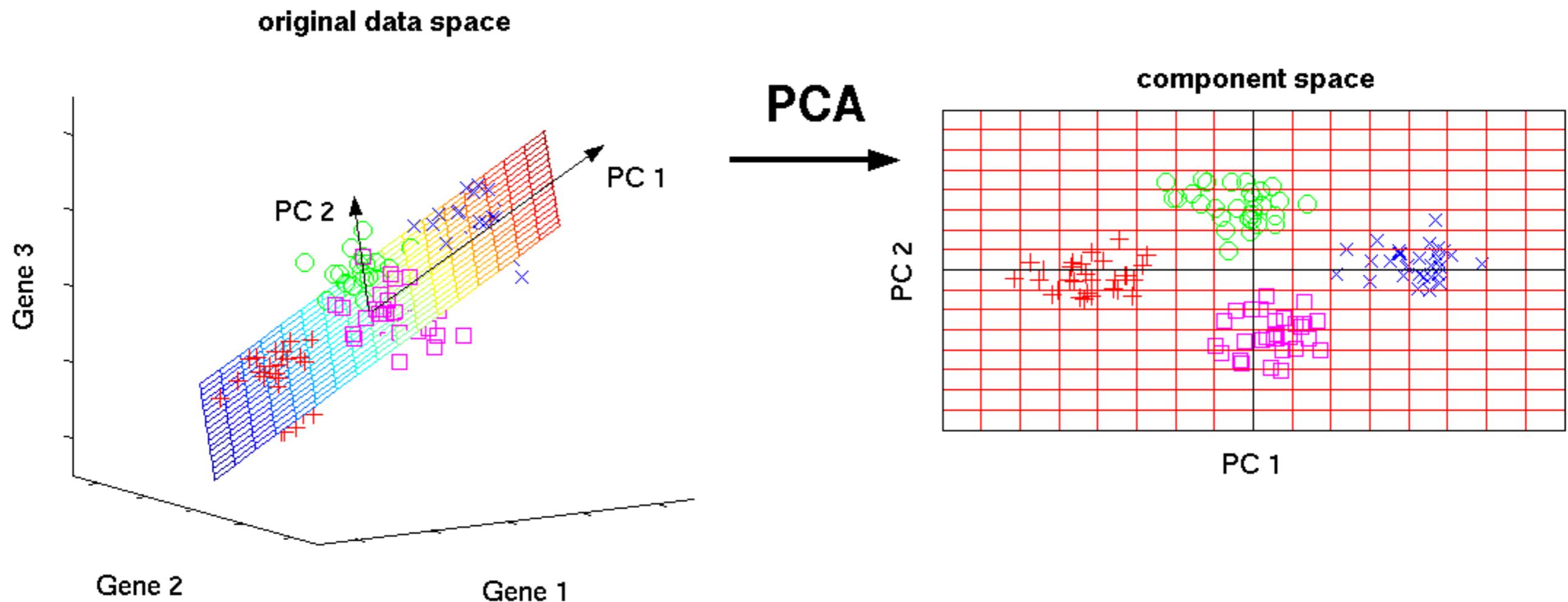
Some trajectory inference methods

| Method | SCUBA pseudotime | Wanderlust | Wishbone | SLICER | SCOUP | Waterfall | Mpath | TSCAN | Monocle | SCUBA |
|--------------------------|---------------------|-------------------------|-------------------------|---------------|---------------------|---------------------|--|---------------------|-------------------------|--------------|
| Visual abstract | | | | | | | | | | |
| 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 | + | + | + | + | - | ± | - | + | + | + |

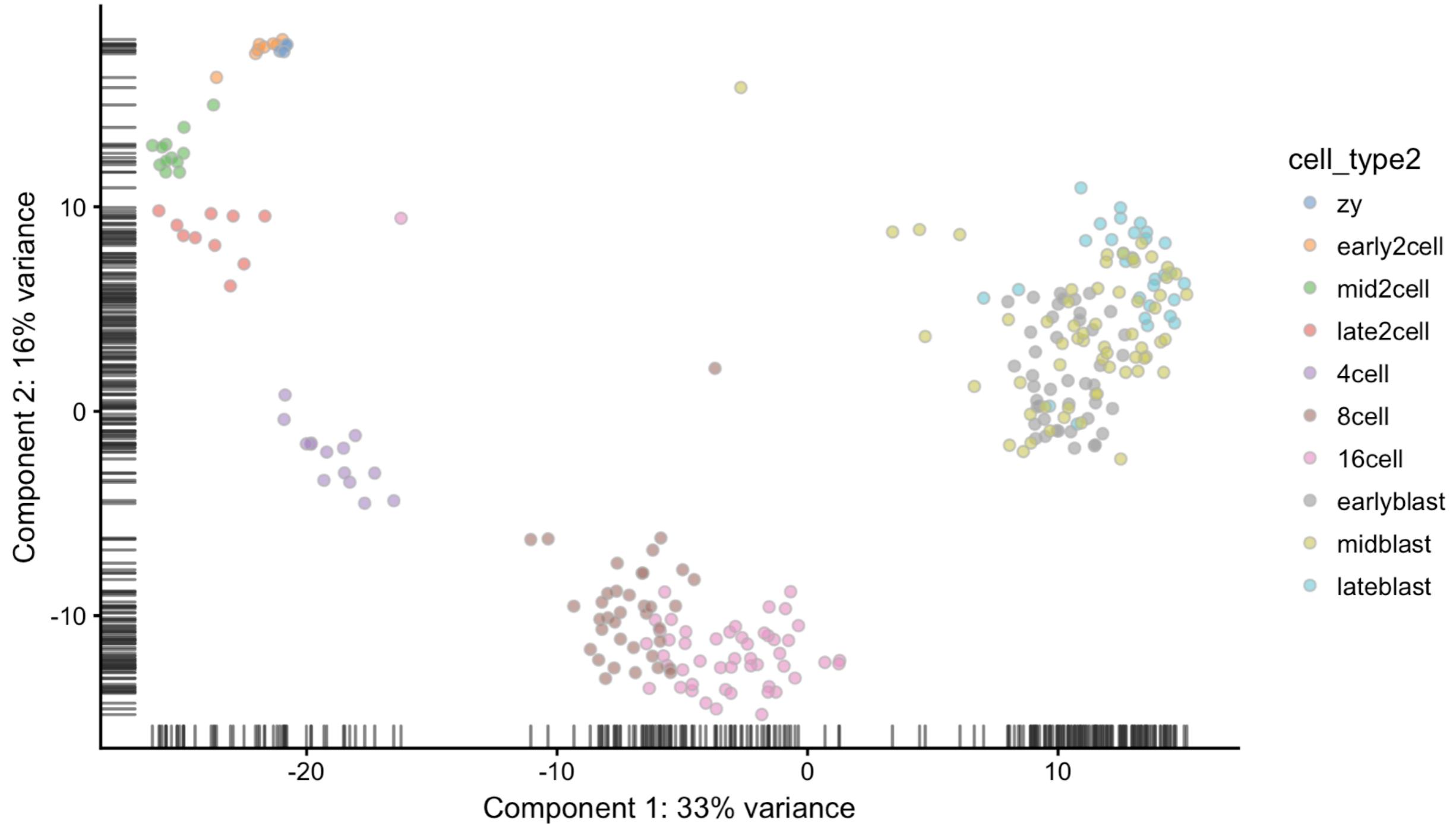
Trajectory inference methods we will use:

