

Imputation of scRNA-seq data

David van Dijk

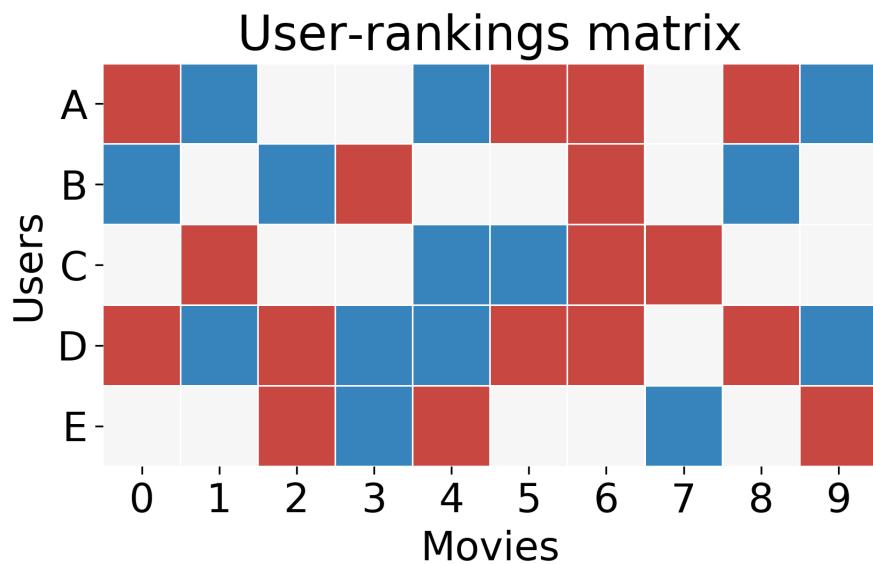
Internal Medicine & Computer Science @ Yale

vandijklab.org

What do we do with missing data?

The Netflix Problem

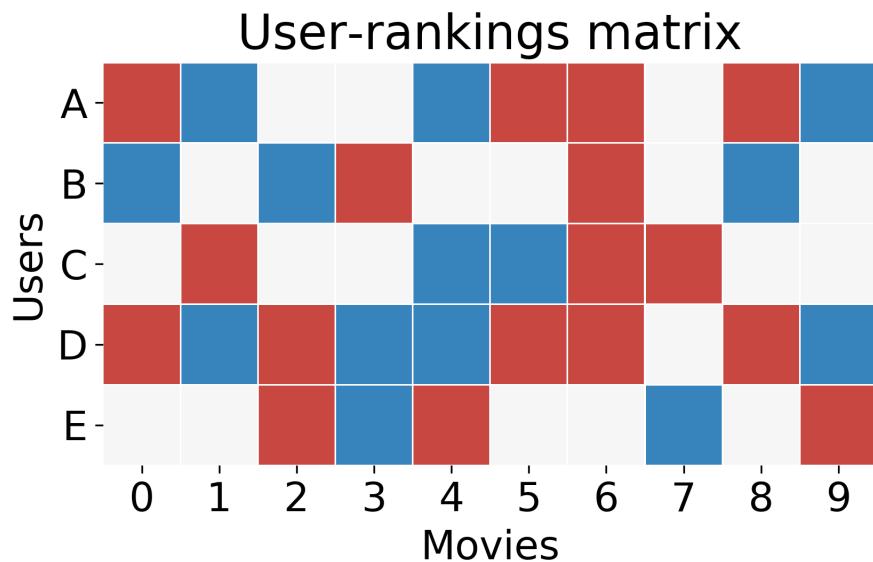
- Rows are users
- Columns are “+1” for like and “-1” for dislike
- We know which values are missing (unwatched)
- Goal: predict missing values to suggest movies to each user



What do we do with missing data?

The Netflix Problem

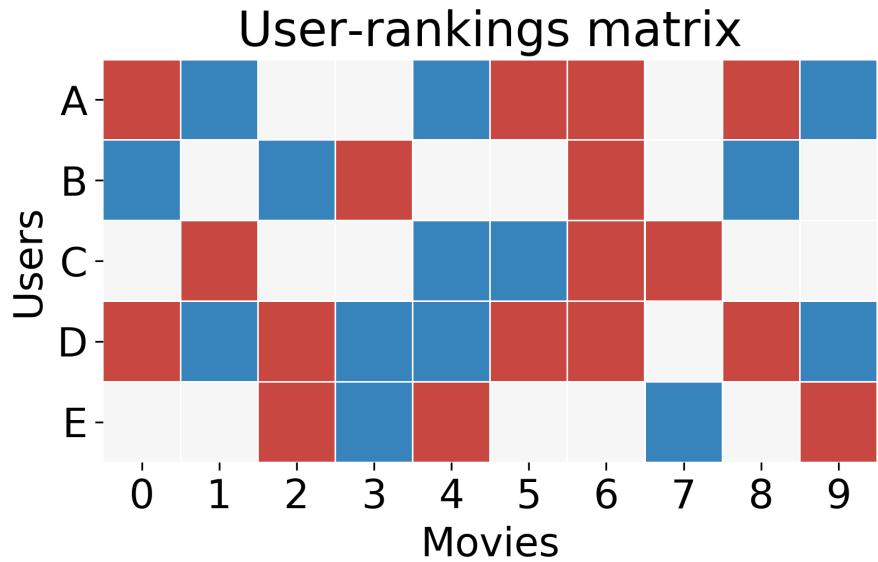
- Rows are users
- Columns are “+1” for like and “-1” for dislike
- We know which values are missing (unwatched)
- Goal: predict missing values to suggest movies to each user



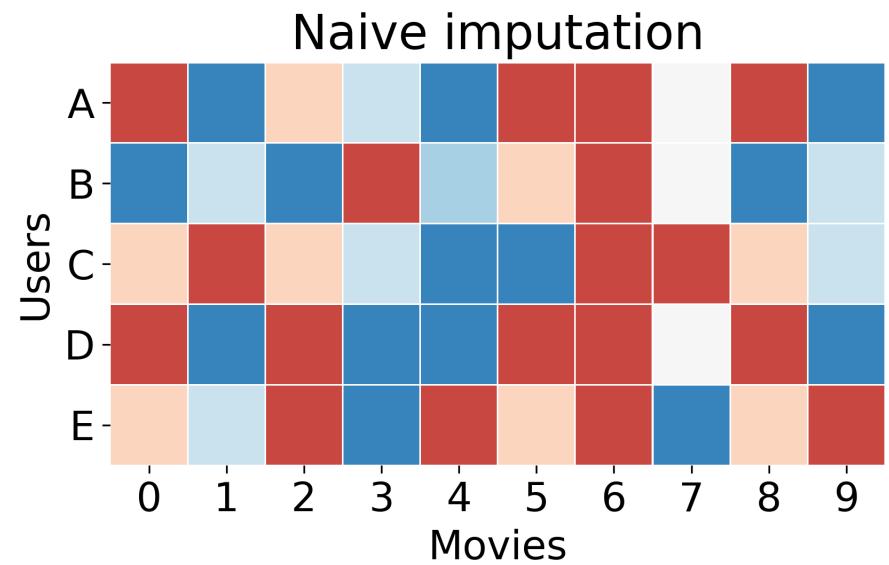
How can we guess missing values?

What do we do with missing data?

Raw data

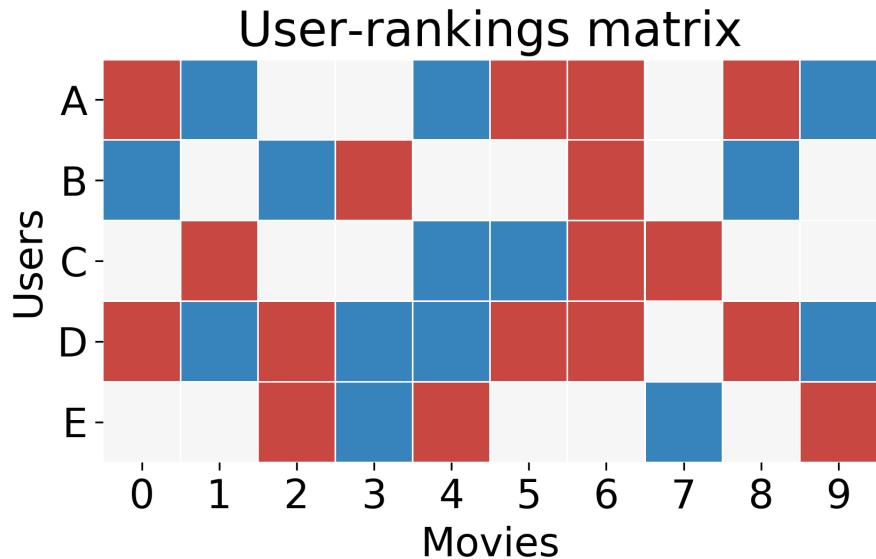


Naïve solution
Use the average ranking

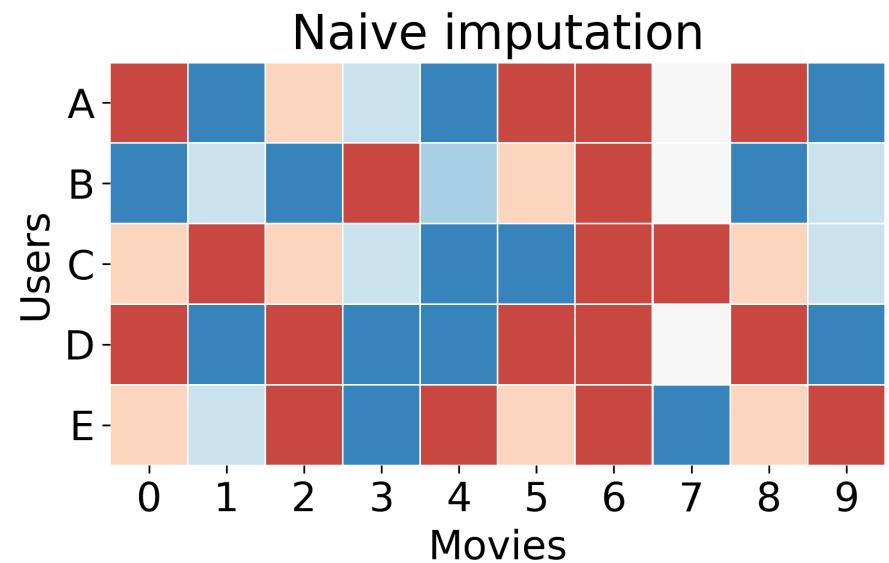


What do we do with missing data?

Raw data



Naïve solution
Use the average ranking

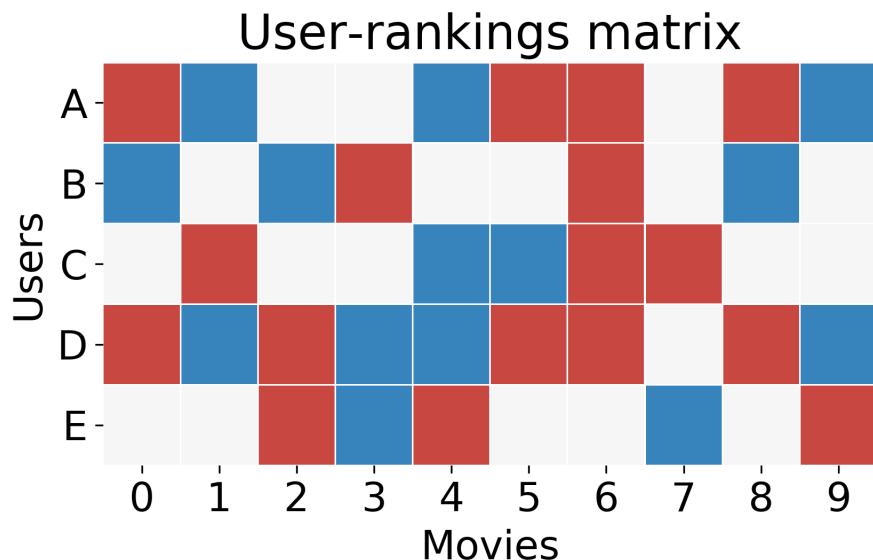


How can we improve on this?

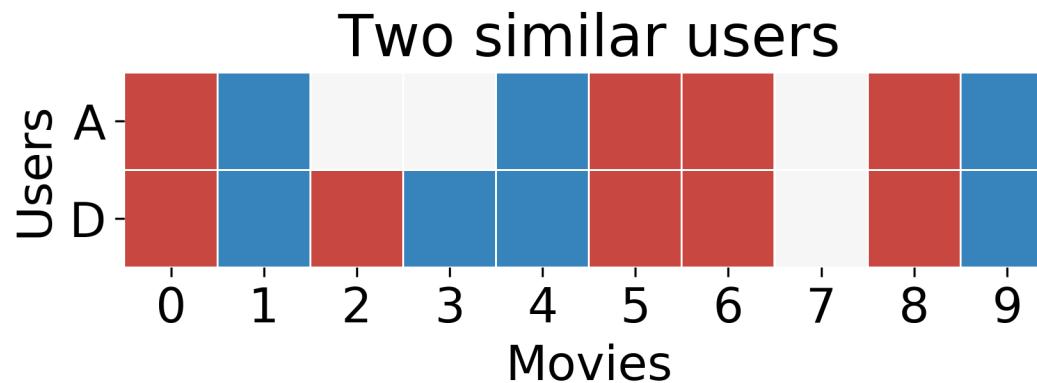
Information redundancy in the user-rankings matrix

Two sources of redundancy

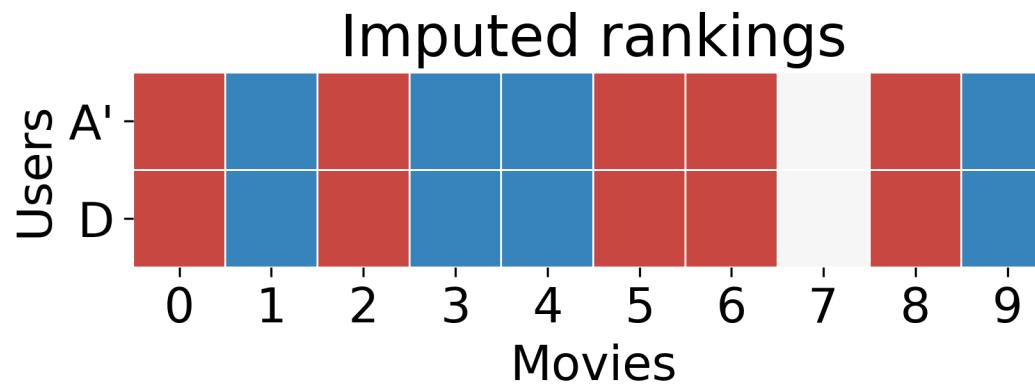
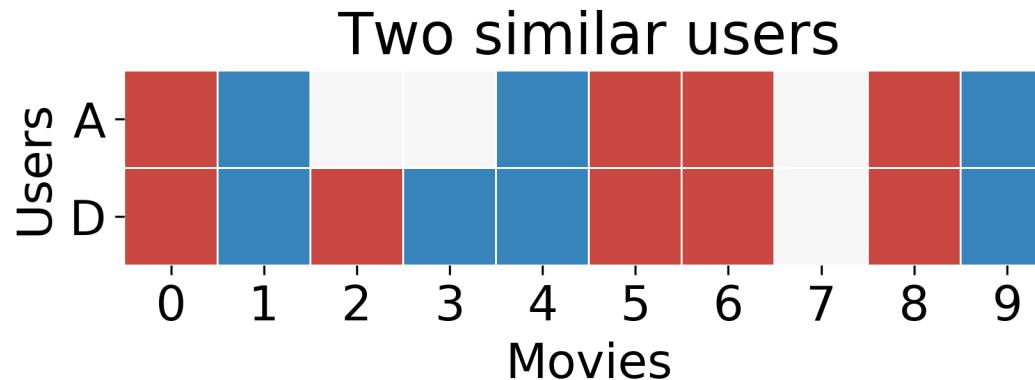
- Similarity between movies (e.g. rom-coms, horror movies)
- Similarity between users (e.g. “Likers of rom-coms”)



