# Dimensionality Reduction for Visualization



Kevin Moon (kevin.moon@usu.edu)
STAT/CS 5810/6655



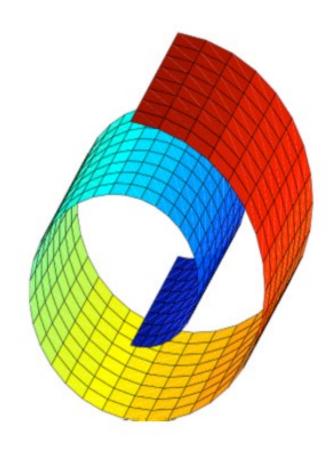
#### Outline



- 1. Manifold Learning
- 2. Visualizing with PCA
- 3. t-SNE
- 4. UMAP
- 5. PHATE
- 6. DIG

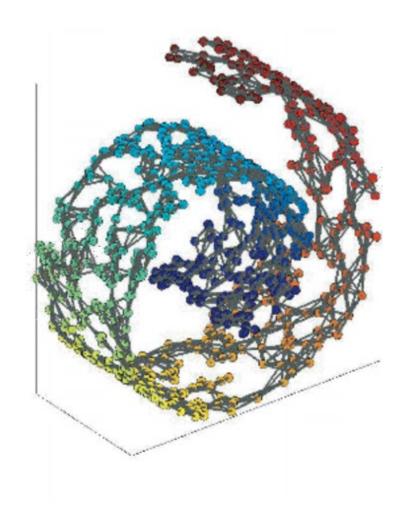
#### Manifold learning





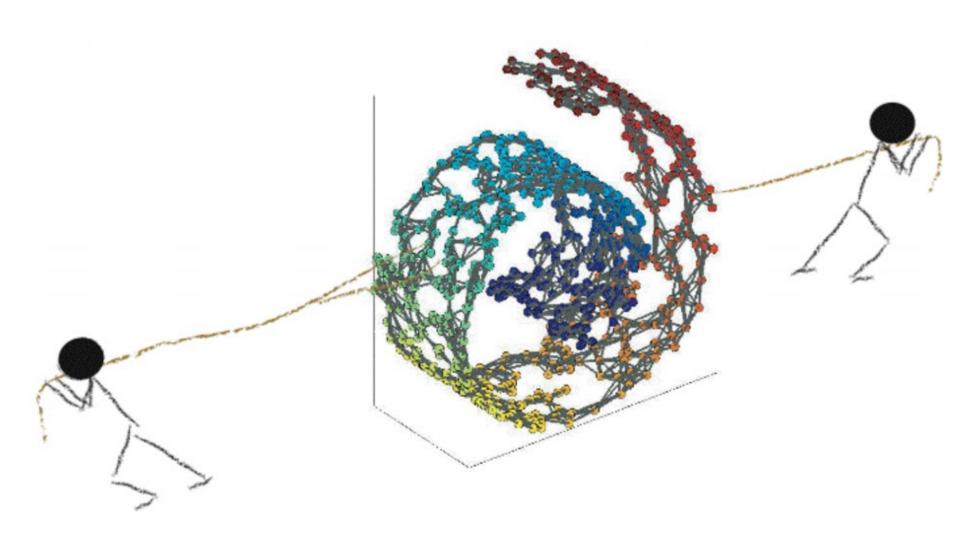
### Manifold learning





### Manifold learning





# Manifold Learning for Data Visualization

#### Data Visualization

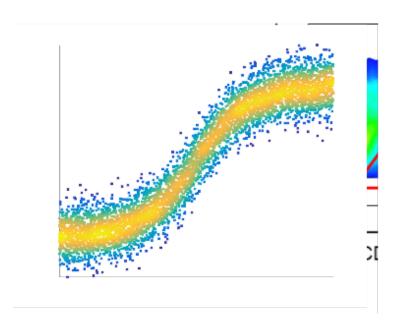


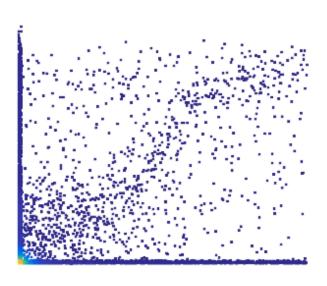
- Humans are very visual
- Data visualization is a necessary tool for exploration
  - Develop intuitive understanding of the structure
  - Generate hypotheses