- Naive PCA
- Diffusion maps
- Monocle
- TSCAN
- SLICER
- Ouija

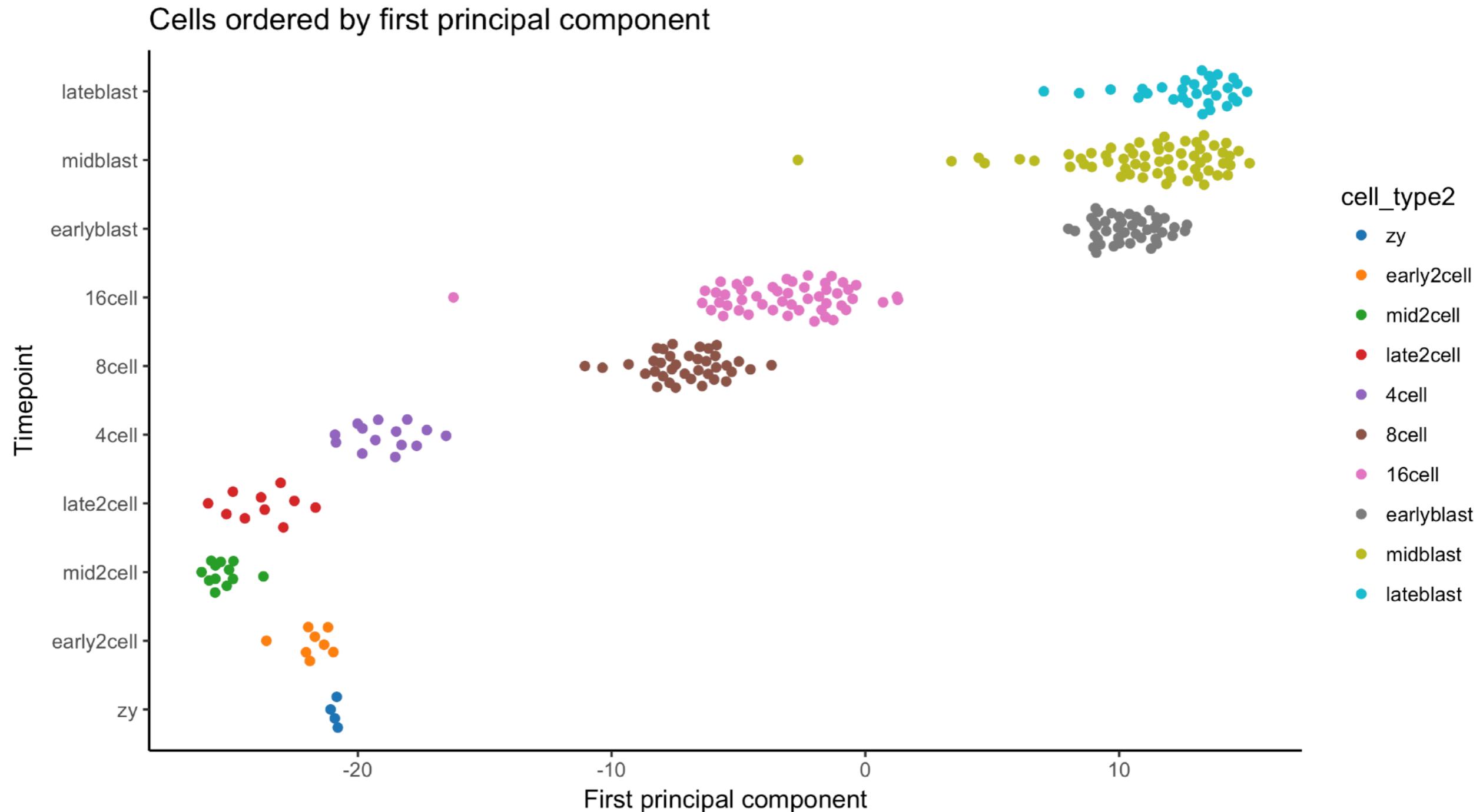
PCA for naive trajectory inference



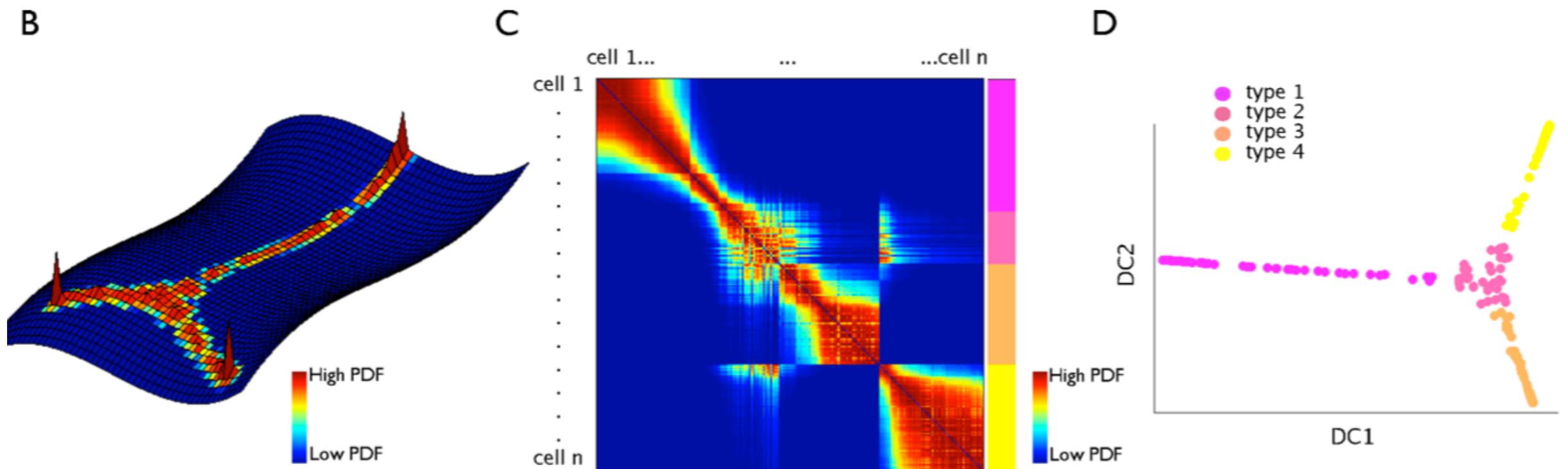
PCA for naive trajectory inference



PCA for naive trajectory inference