Identifying similar users to impute missing values

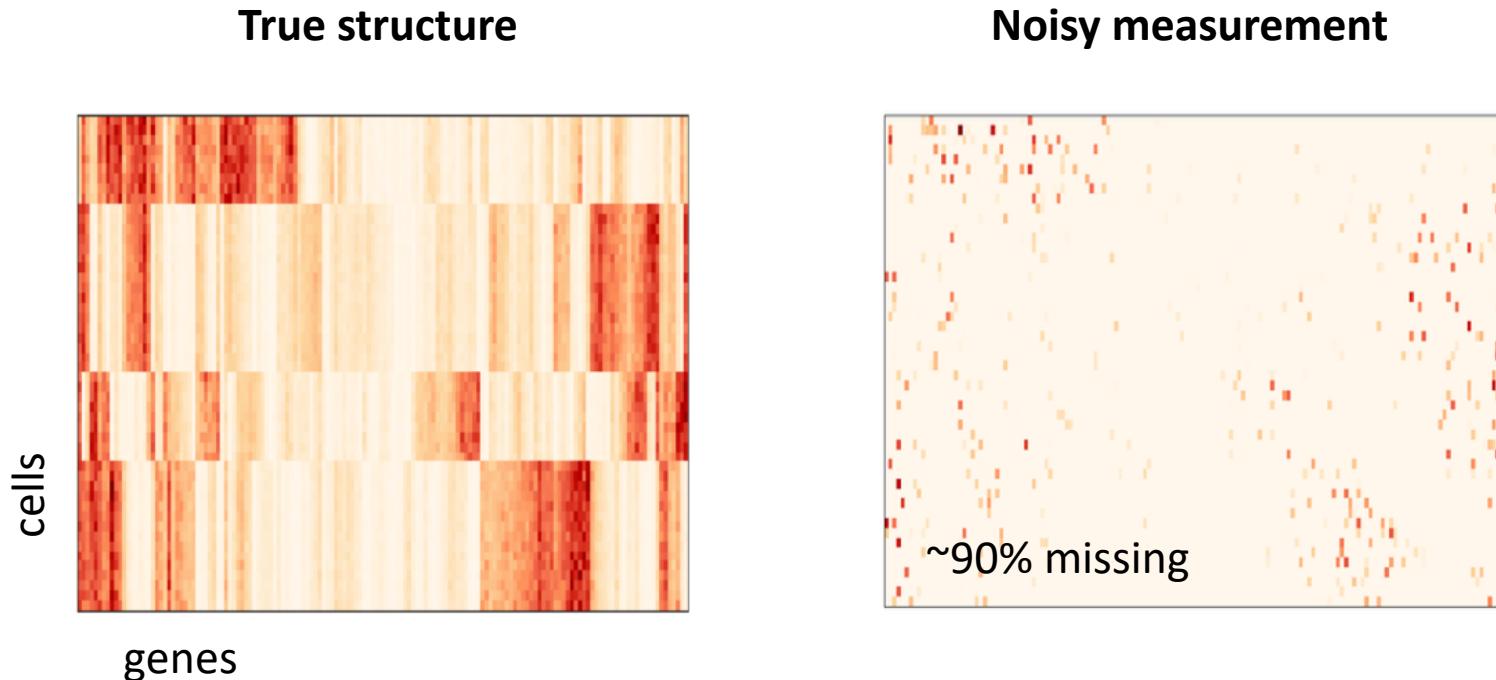


Identifying similar users to impute missing values

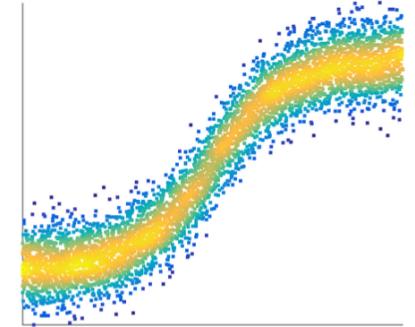


Problem: scRNA-seq data is sparse and noisy

How do we know which values to impute?



We only capture ~5-15% of the transcripts per cell
due to inefficient RT in small volumes

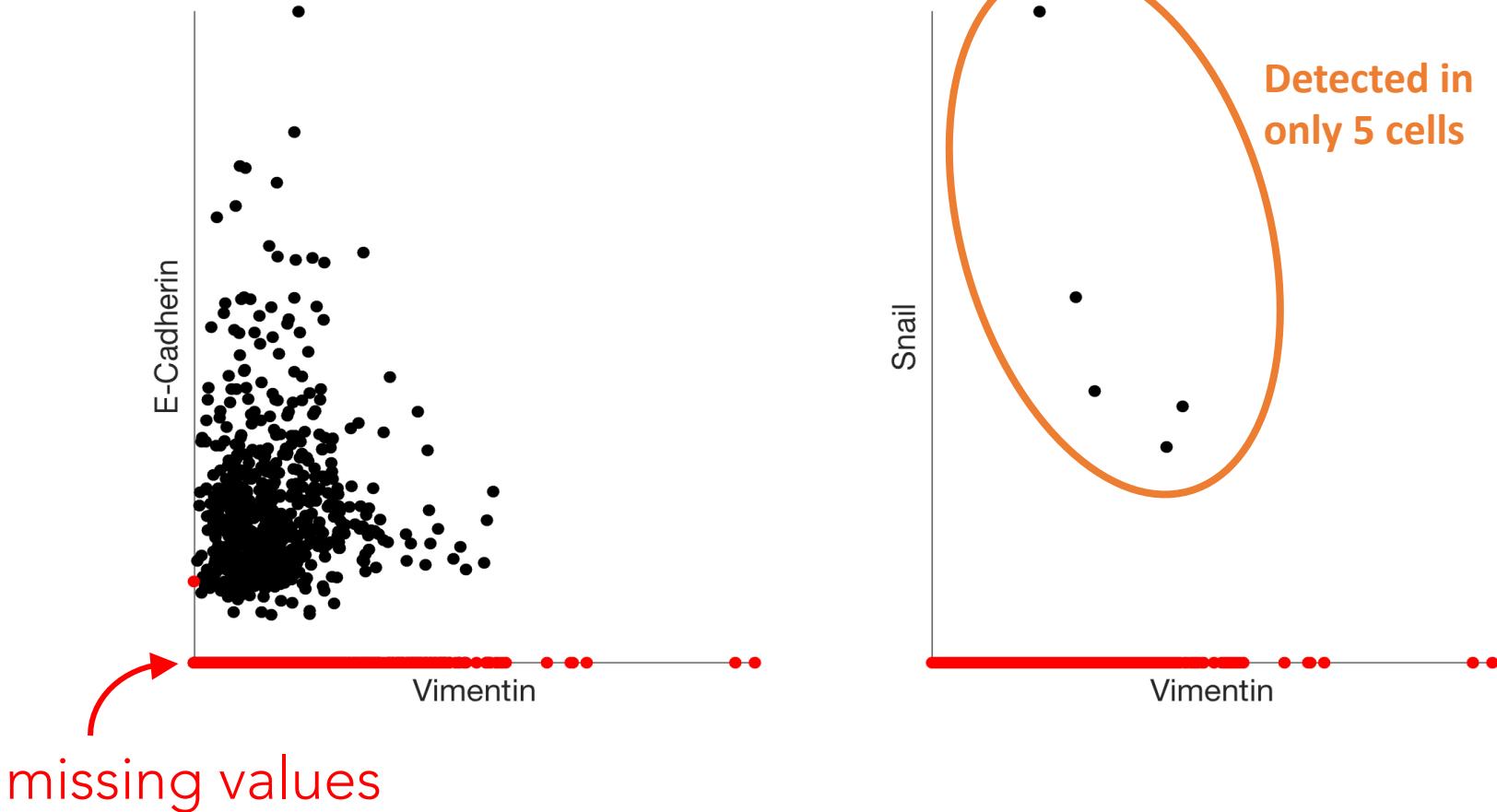


MAGIC

Markov Affinity-based Graph Imputation of Cells
(van Dijk et al. *Cell*, 2018)

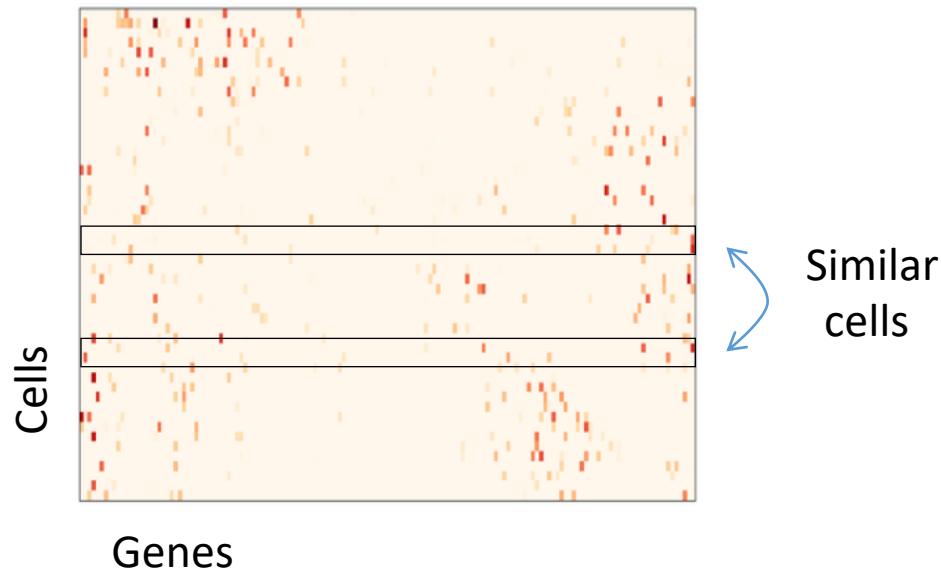
Main idea: cells learn values from their neighbors, gaining information from similar cells while retaining their individual identity.

HMLE breast cancer cell line, TGF- β induced EMT



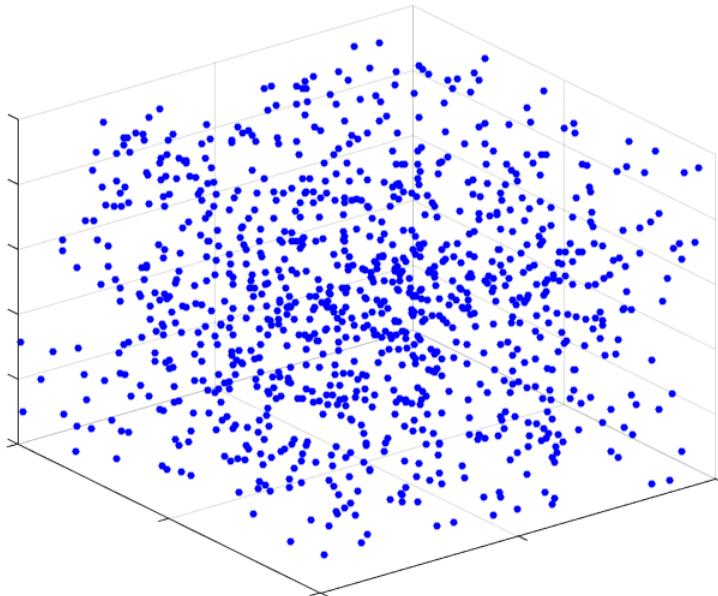
dropout obscures gene-gene relationships

Solution: Recover data by learning from similar cells



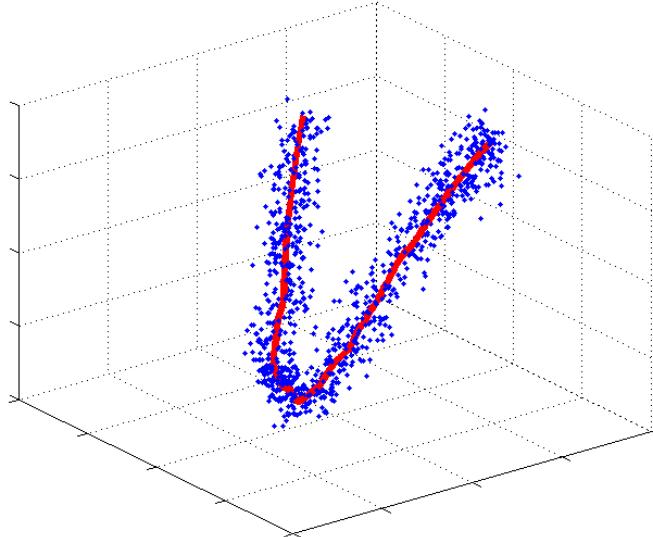
- Biologically similar cells will have similar expression
- Cells “exchange” information

Low dimensional structure



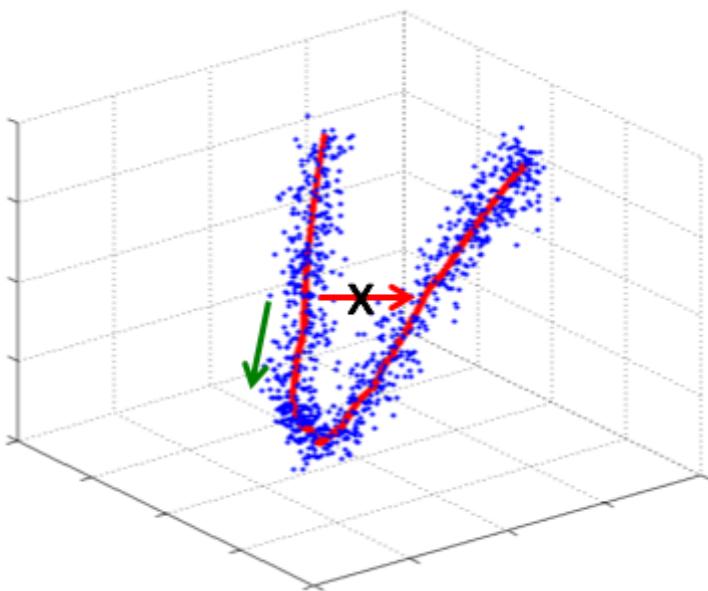
- Many features (e.g. genes) but high degree of redundancy
- Features are correlated
- States (e.g. cell types) are constrained

Manifold

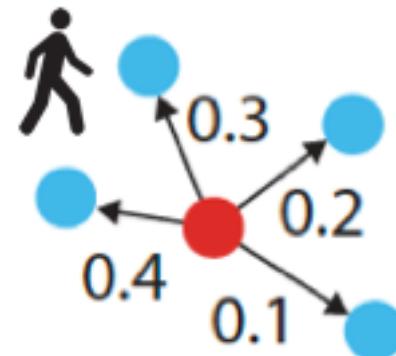


- Locally Euclidean
- Smooth
- Low dimensional
- Data has a “shape”