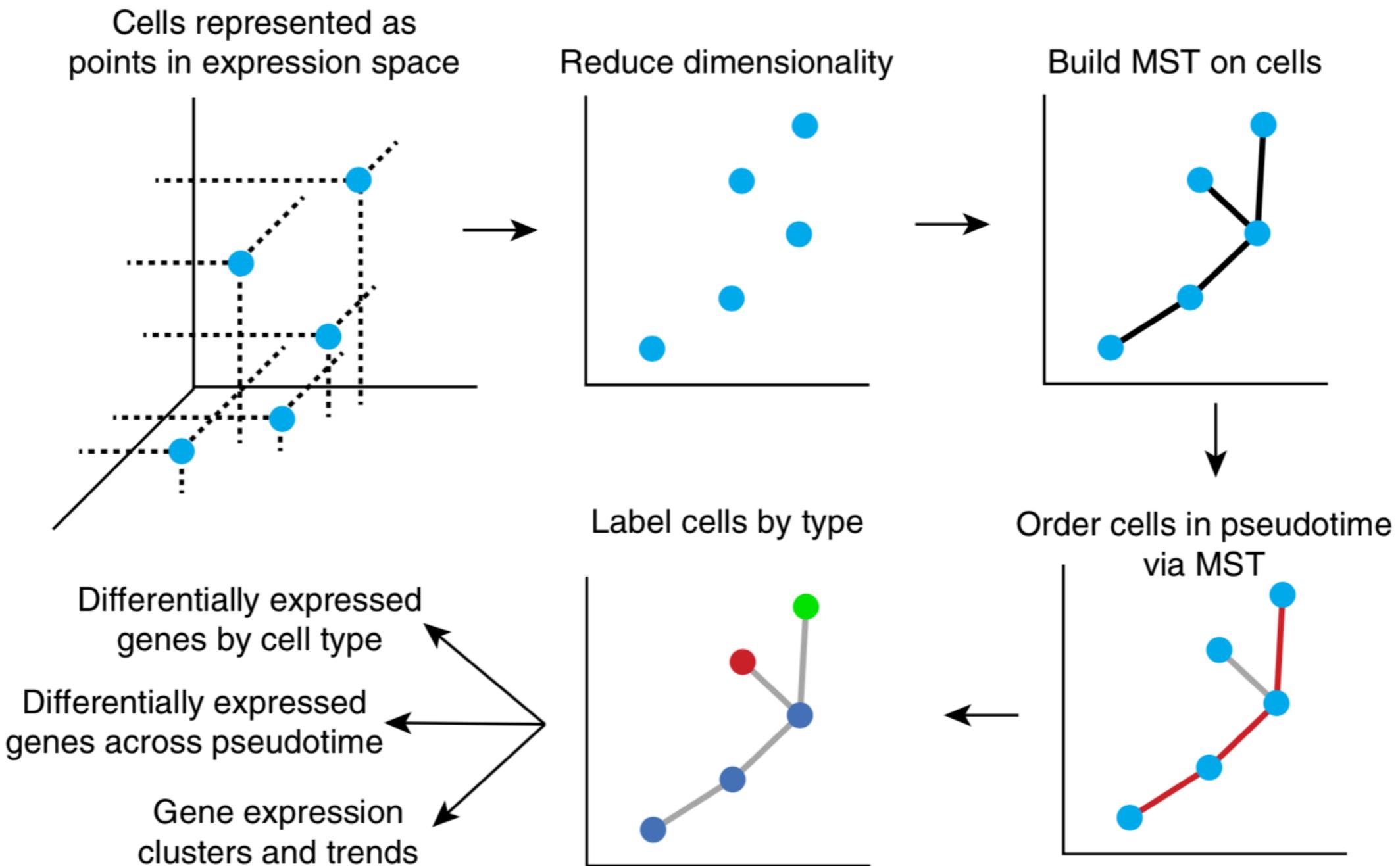
Diffusion Maps



see also “dpt”: diffusion pseudotime

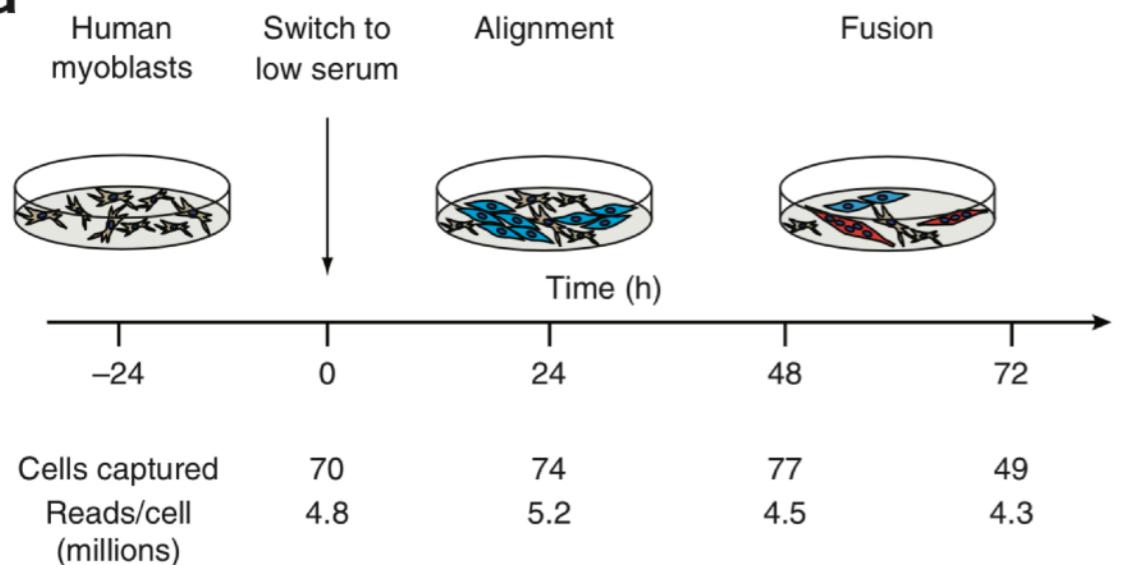
Monocle

a



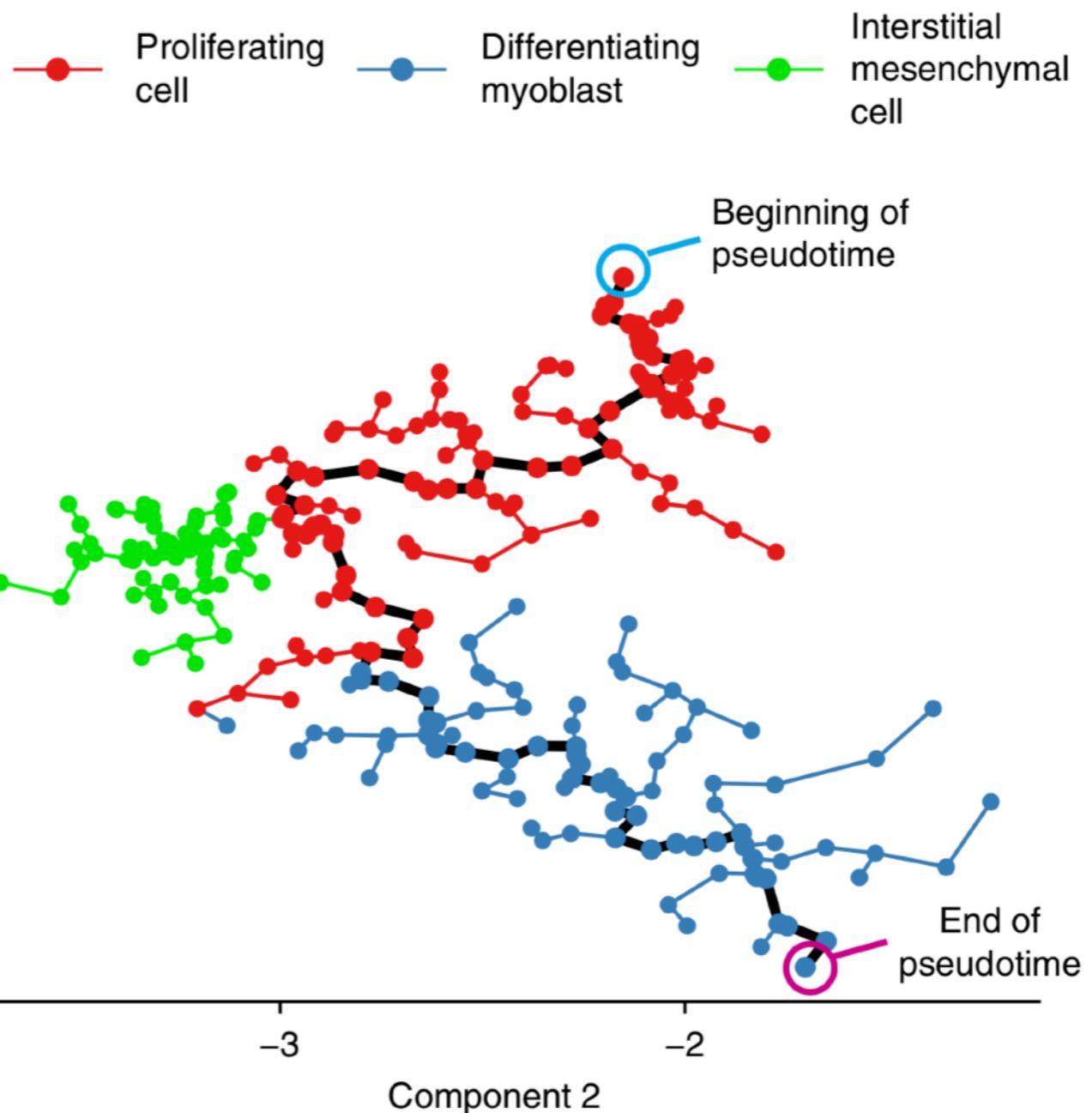
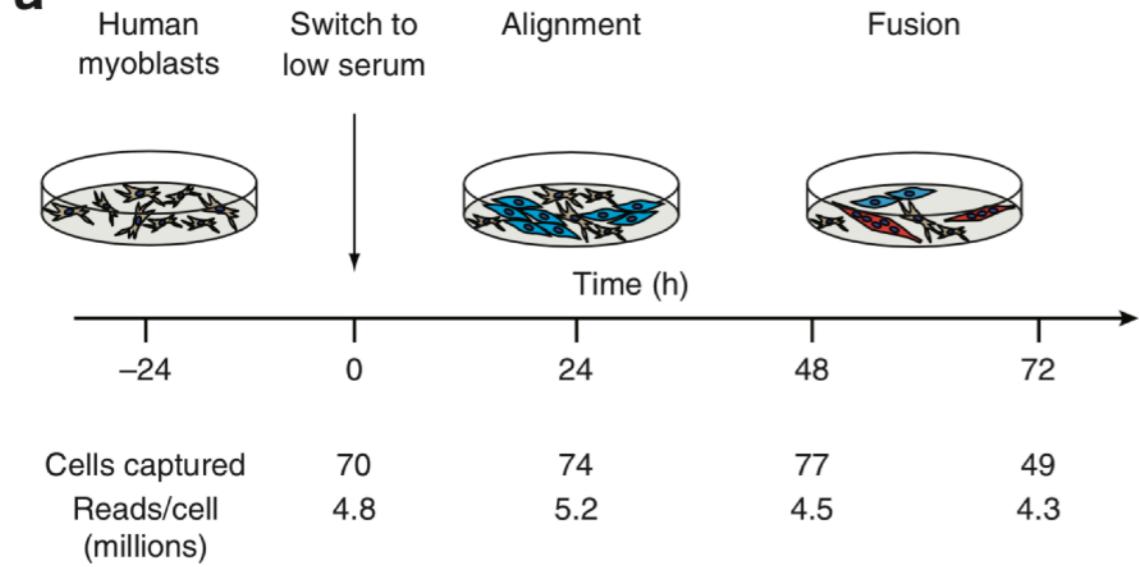
Monocle

a

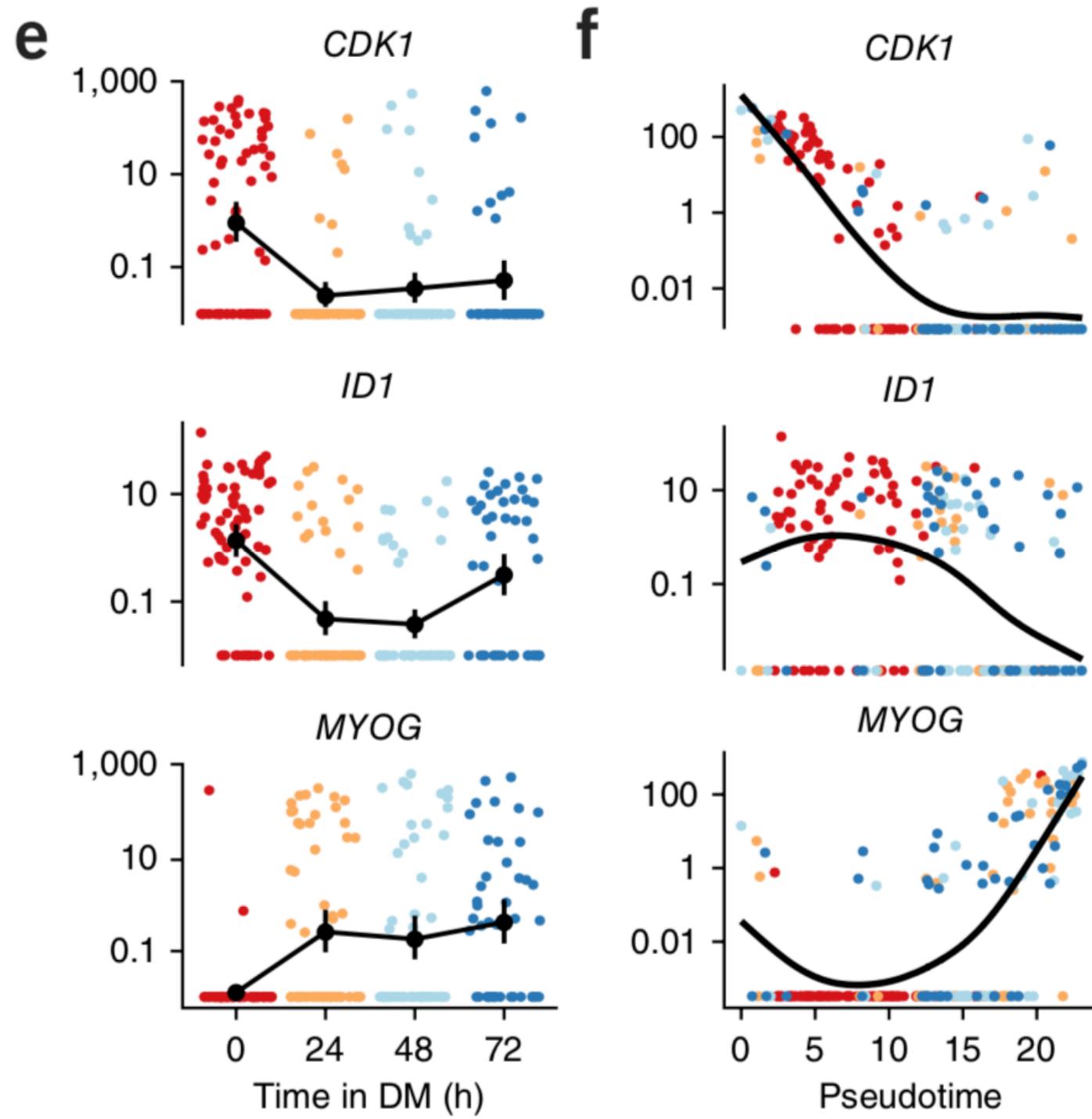


Monocle

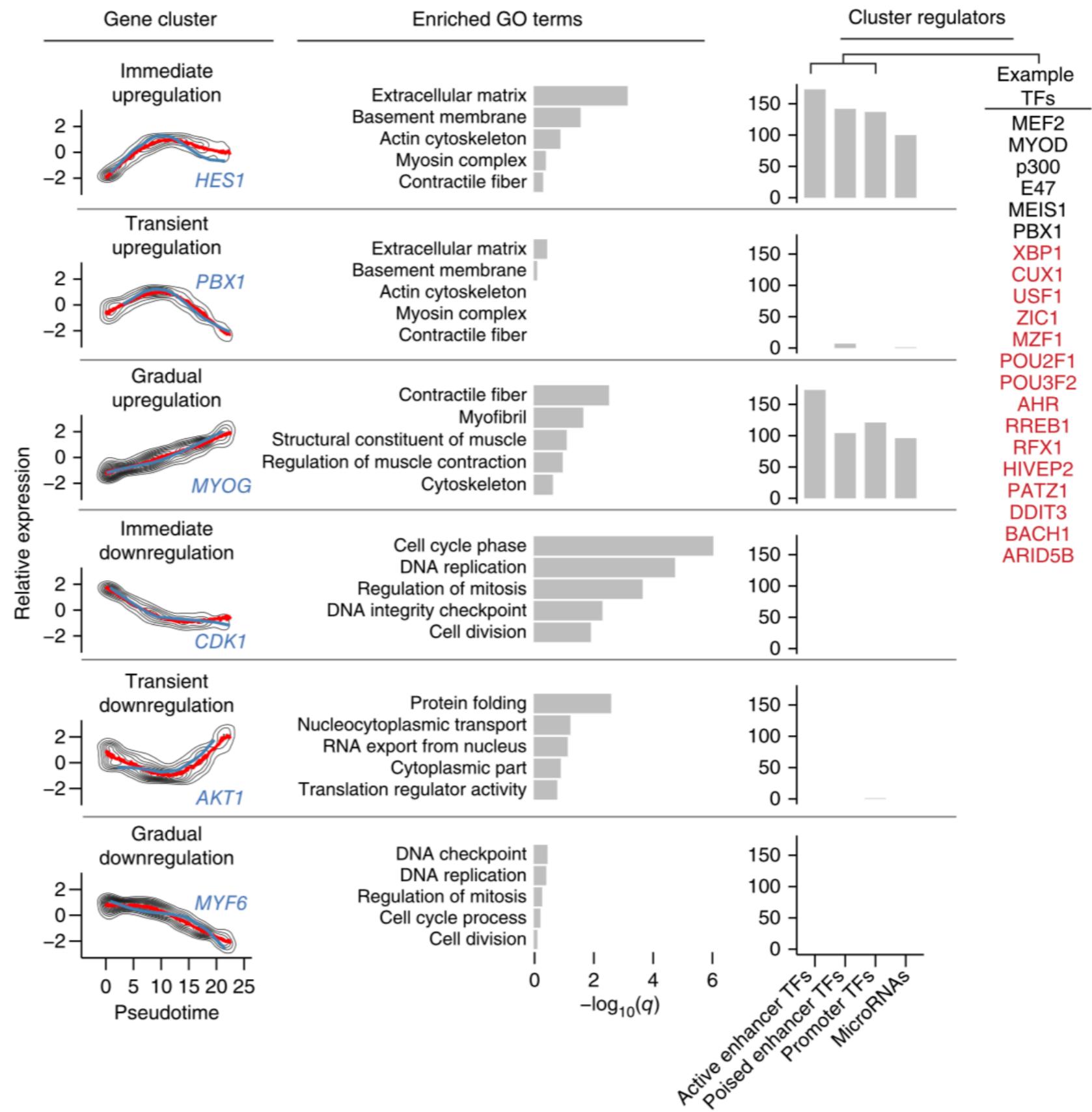
a



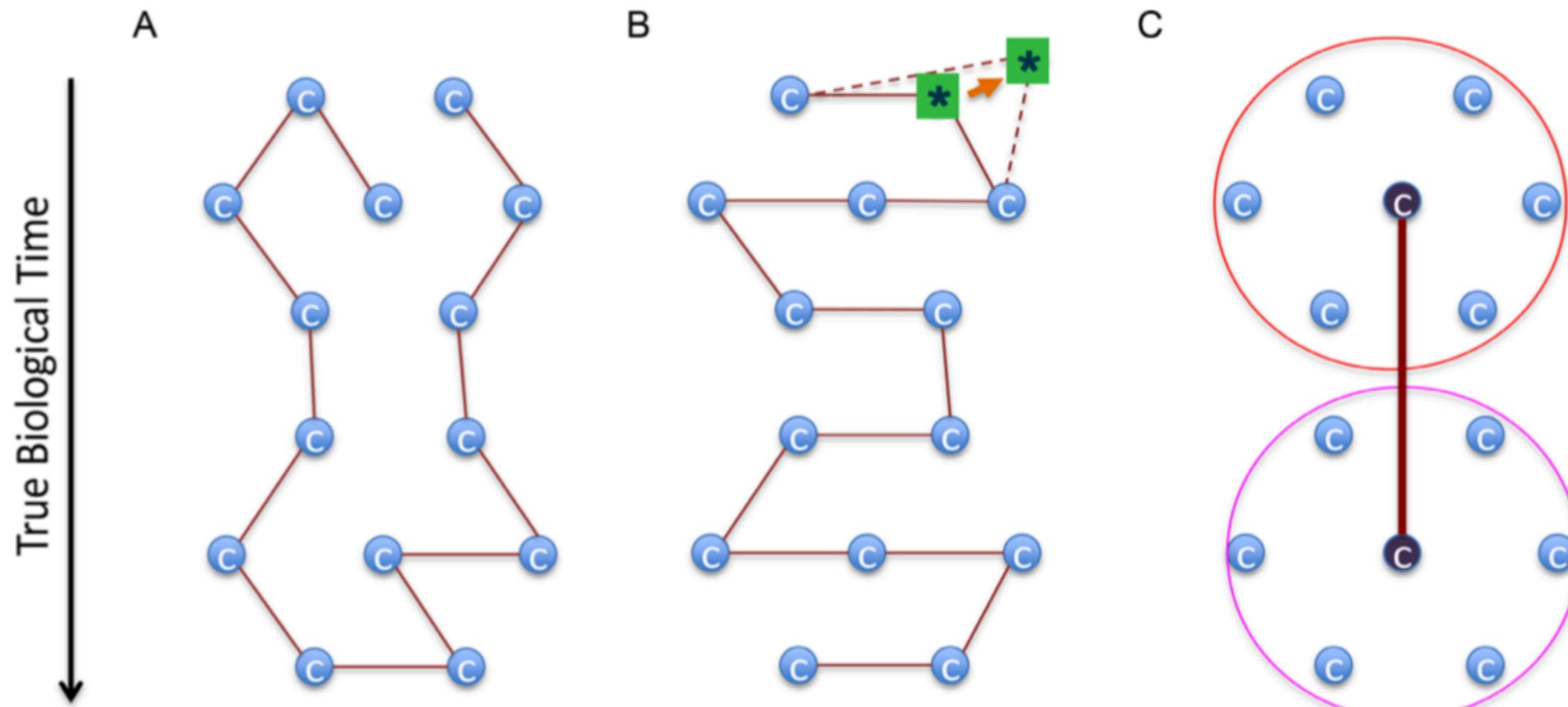
Monocle



Monocle

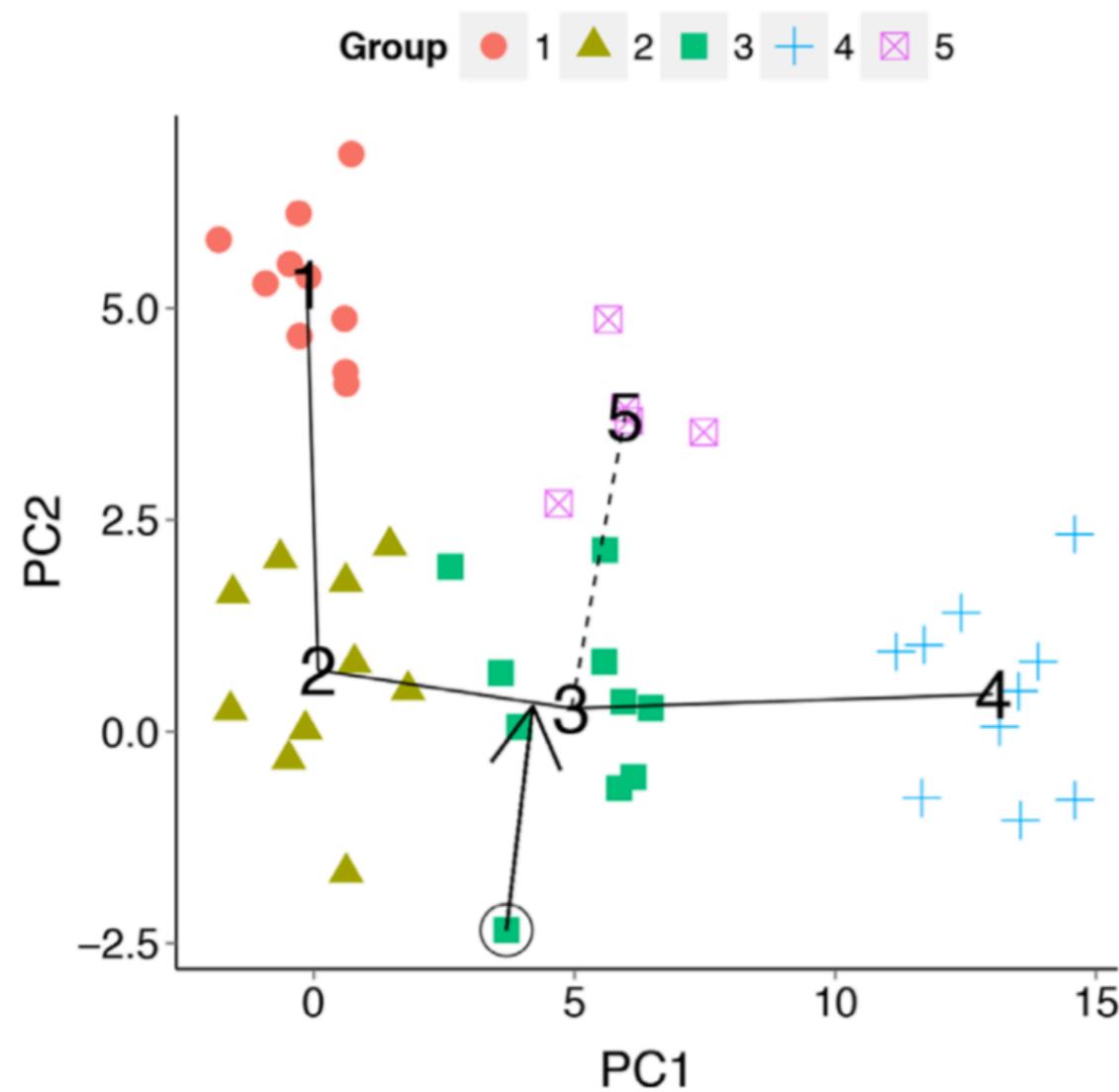


TSCAN