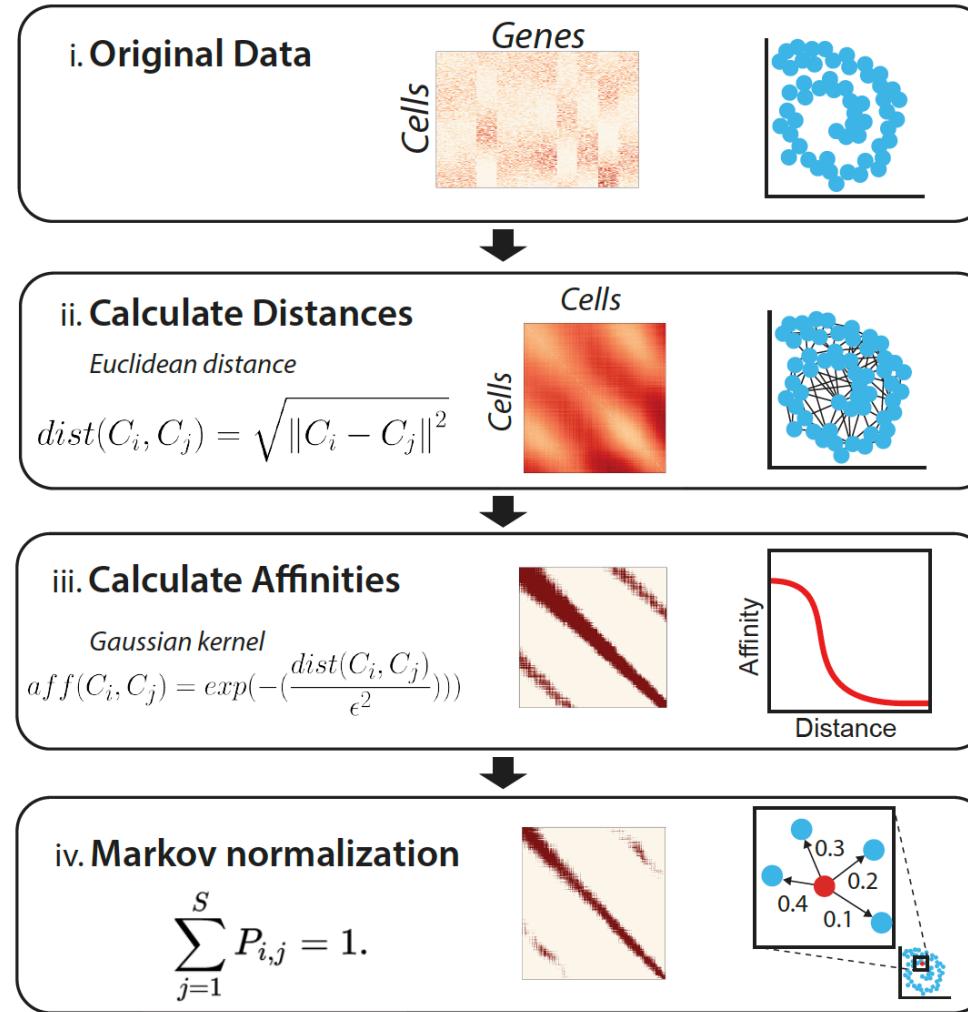
Manifold learning via data diffusion



- Biological neighborhood = manifold
- Walk via small local steps
- Random walk = diffusion



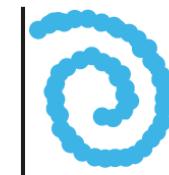
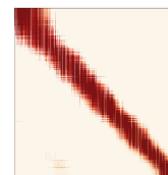
The Algorithm



Diffusion and imputation step

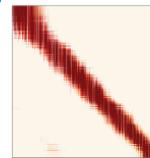
v. Exponentiate markov
matrix

$$[\quad]^t$$



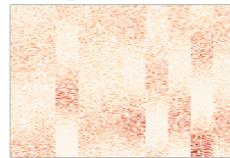
vi. Impute gene expression

Exp. Markov Mat.



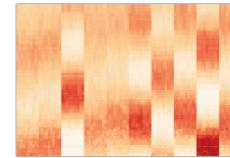
X

Original Data

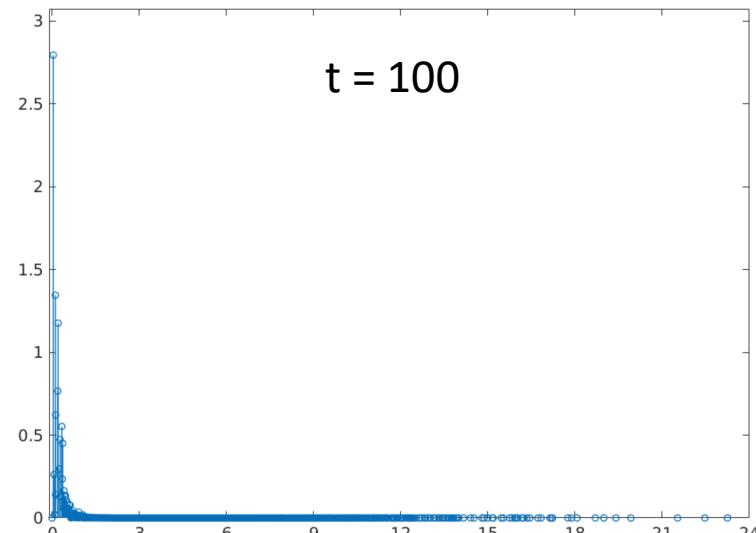
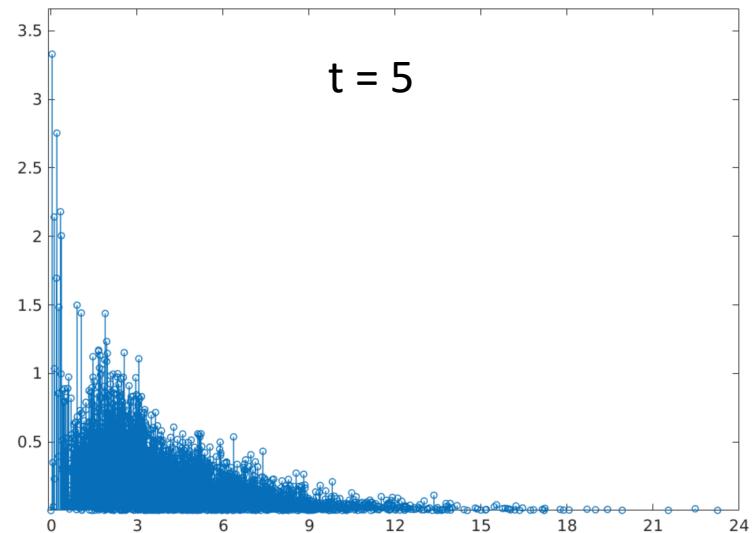
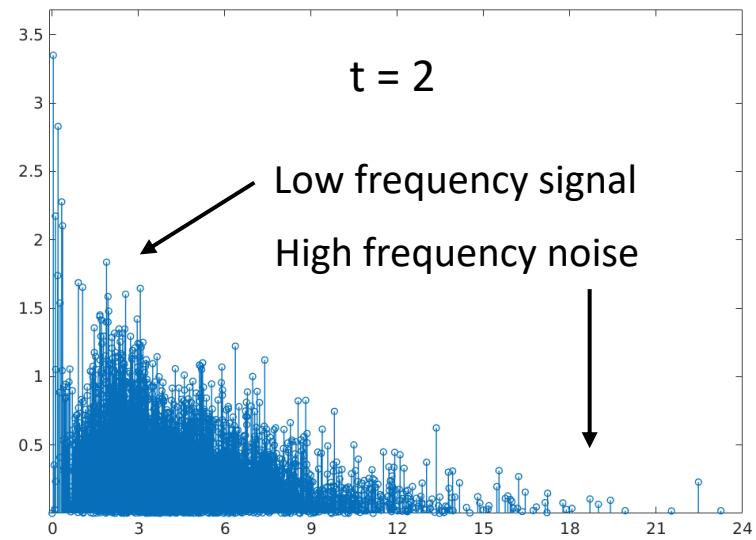
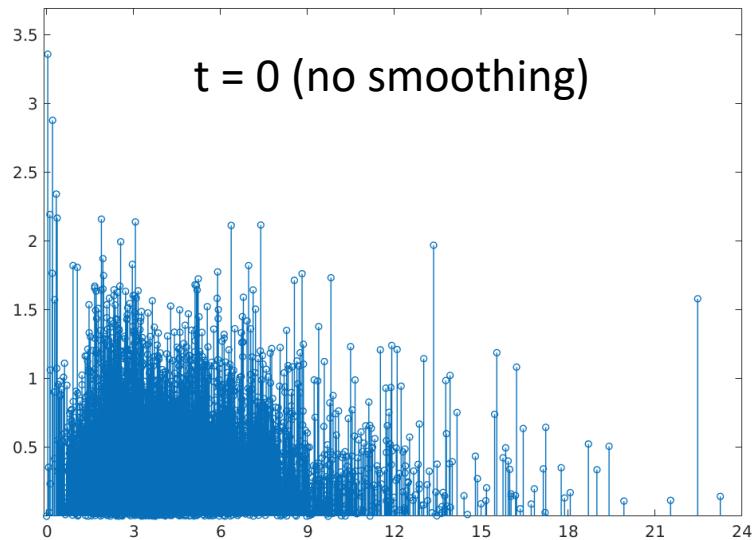


=

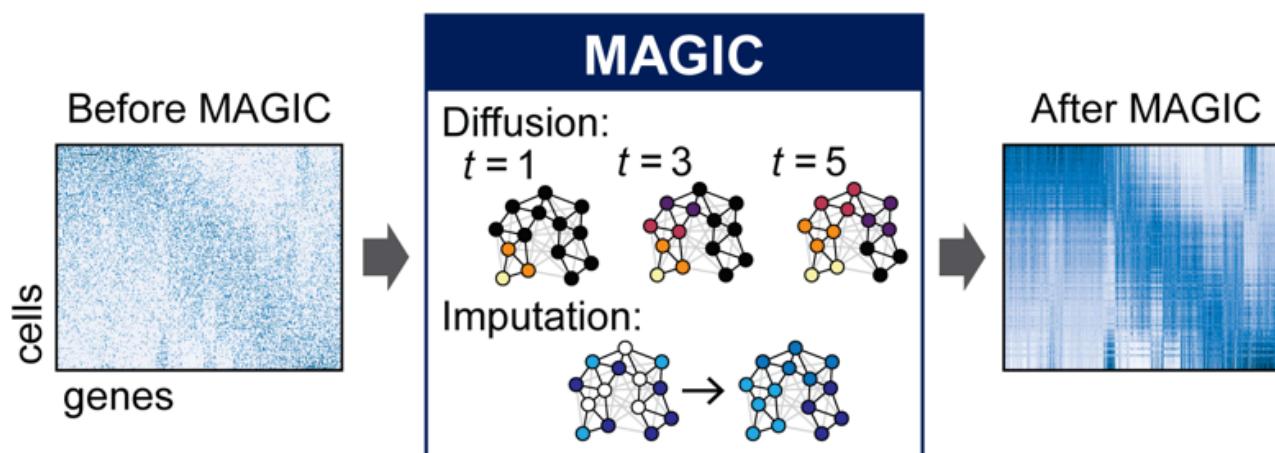
Imputed Data

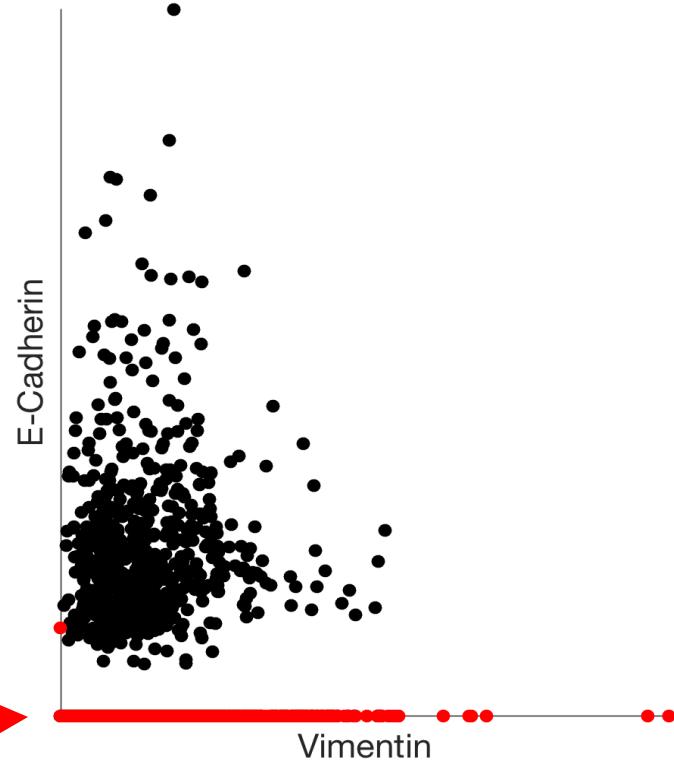


Graph diffusion = low pass filter on graph spectrum

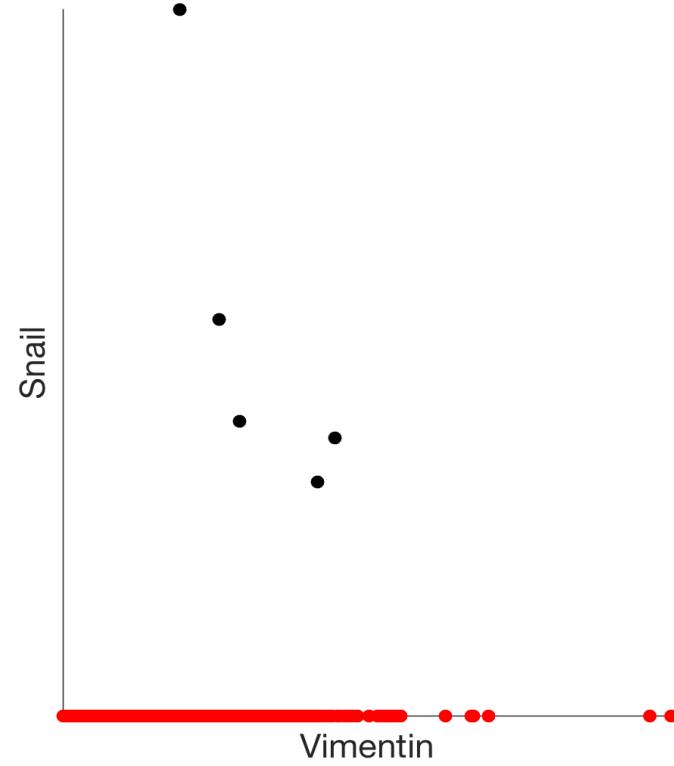


Markov Affinity-based Graph Imputation of Cells (MAGIC)