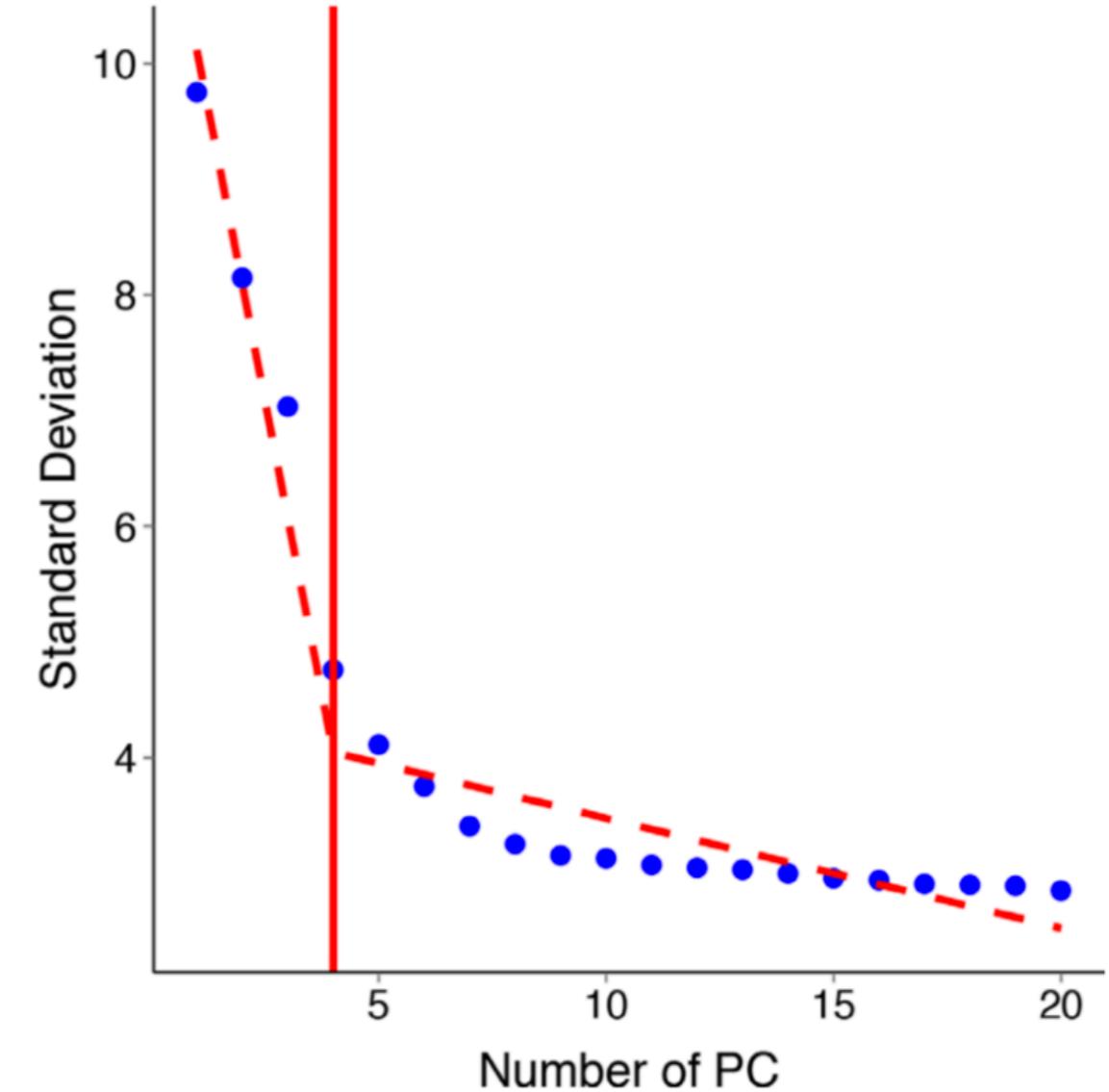


TSCAN

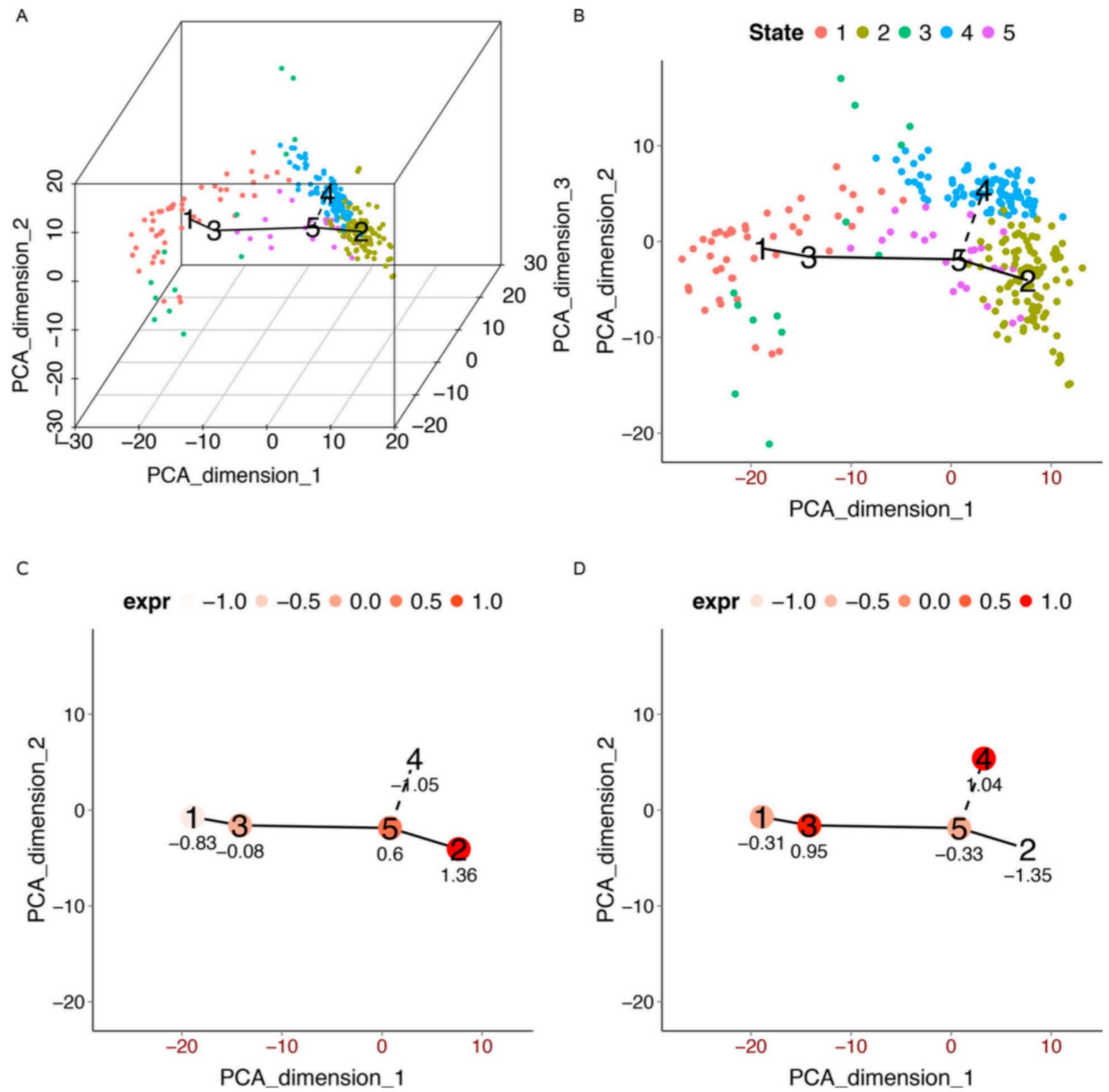
D



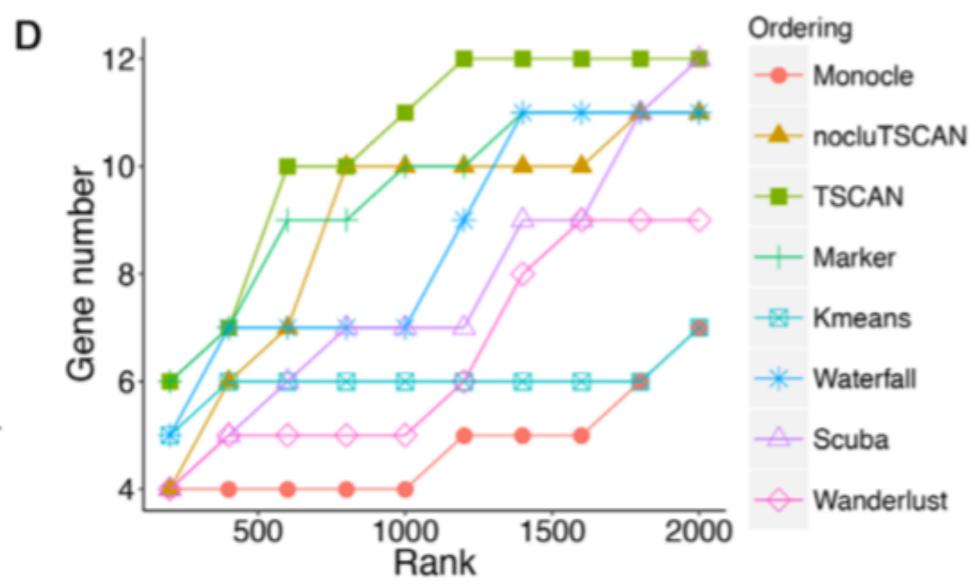
E



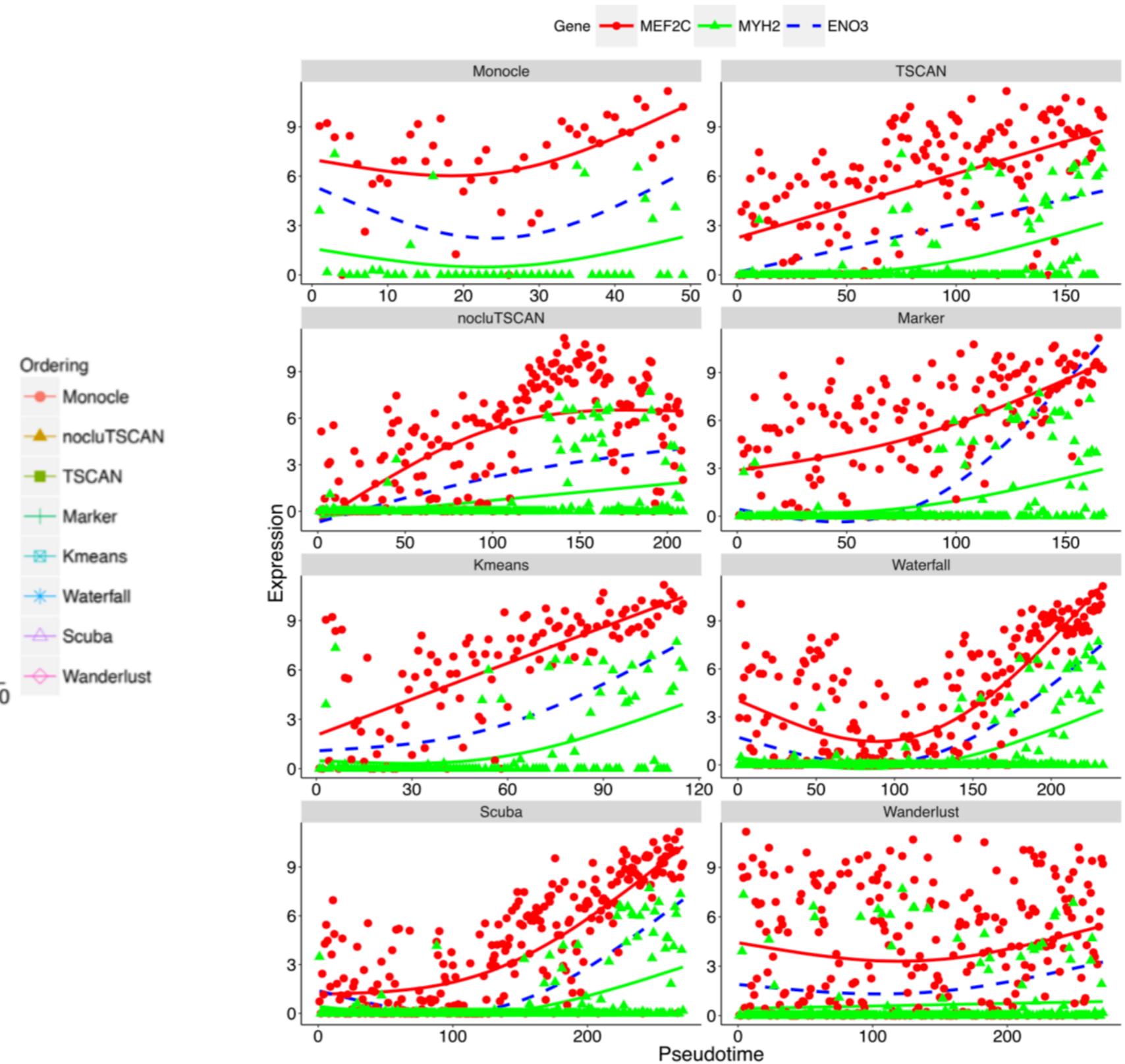
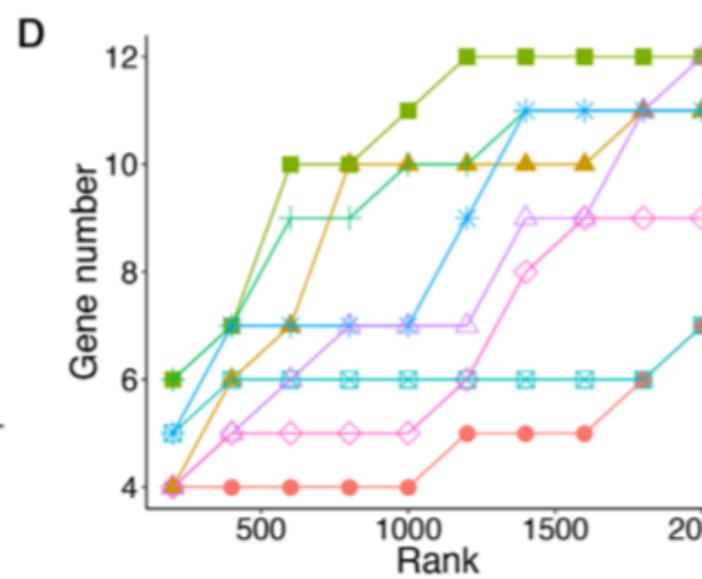
TSCAN



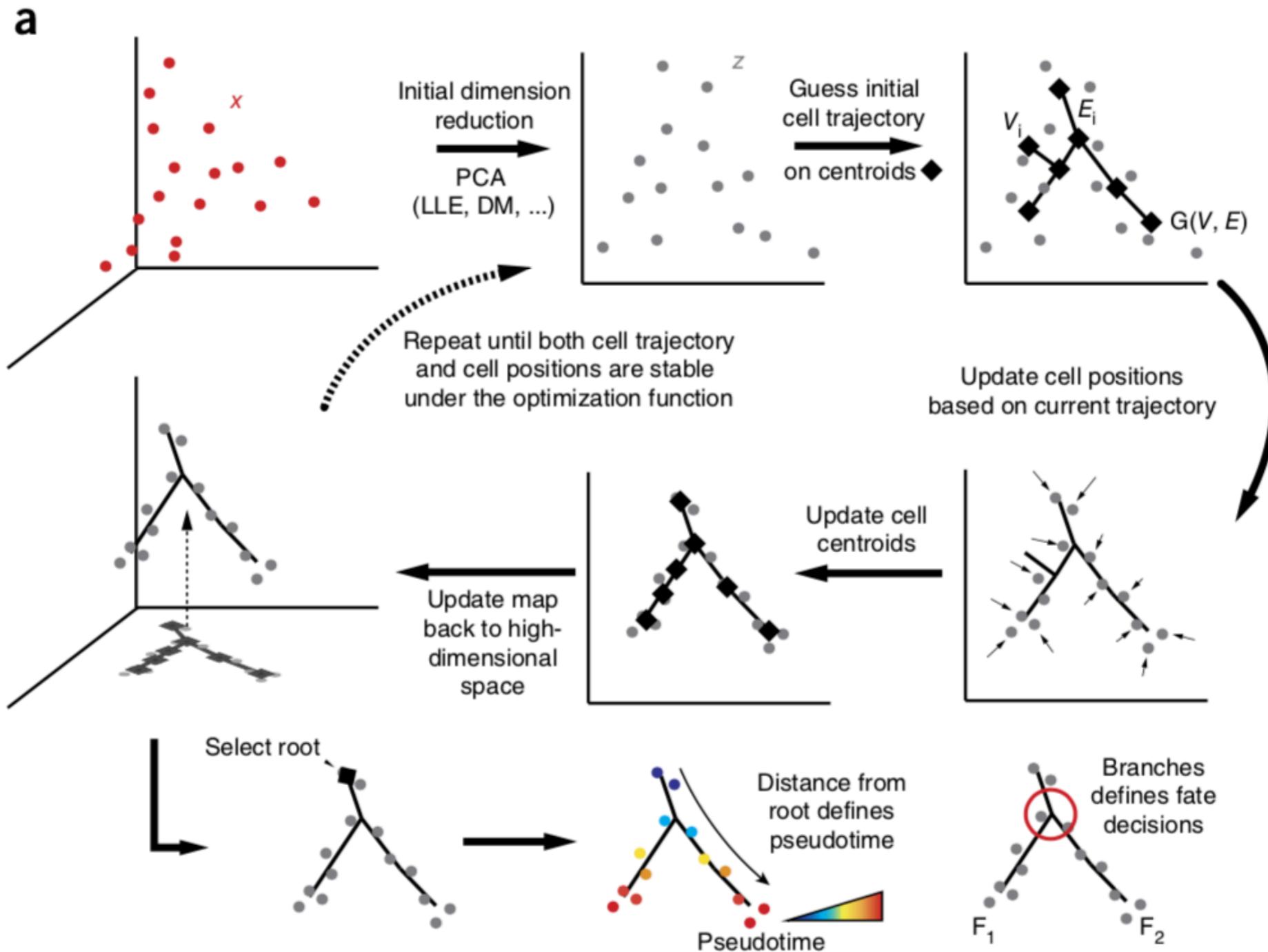
TSCAN



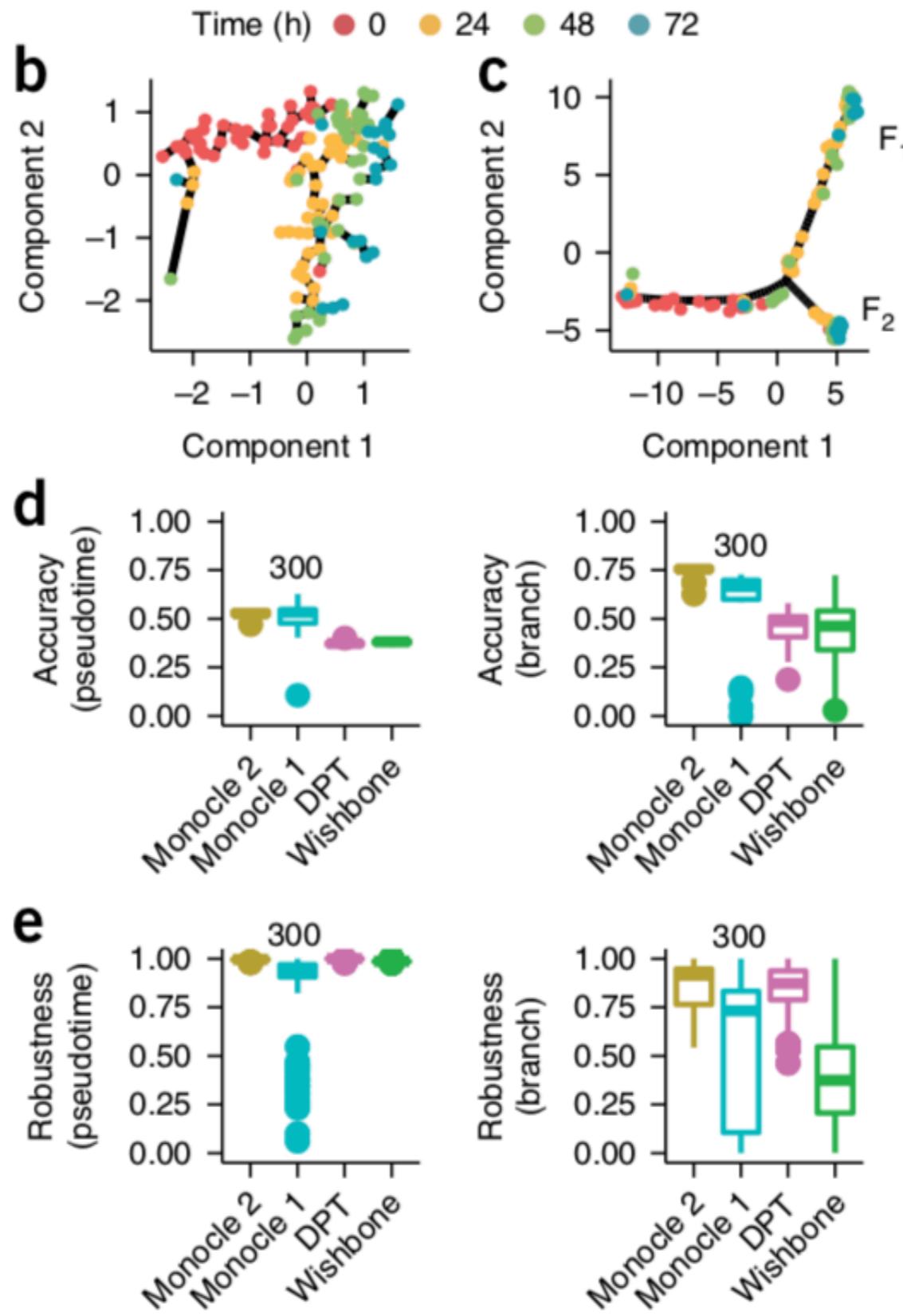
TSCAN



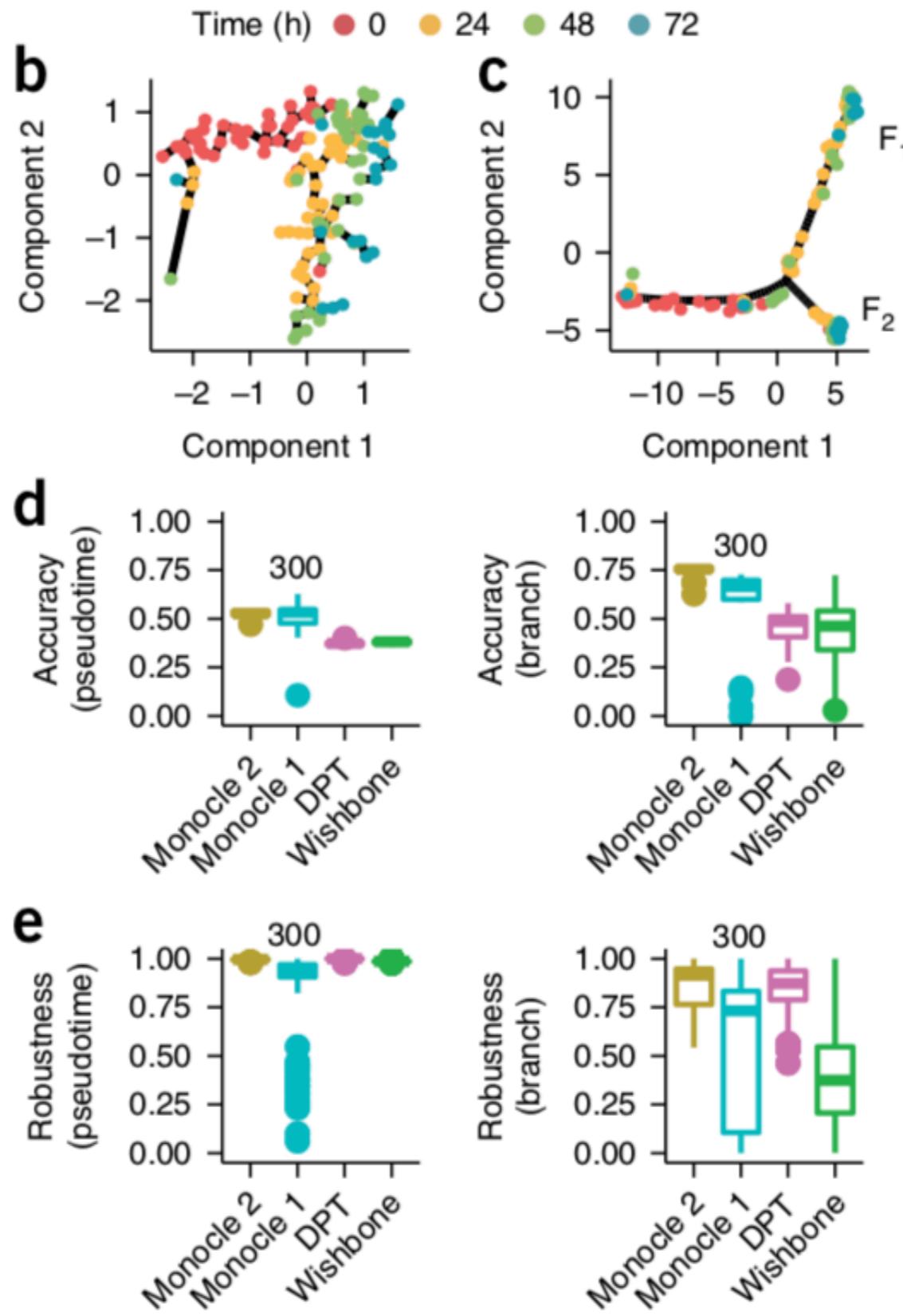
Monocle2



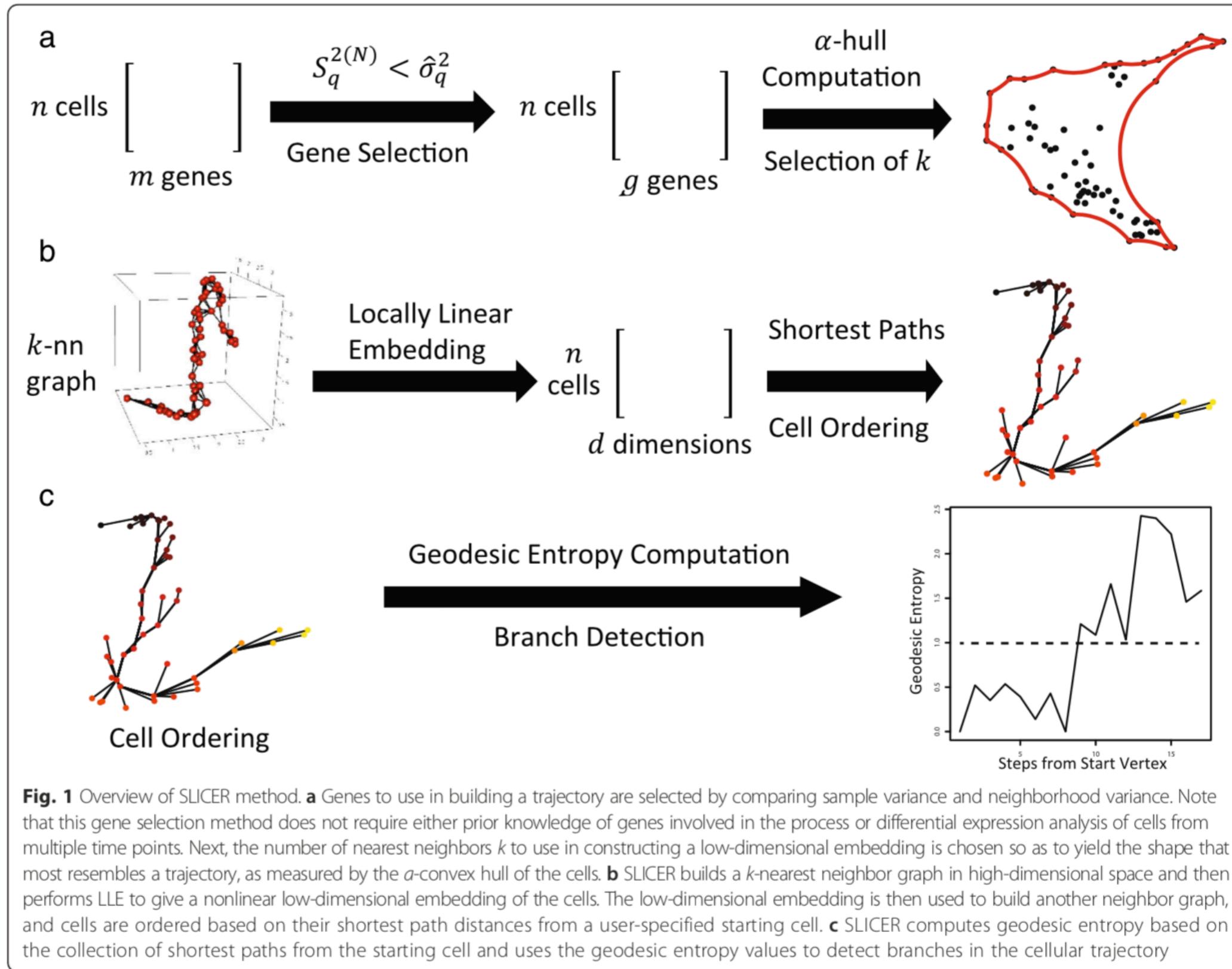
Monocle2



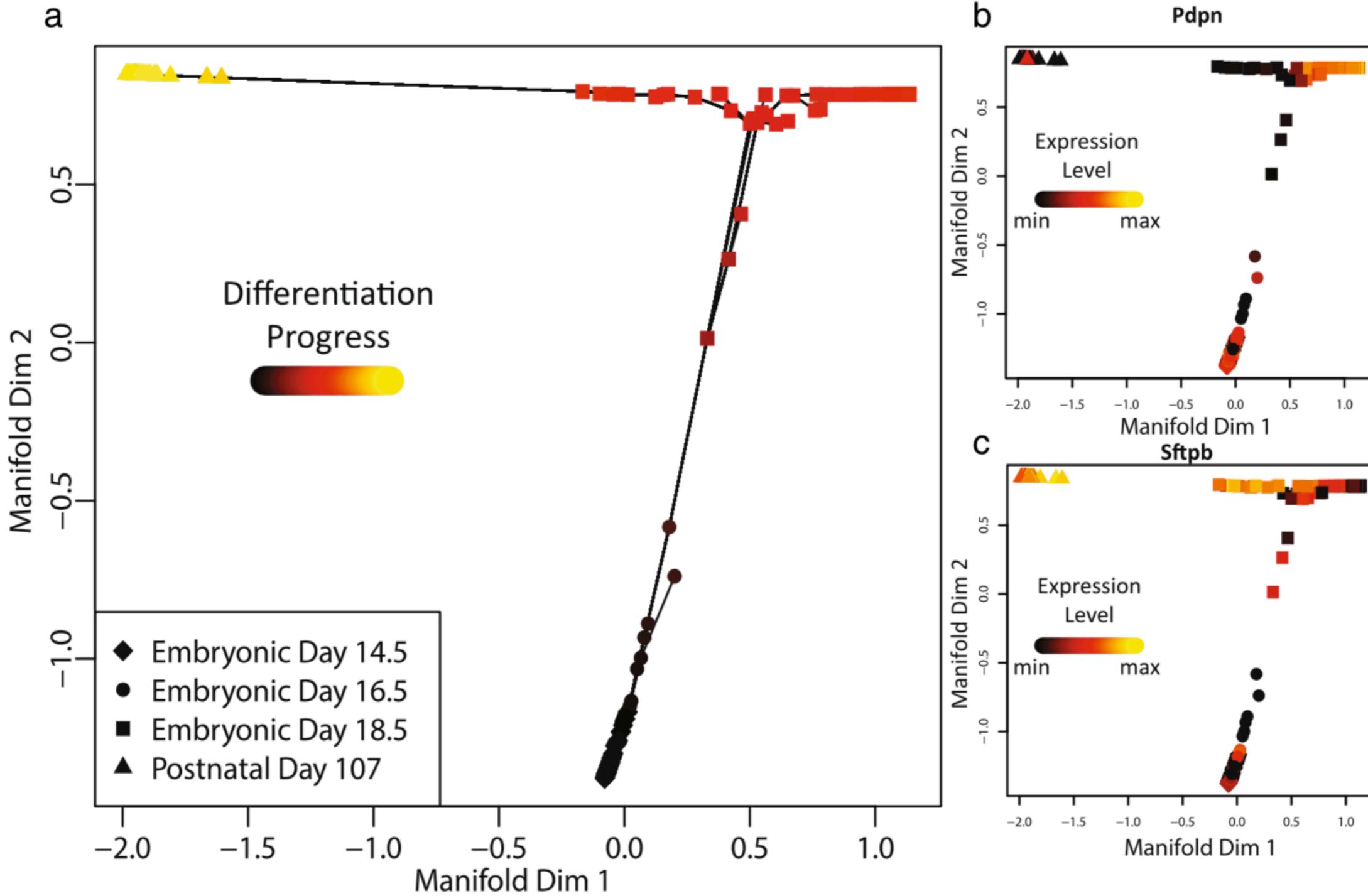
Monocle2



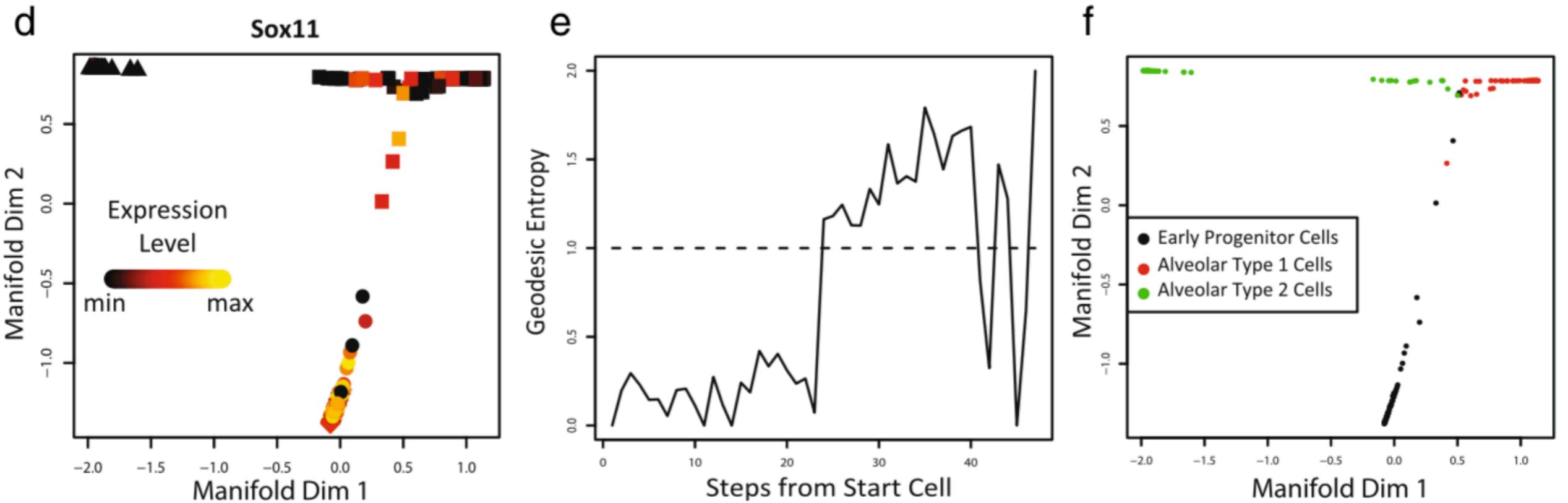
SLICER



SLICER

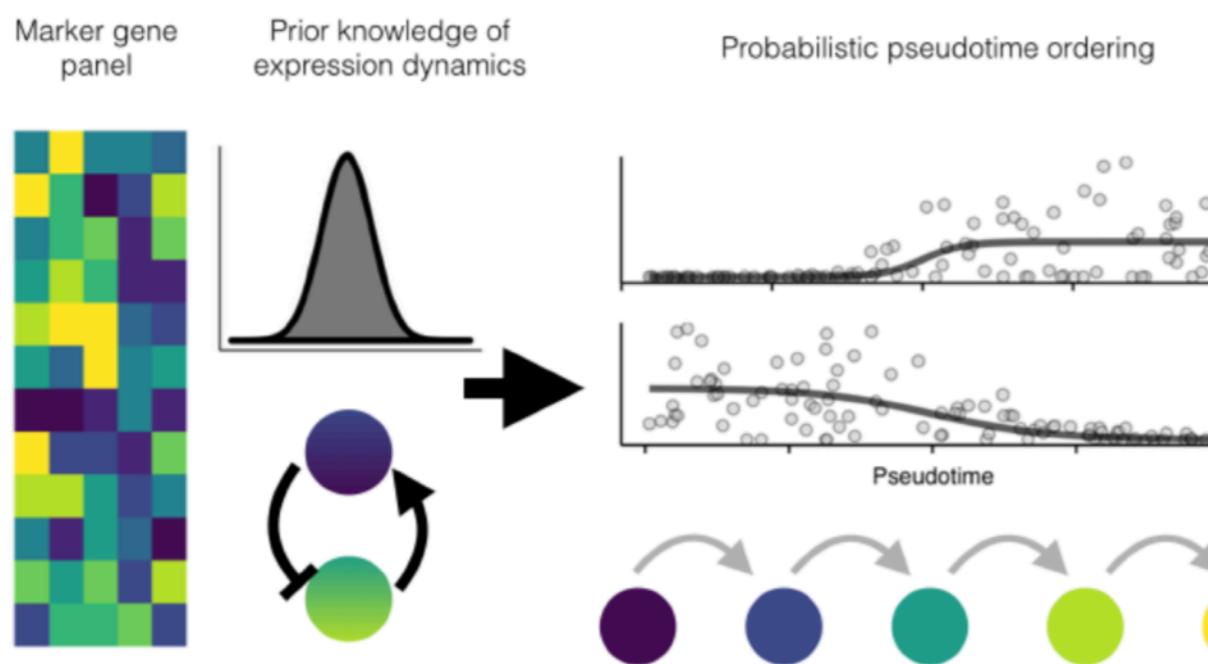


SLICER



Ouija: Using knowledge of marker genes to guide pseudo time inference

A

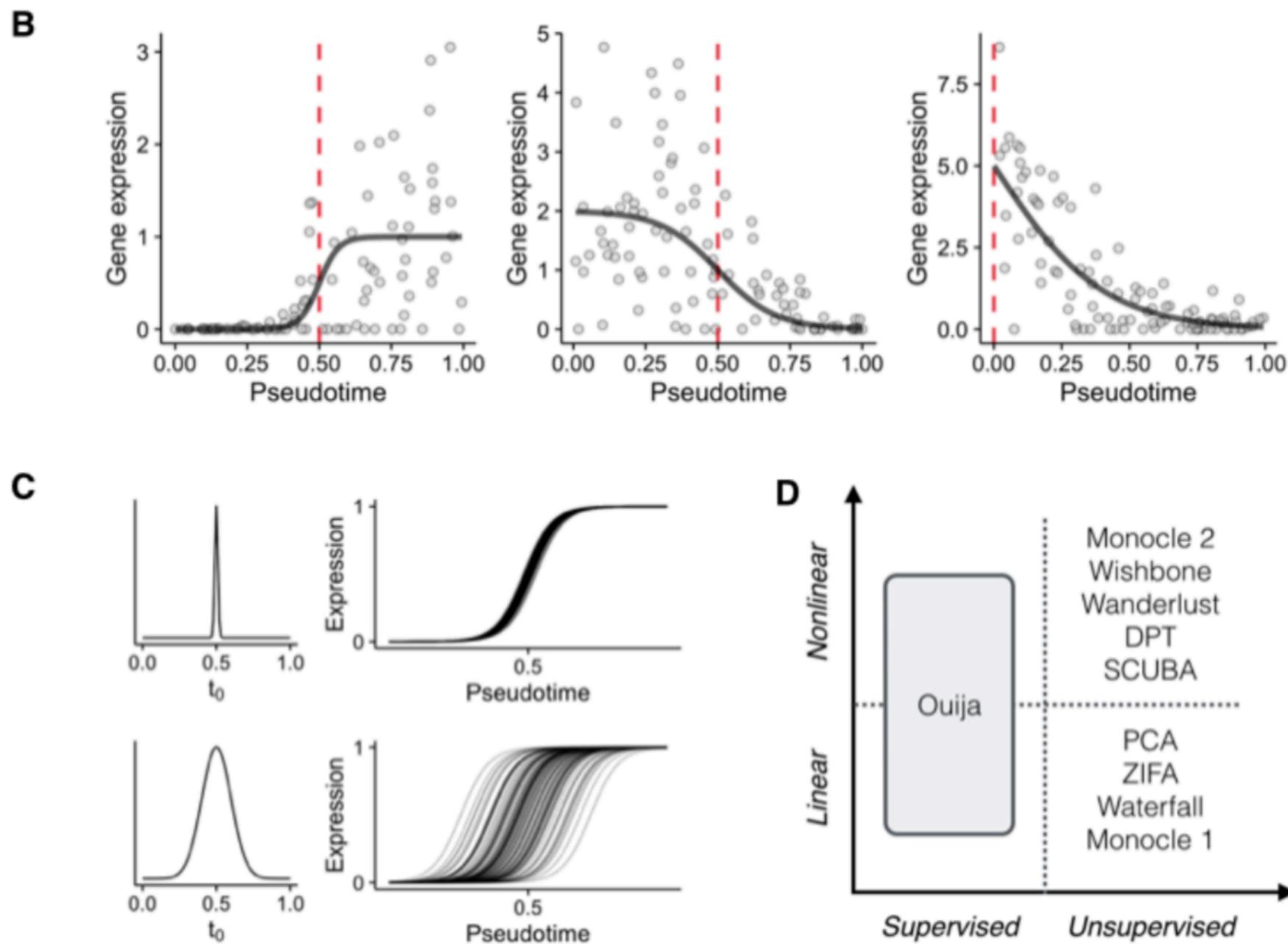


A

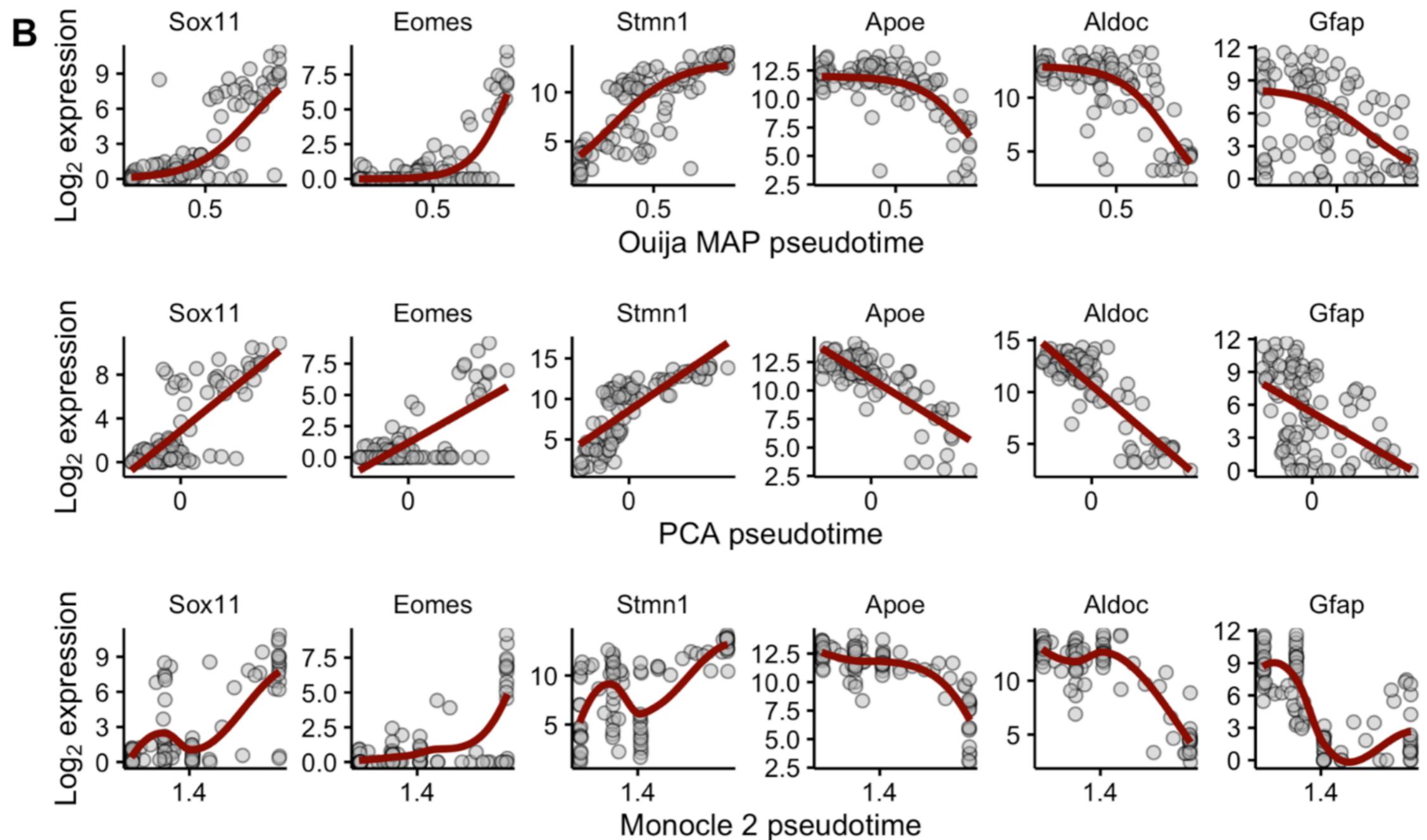
$$\text{Marker gene panel} = f \left(\begin{array}{c} \text{Interpretable mapping} \\ \text{Pseudotimes} \\ \text{Gene behaviour factor matrix} \end{array} \right) + \varepsilon$$

Heteroskedastic noise

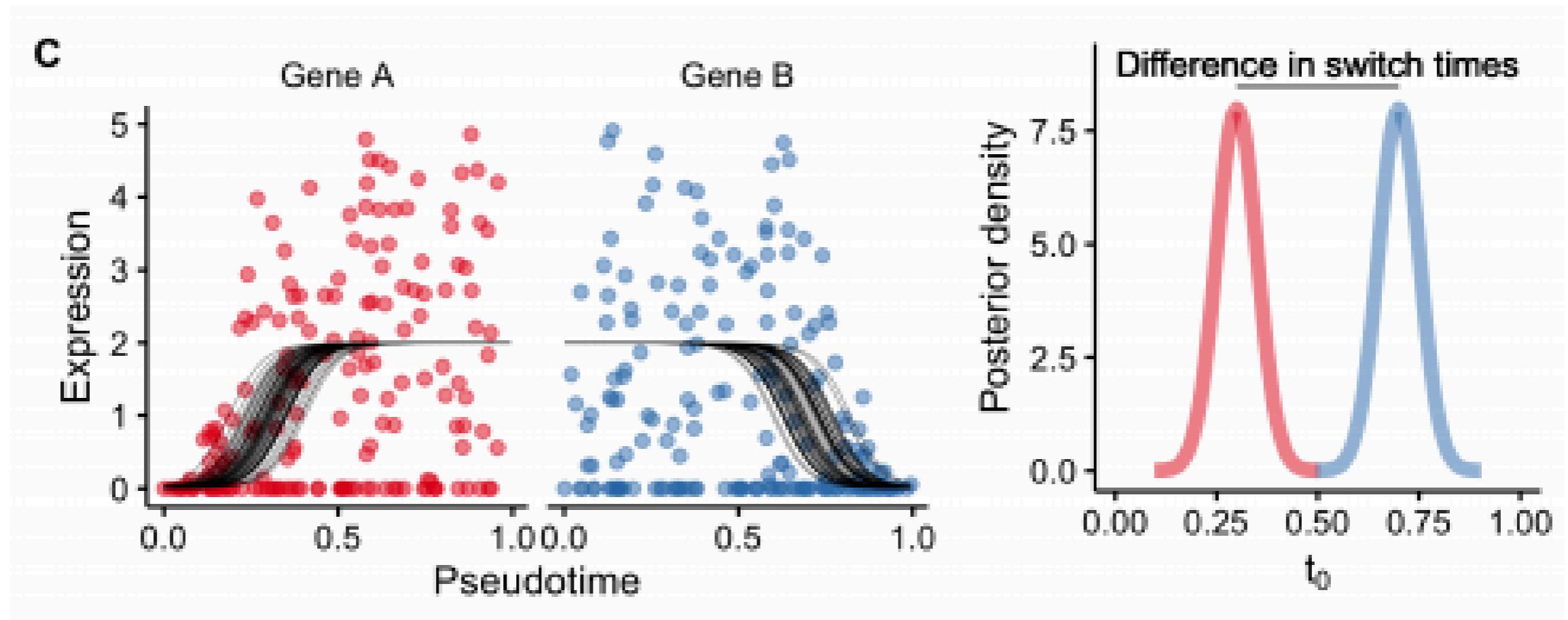
Ouija: Using knowledge of marker genes to guide pseudo time inference



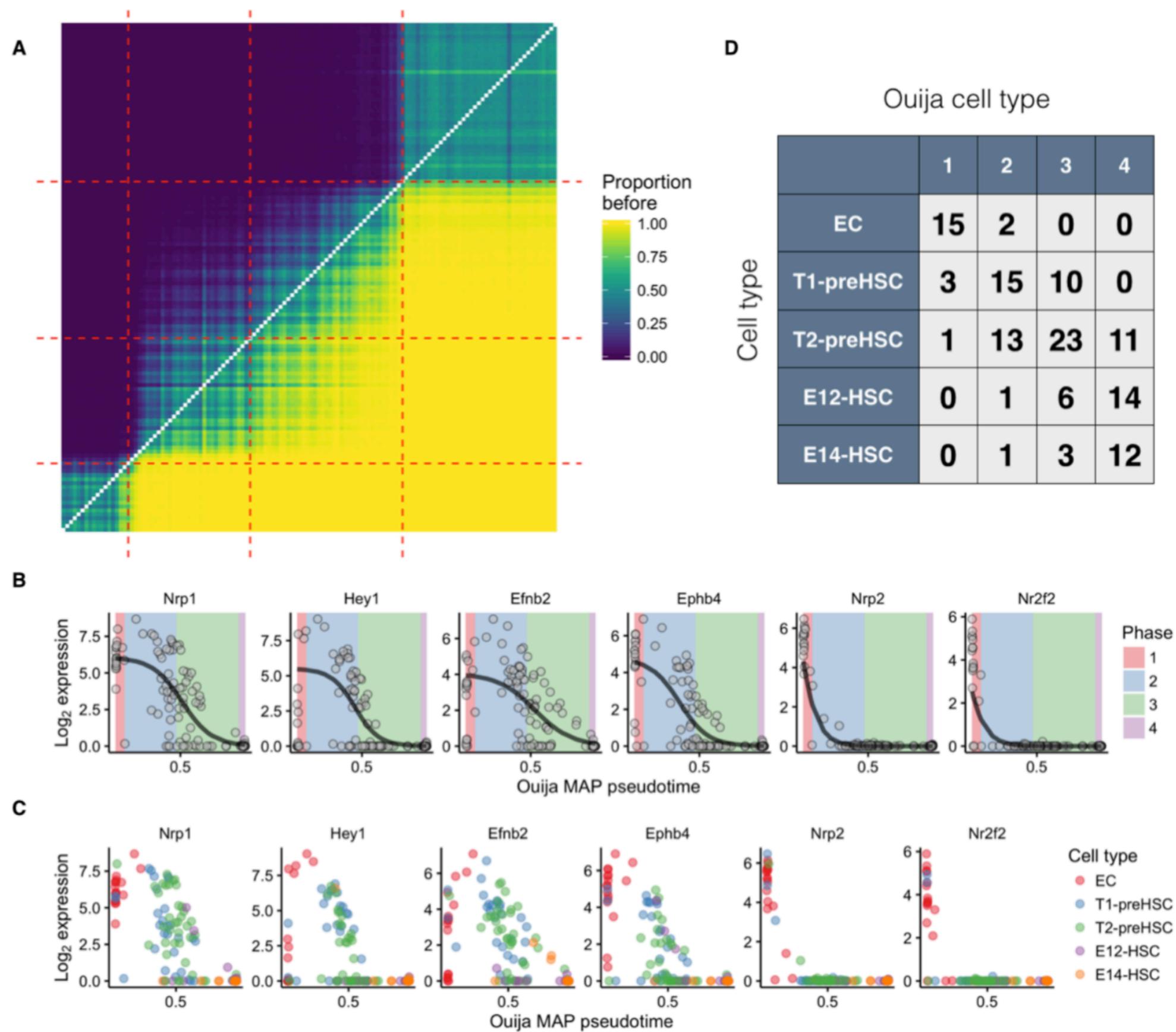
Ouija vs PCA vs Monocle on data from neurogenesis in mice