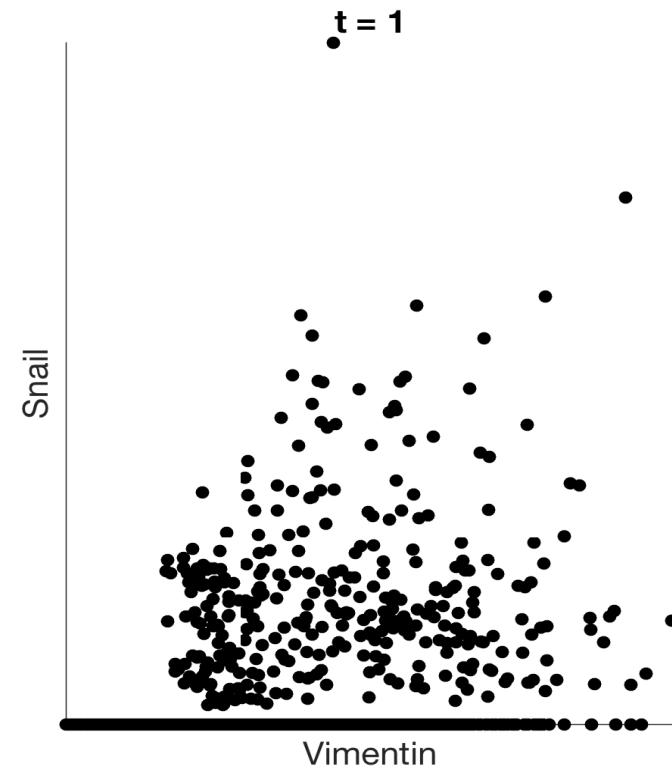
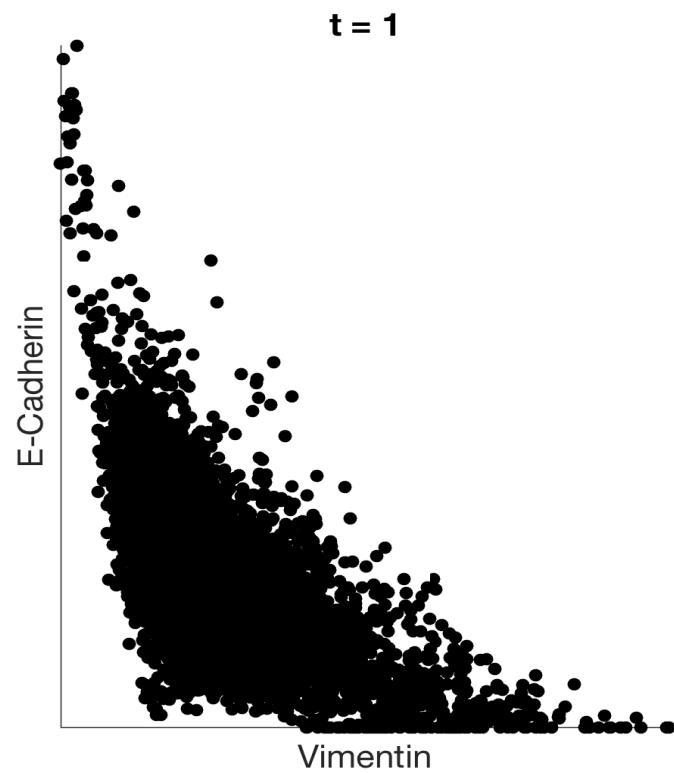




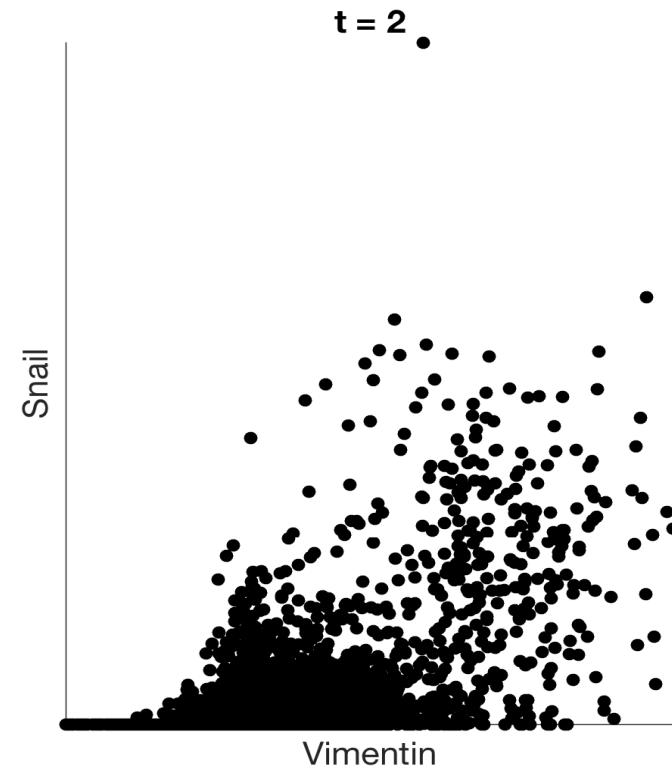
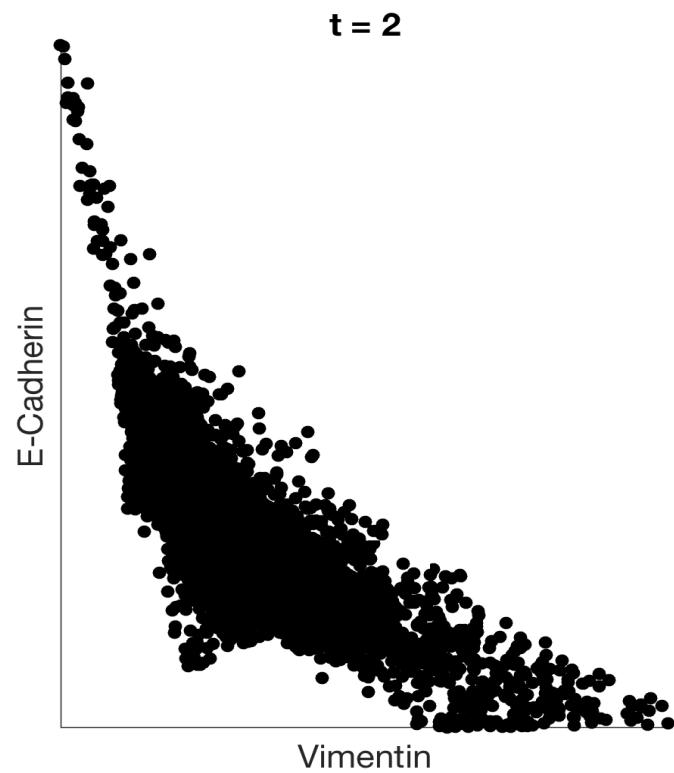
missing values



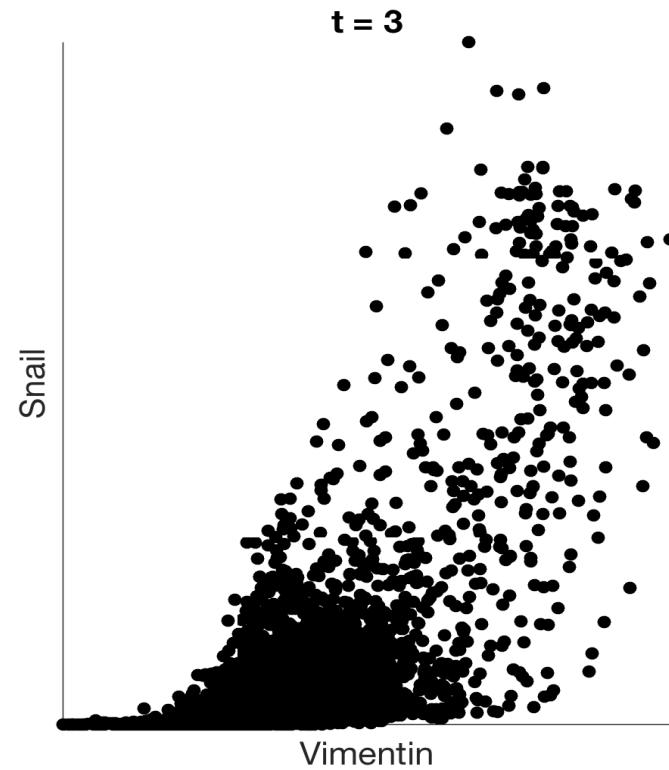
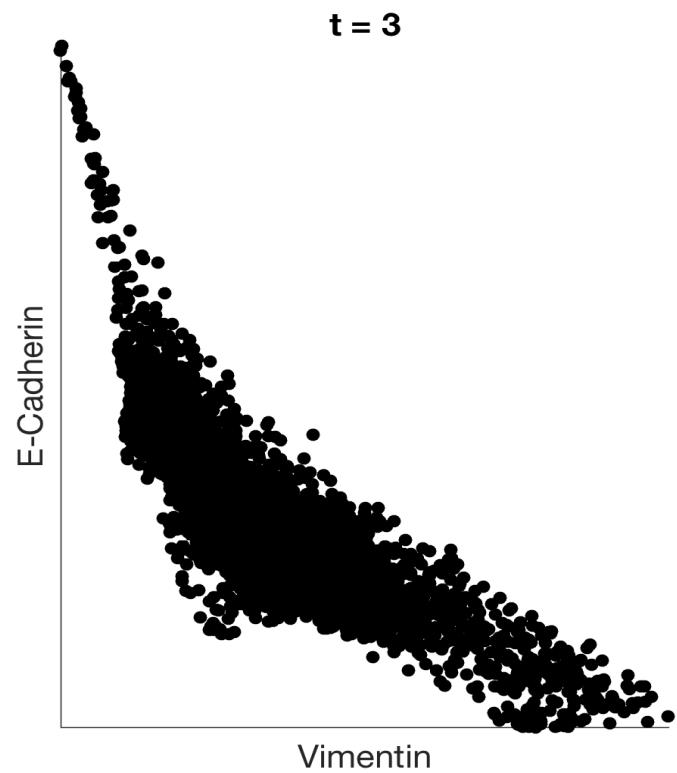
original data with dropout



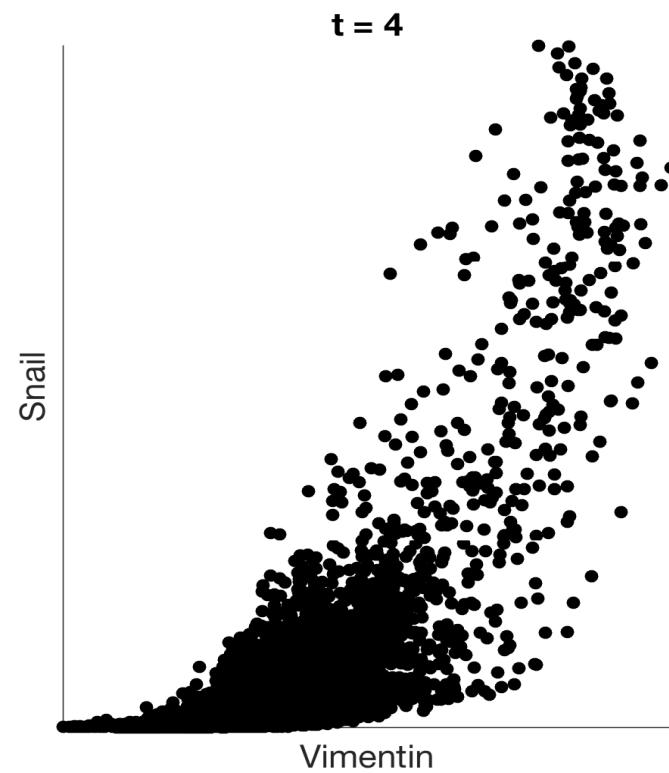
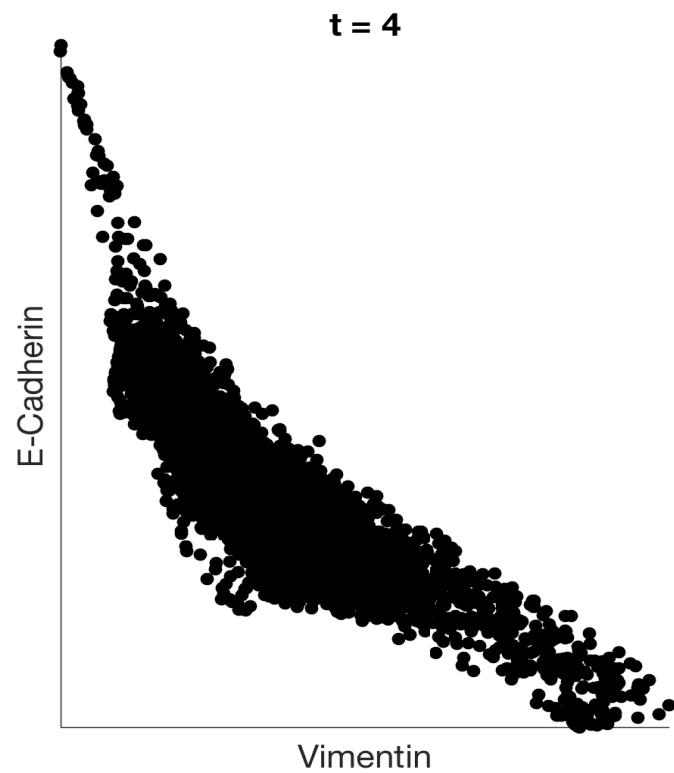
imputation with MAGIC



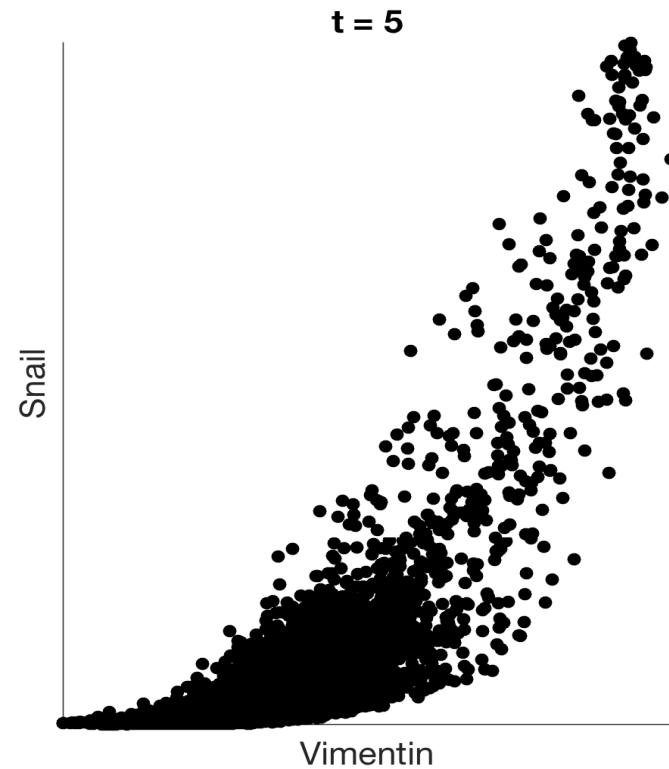
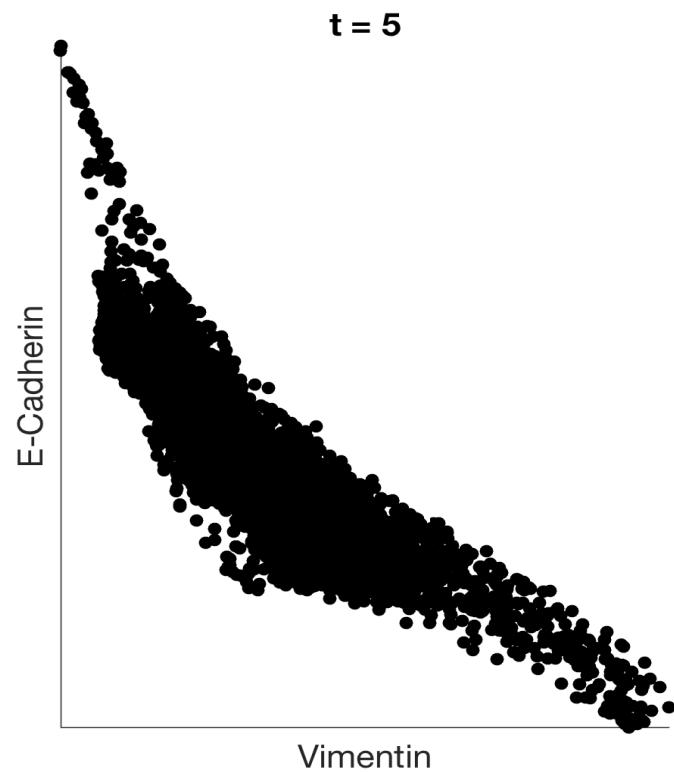
imputation with MAGIC



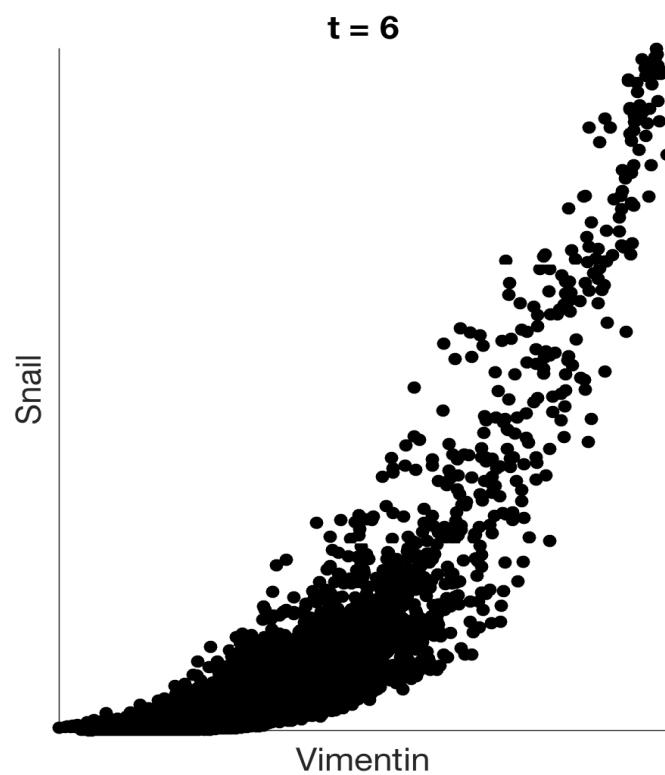
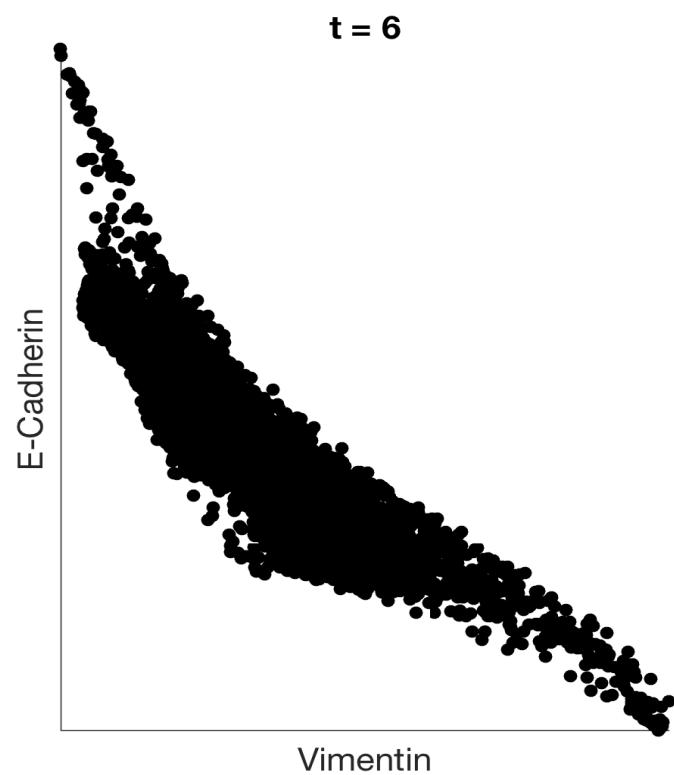
imputation with MAGIC



imputation with MAGIC

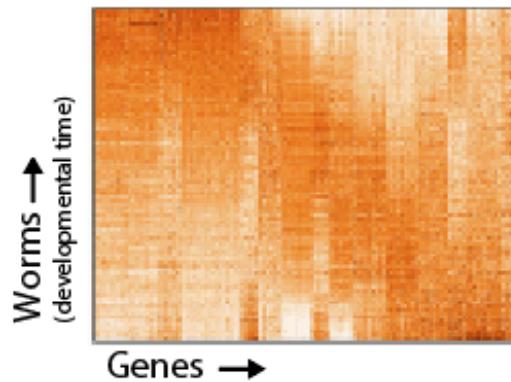


imputation with MAGIC

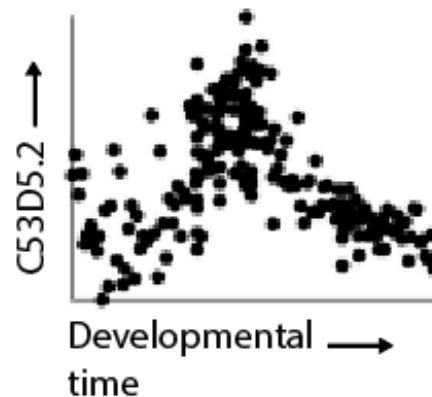
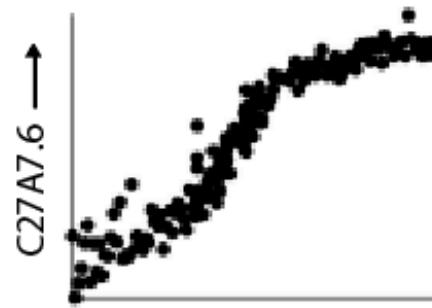


imputation with MAGIC

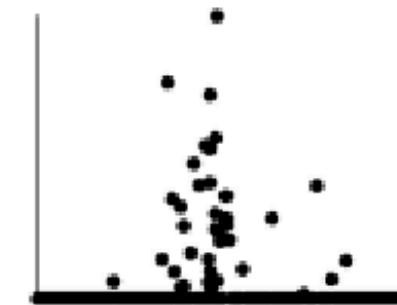
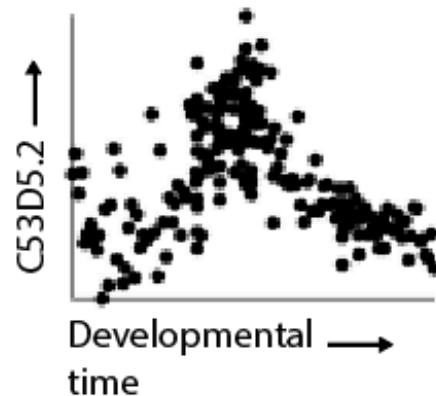
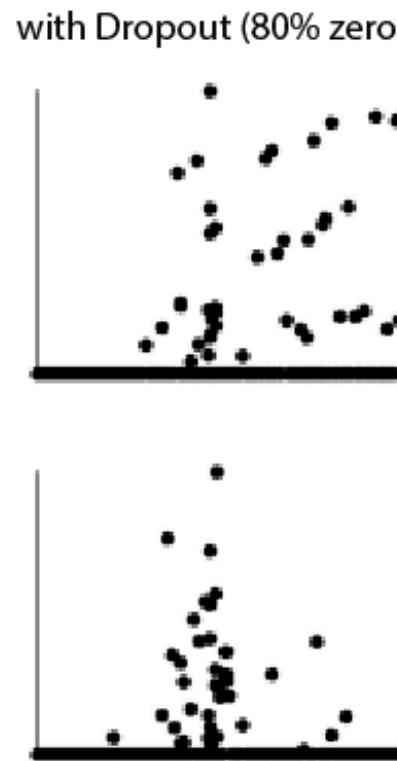
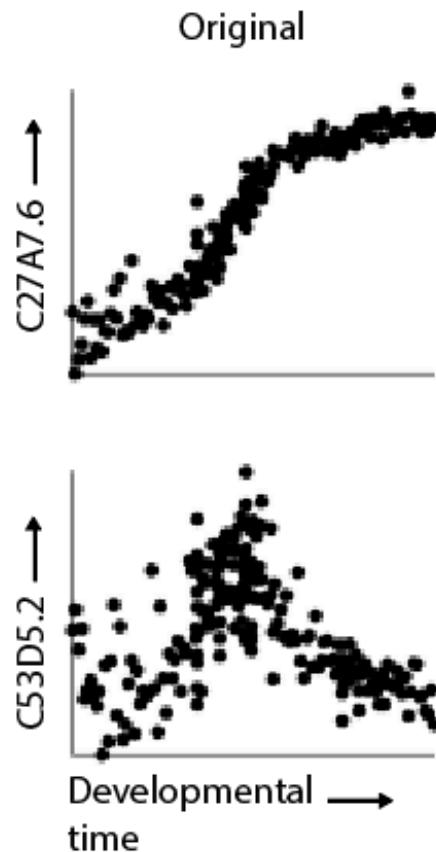
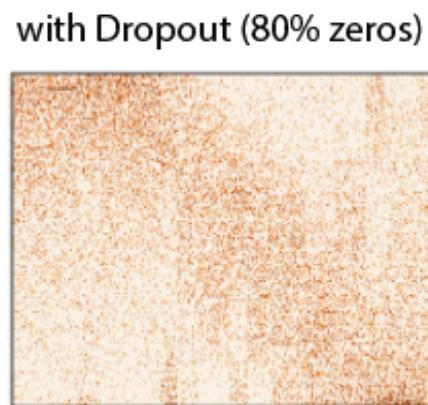
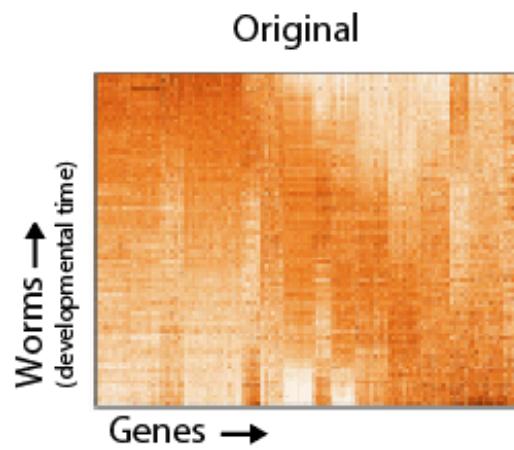
Original



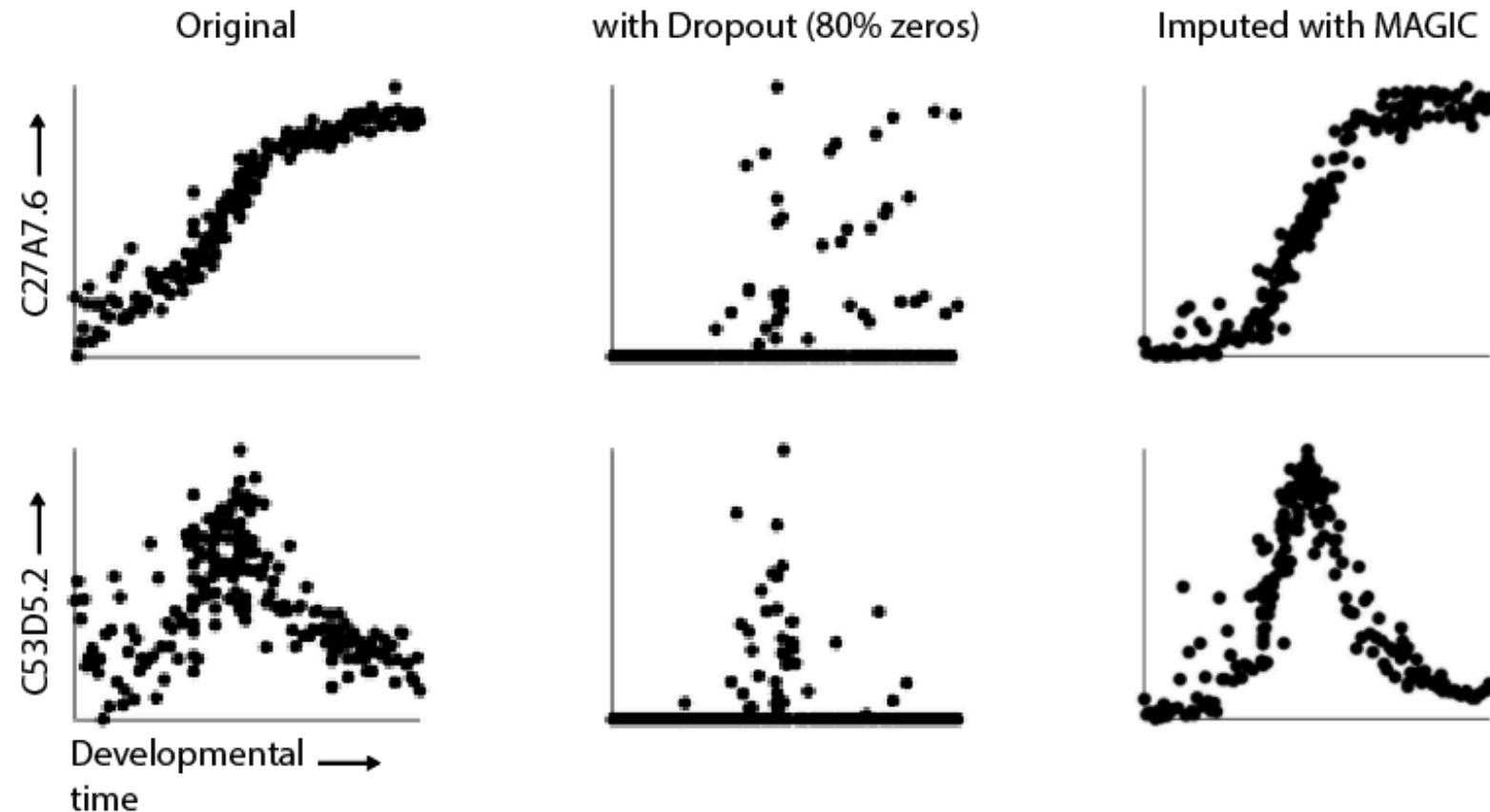
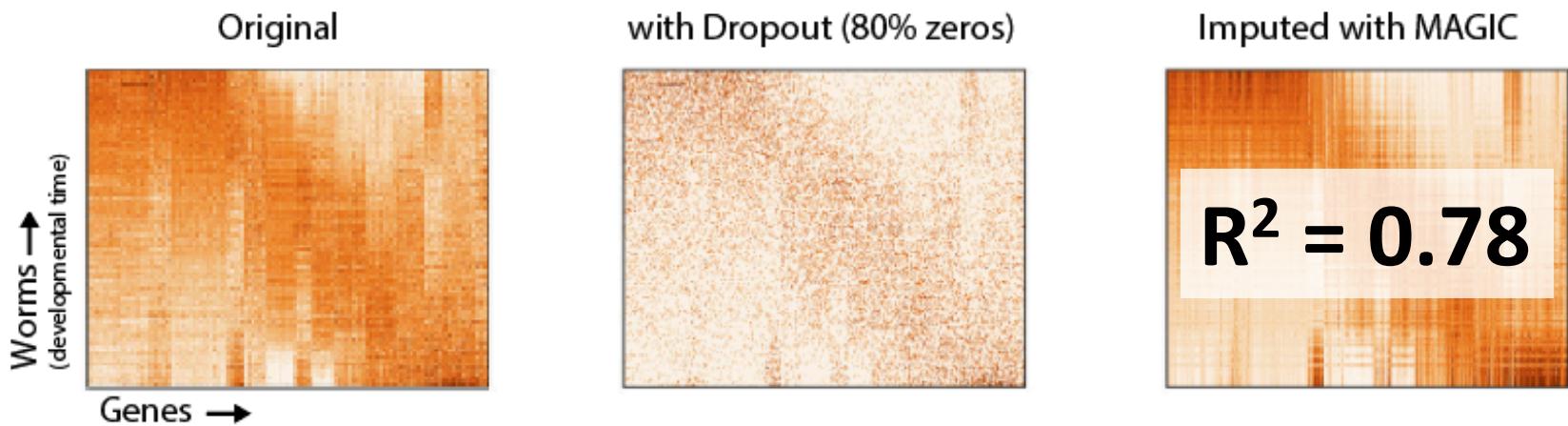
Original