Ouija: difference in switch times (gene regulation)



Ouija: metastable states



Trajectory inference methods we have looked at:

- Naive PCA
- Diffusion maps
- Monocle
- TSCAN
- SLICER
- Ouija

Robust pseudotime analysis

- Use multiple approaches
- Consider bootstrapping
- Use orthogonal data (cell markers, etc)
- Use prior biological knowledge
- Be sceptical

Gotchas

- Trajectory methods will return a trajectory, even where none is present
 - Compare with clustering methods (e.g. flowSOM)
- Particular assumptions for a method may introduce bias that prevent the method from finding the correct but unexpected trajectory in the data

Visualisation

Acknowledgements

Stegle Group @ EMBL-EBI:

Oliver Stegle, Florian Buettner, Anna Cuomo, Raghd Rostom, Marc Jan Bonder, Yuanhua Huang, Ricard Argelaguet + many more great colleagues

scater authors:

Aaron Lun, Kieran Campbell, Quin Wills

Hemberg Group @ Sanger Institute:

Vlad Kiselev, Tallulah Andrews, Martin Hemberg

Teichmann Group @ Sanger Institute:

Sarah Teichmann, Valentine Svensson + more

Bioconductor and open-source software developers



Australian Government

National Health and Medical Research Council

EMBL-EBI



Useful resources

- **scater:**
McCarthy et al. *Bioinformatics*. 2017; doi:10.1093/bioinformatics/btw777
<http://bioconductor.org/packages/scater/>
- F1000 Research "A **step-by-step workflow** for low-level analysis of single-cell RNA-seq data":
<http://f1000research.com/articles/5-2122/v1>
- **Hemberg Lab scRNA-seq course:**
<http://hemberg-lab.github.io/scRNA.seq.course/>
- **conquer:** a repository for processed, QC'd single-cell datasets.
<http://imlspenticton.uzh.ch:3838/conquer/>
- **scRNA tools database:** <https://www.scrna-tools.org/>
- **(Long!) list** of single-cell tools and software:
<https://github.com/seandavi/awesome-single-cell>