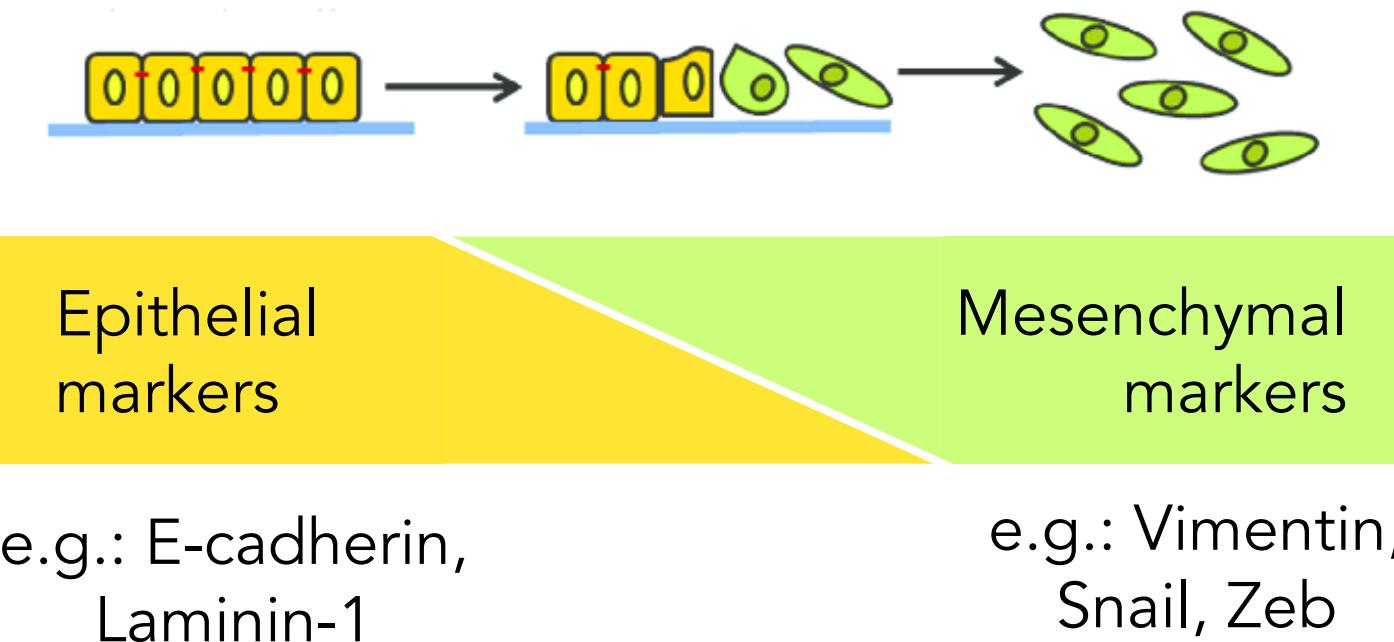
(Francesconi, 2014, Rockman, 2010)



(Francesconi, 2014, Rockman, 2010)

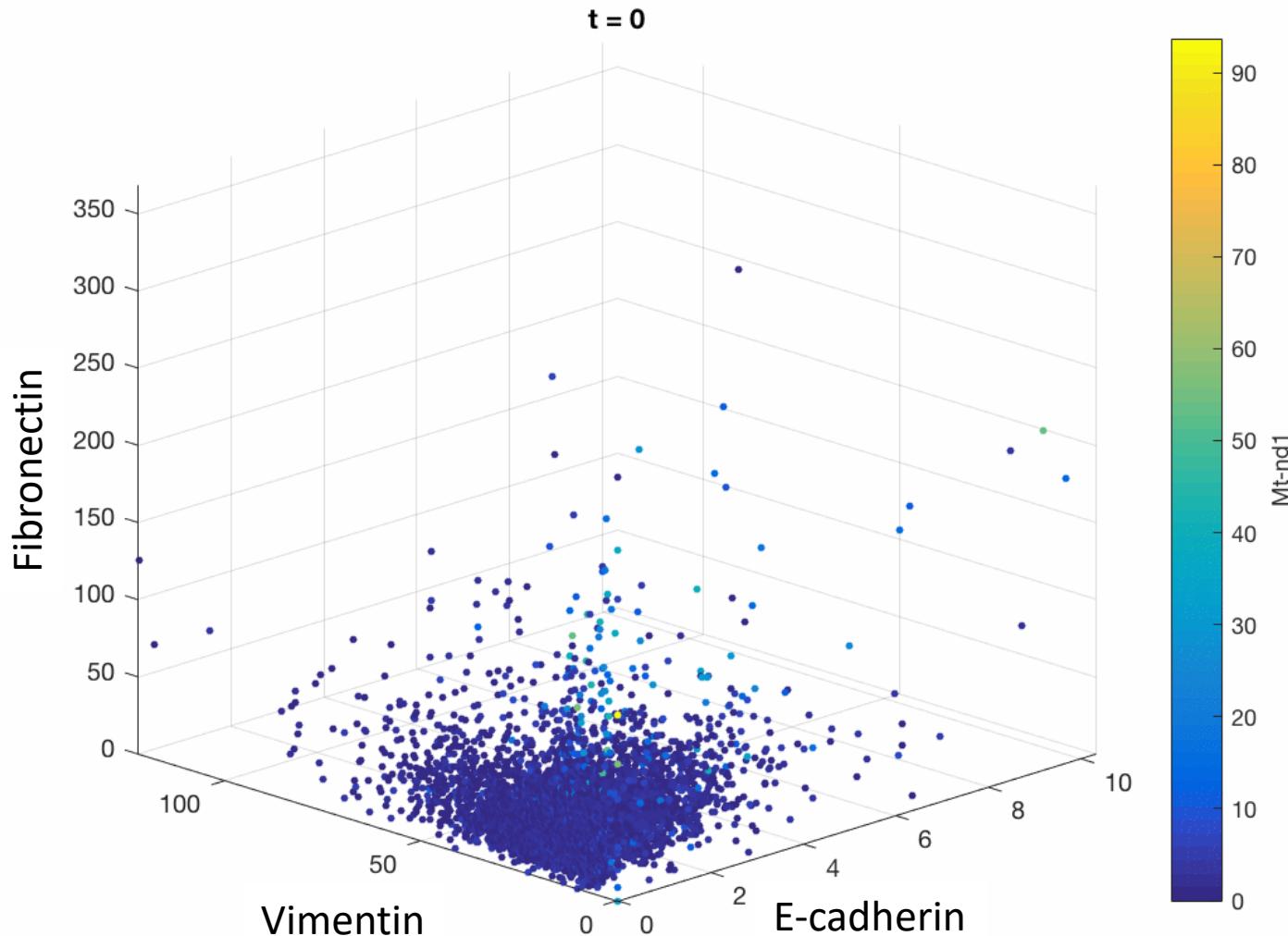


TGF- β induced EMT (epithelial-to-mesenchymal transition)



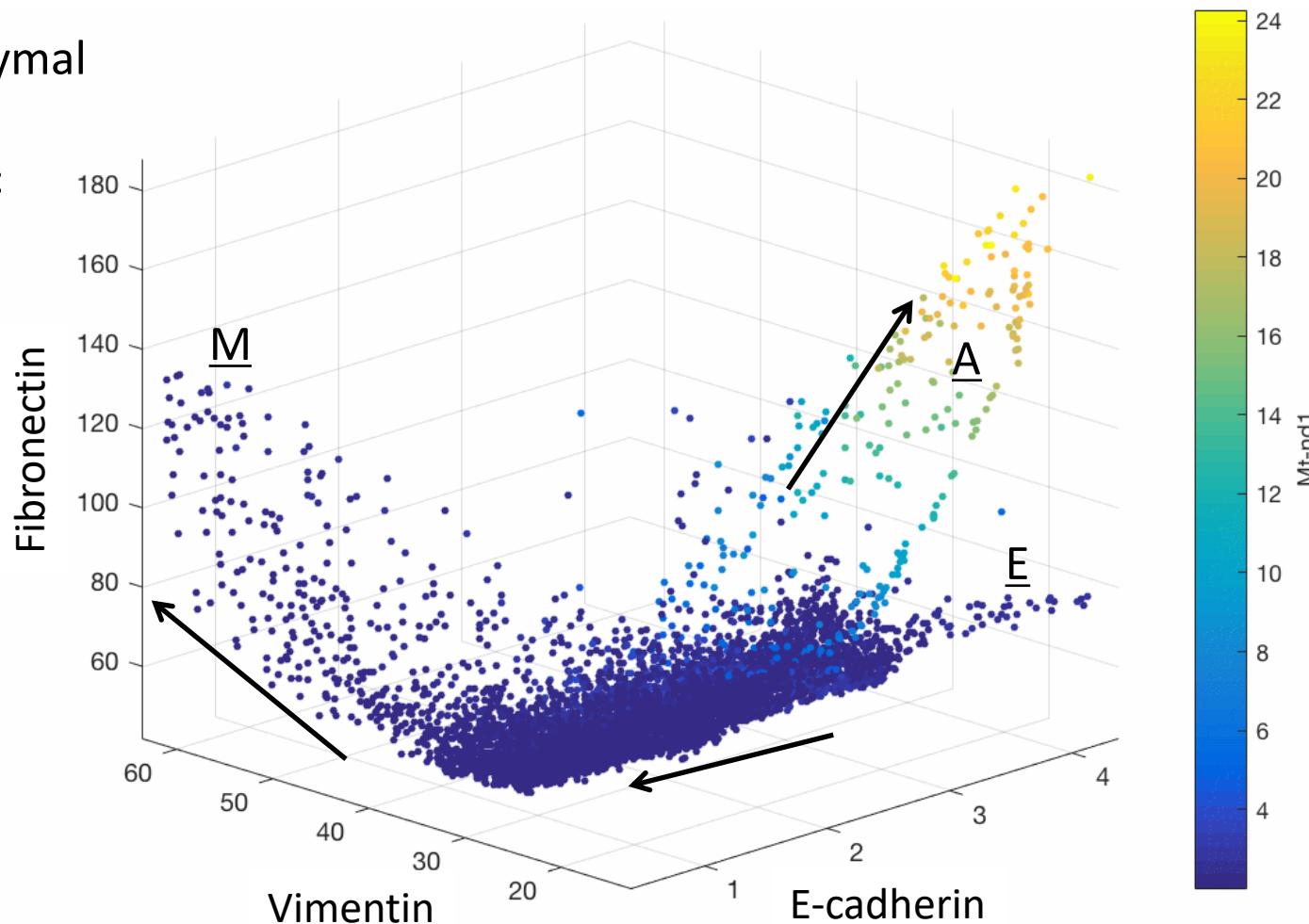
With Weinberg Lab (Whitehead):
Diwakar Pattabiraman, Brian Bierie, Christine Chaffer

MAGIC reveals EMT continuum

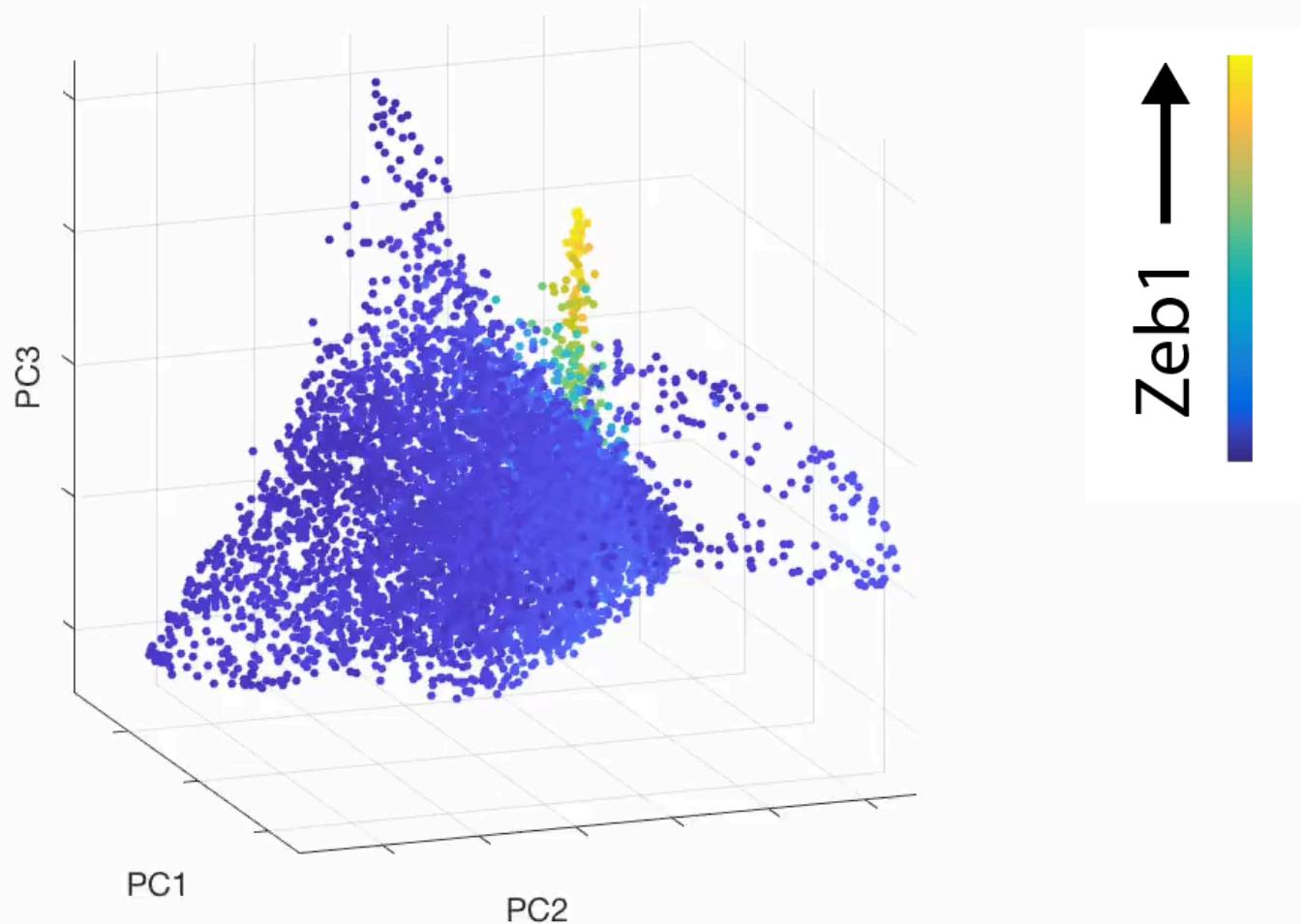


MAGIC reveals EMT continuum

Mesenchymal
Epithelial
Apoptotic

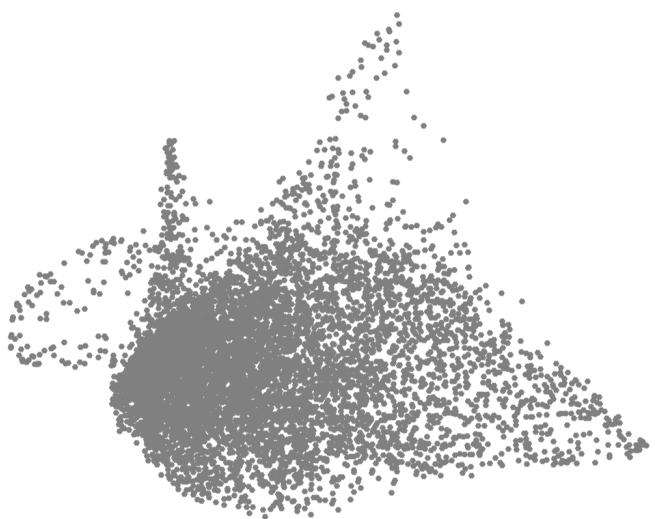
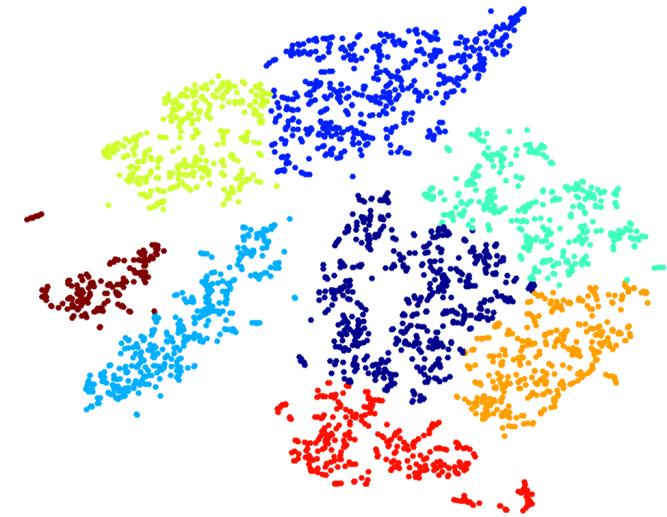


MAGIC reveals complex cellular state-space

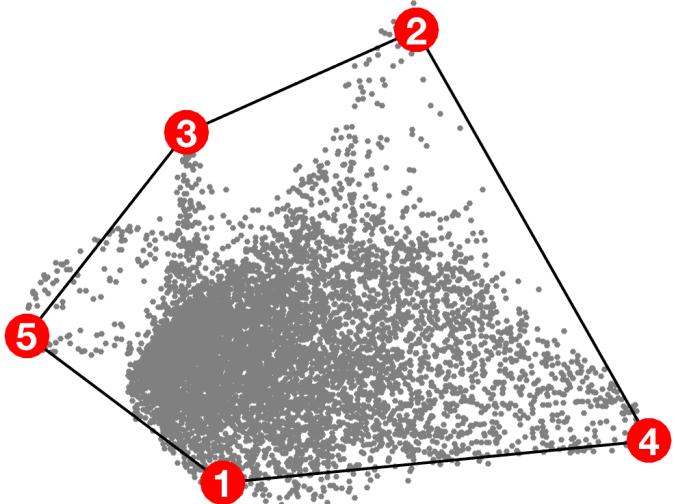




Clustering

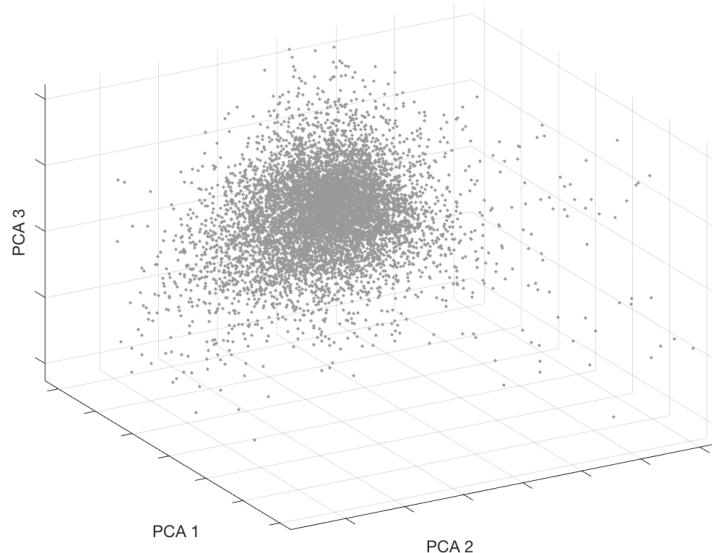


Archetypal analysis

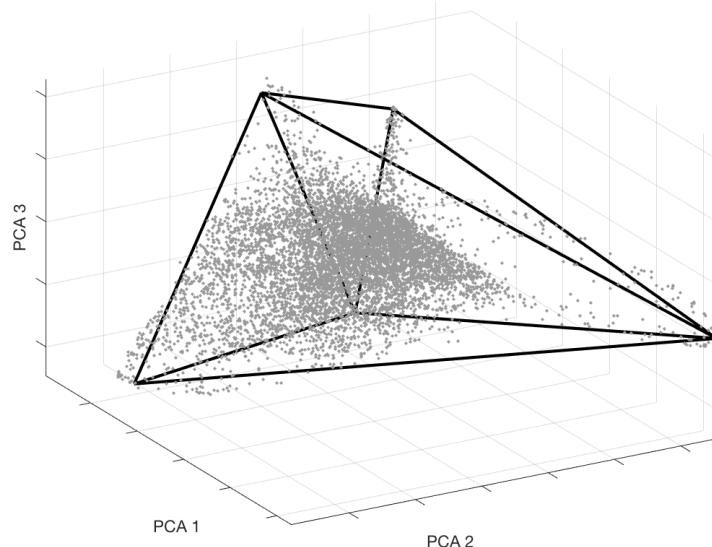


Archetypal Analysis

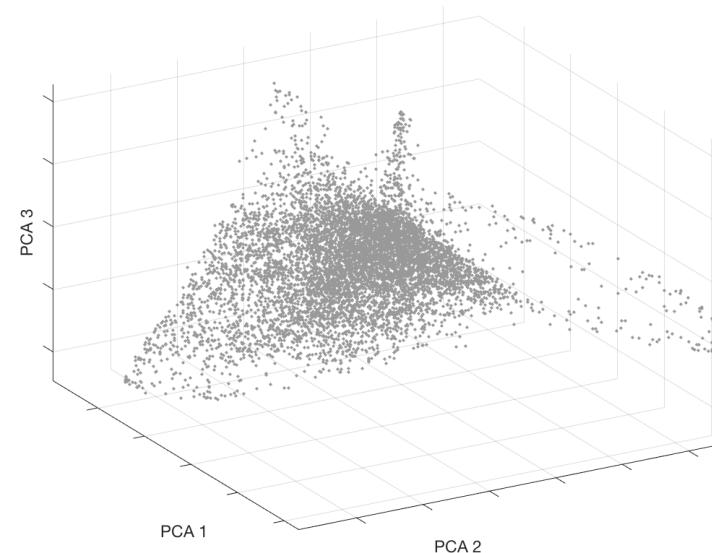
Data before MAGIC



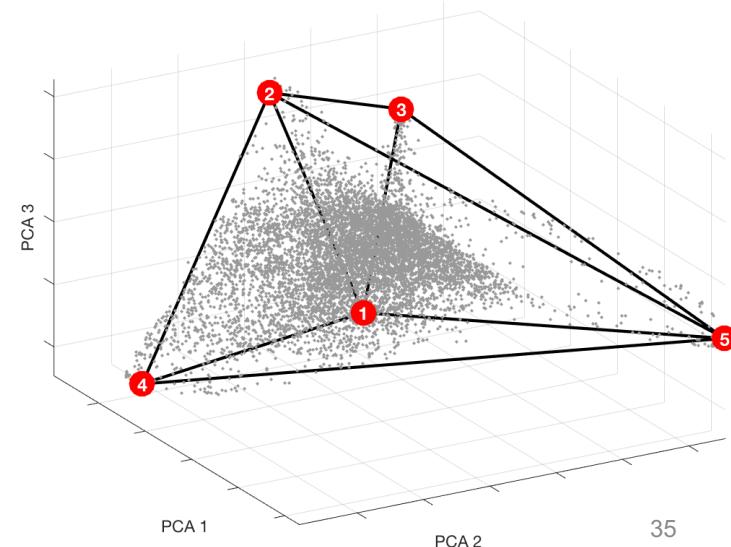
Fit polytope

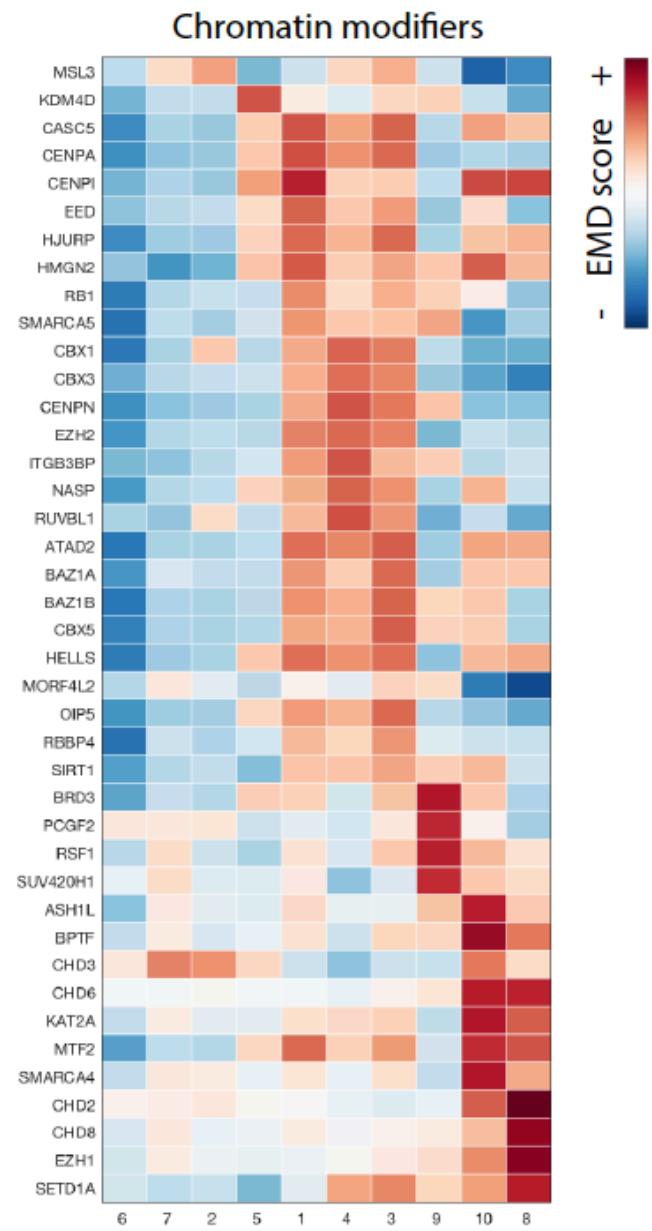
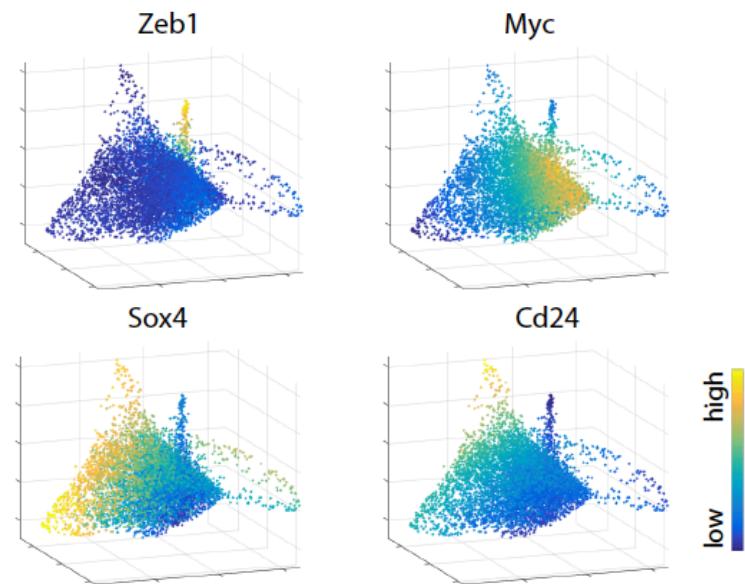
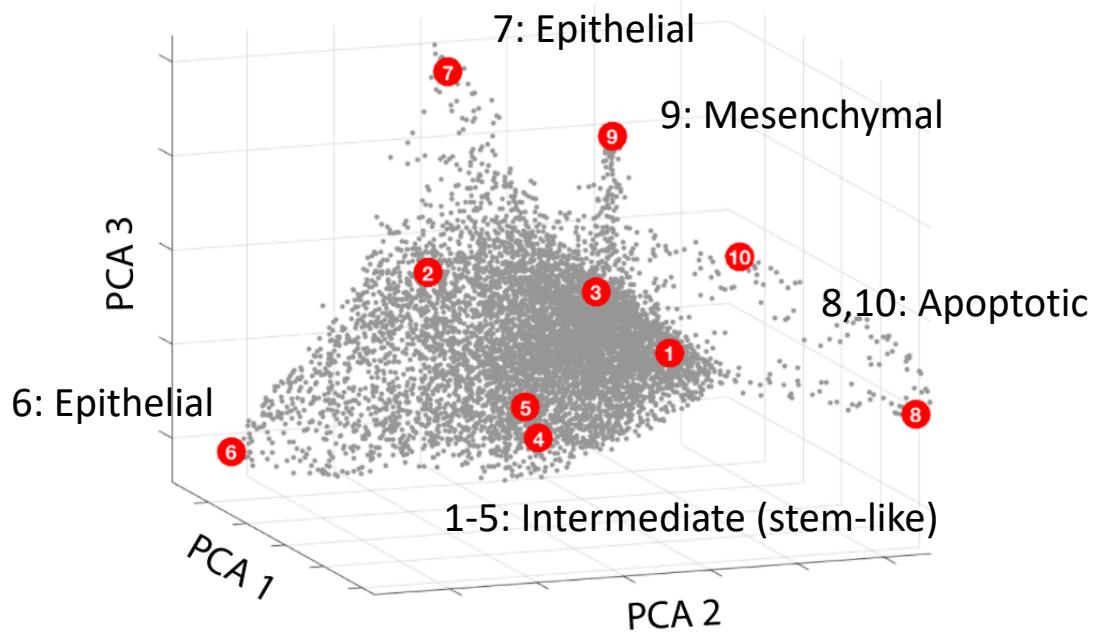


Data after MAGIC

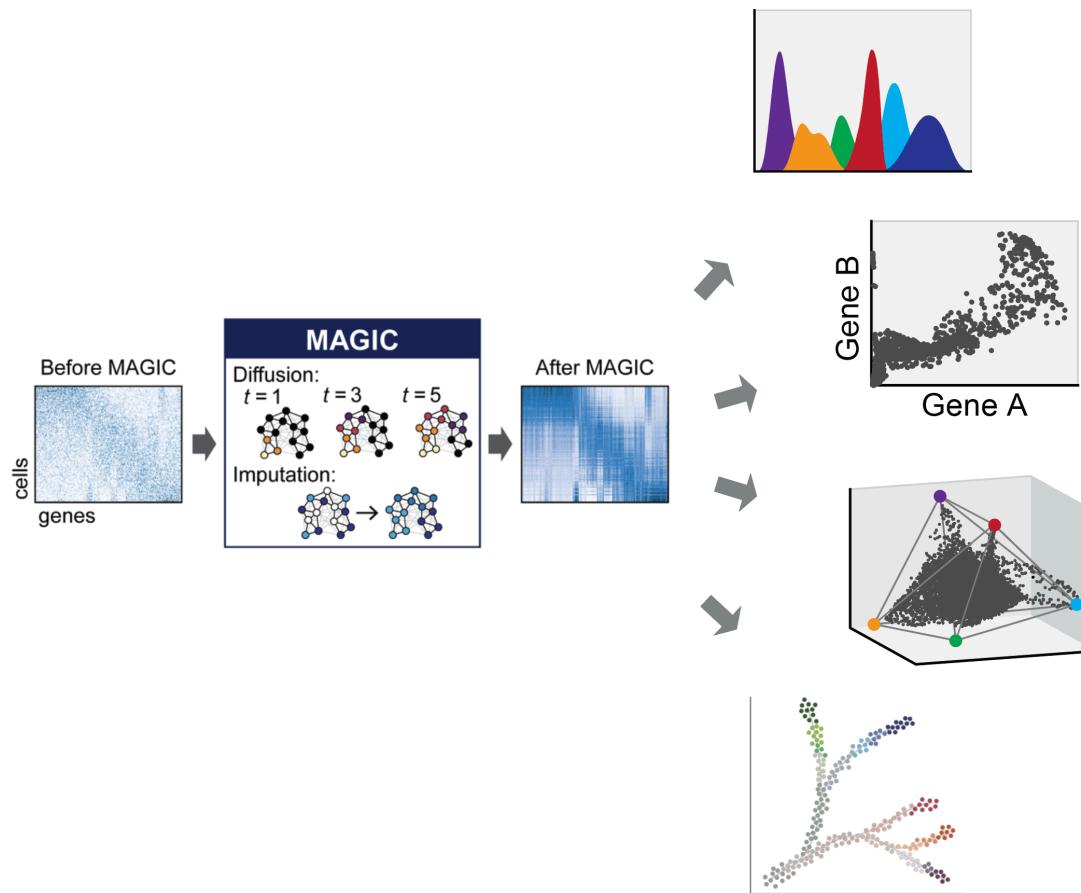


Infer archetypes



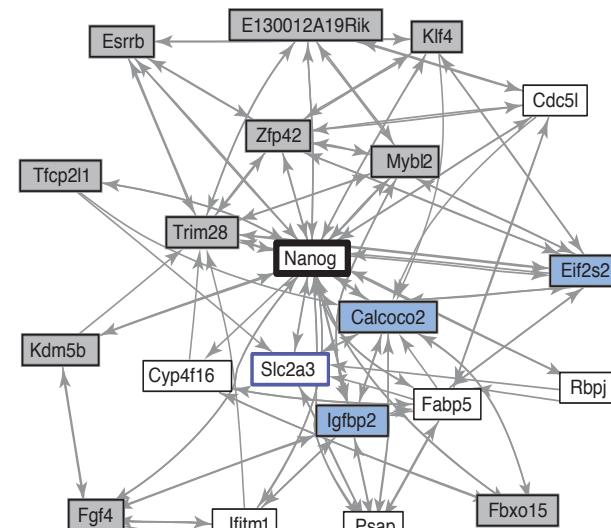


MAGIC enables downstream analyses





scRNA-seq data

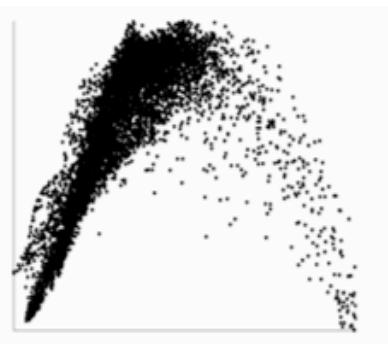


regulatory network

Using “snapshot” scRNA-seq data to predict transcription factor-target interactions involved in the epithelial—mesenchymal transition.

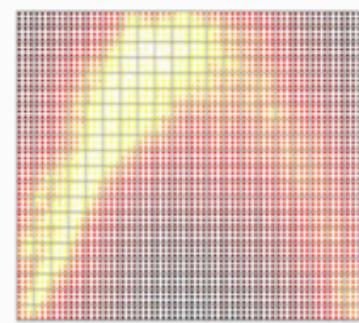
Data after MAGIC

Gene Y



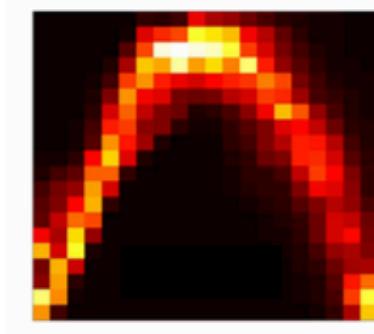
Gene X

Density estimation

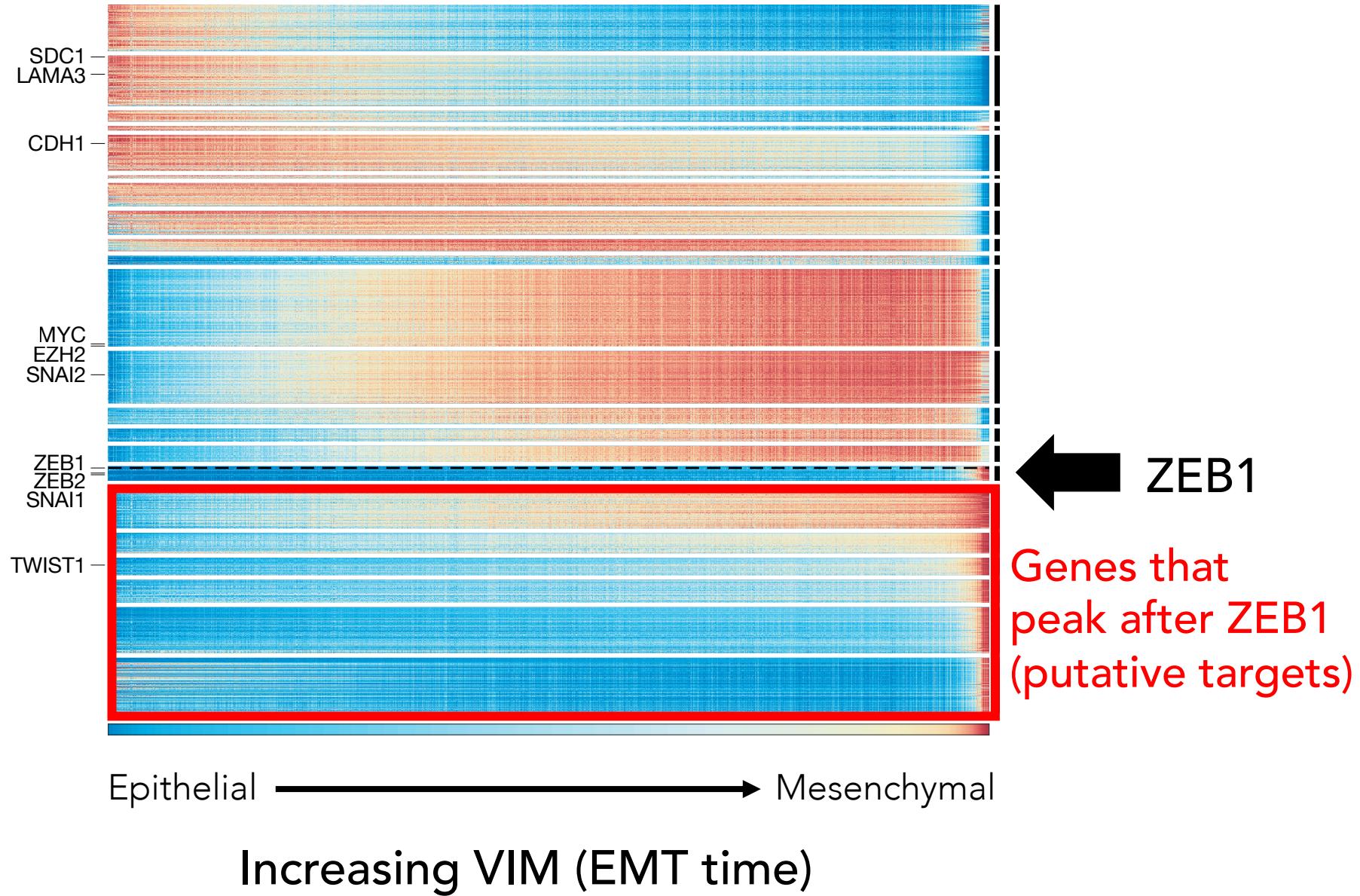


low high

Mutual information

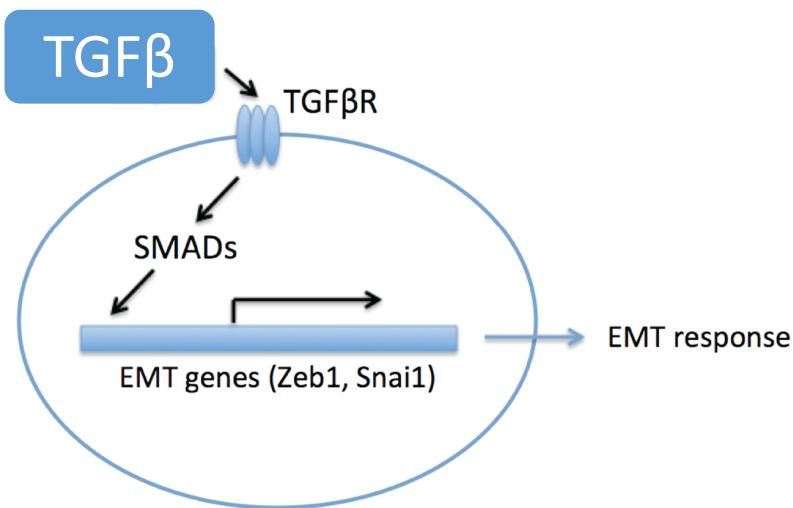


quantifying gene-gene relationships

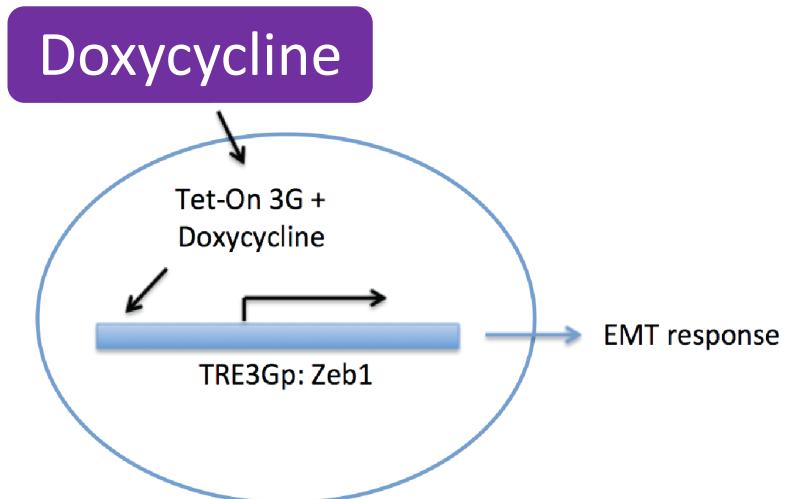


Validating predictions with ZEB1 induction

TGF β induction of EMT



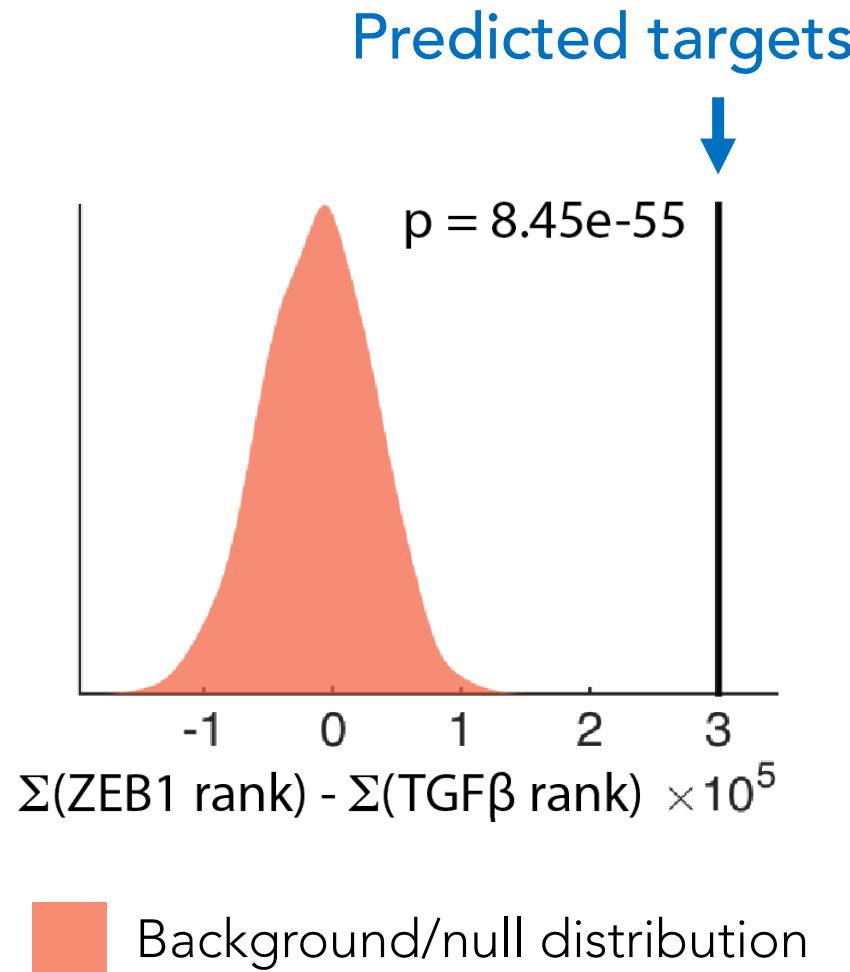
ZEB1 induction of EMT



With Weinberg Lab (Whitehead):

Diwakar Pattabiraman, Brian Bierie, Christine Chaffer

Validating predictions with ZEB1 induction



Summary

- scRNA-seq data is sparse and noisy
- Sparsity obscures biology
- Imputation via manifold learning (MAGIC)
- Imputed data enables downstream analyses