#### Challenges in Visualizing (Biomedical) Data



- High dimensions
- Noise/artifacts
- Nonlinearities
  - Most biological processes are NOT linear





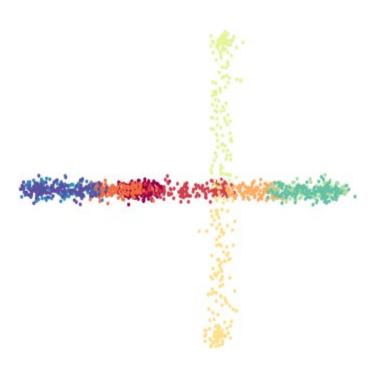
#### Visualization with PCA





1400 points, 60 dimensions

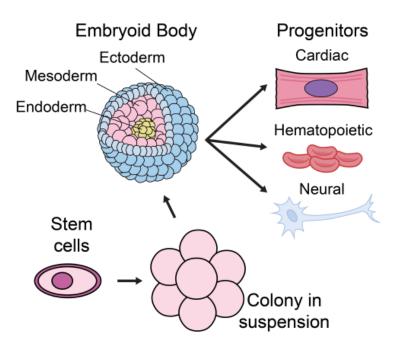
PCA



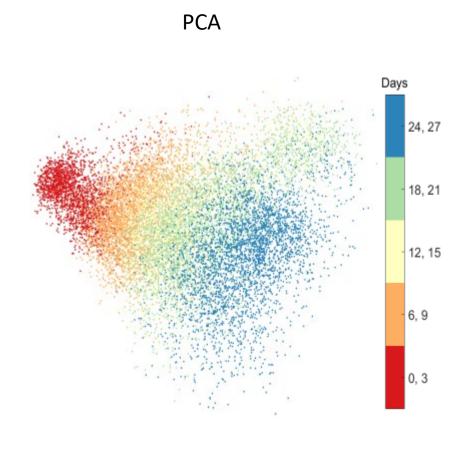
#### Visualization with PCA



Newly generated scRNA-seq data (27 days)



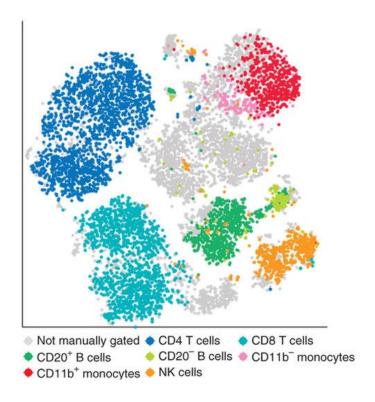
 $\approx$ 31k cells, more than 17k genes



#### t-SNE (van der Maaten & Hinton, JMLR, 2008)



- t-distributed stochastic neighbor embedding (t-SNE)
- Widely used on single cell data
  - Biology version (Amir et al, Nature Biotech, 2013) has 1400+ citations
- Attempts to preserve local relationships in both the high and low-dimensional spaces
  - Designed for separating clusters



(Amir et al, 2013)

### t-SNE



**Basic idea**: neighbors in the high-dimensional space should be neighbors in low-dimensions

Compute affinities/probabilities from distances:

$$p_{j|i} = rac{\exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2/2\sigma_i^2)}{\sum_{k 
eq i} \exp(-\|\mathbf{x}_i - \mathbf{x}_k\|^2/2\sigma_i^2)}$$

Symmetrize:

$$p_{ij} = rac{p_{j|i} + p_{i|j}}{2N}$$

Probabilities in low dimensions:

$$q_{ij} = rac{(1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2)^{-1}}{\sum_k \sum_{l 
eq k} (1 + \|\mathbf{y}_k - \mathbf{y}_l\|^2)^{-1}}$$





Compute affinities/probabilities from distances:

$$p_{j|i} = rac{\exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2/2\sigma_i^2)}{\sum_{k 
eq i} \exp(-\|\mathbf{x}_i - \mathbf{x}_k\|^2/2\sigma_i^2)}$$

• Symmetrize:

$$p_{ij} = rac{p_{j|i} + p_{i|j}}{2N}$$

Probabilities in low dimensions:

$$q_{ij} = rac{(1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2)^{-1}}{\sum_k \sum_{l 
eq k} (1 + \|\mathbf{y}_k - \mathbf{y}_l\|^2)^{-1}}$$

Minimize KL divergence between them:

$$ext{KL}\left(P \parallel Q
ight) = \sum_{i 
eq j} p_{ij} \log rac{p_{ij}}{q_{ij}}$$

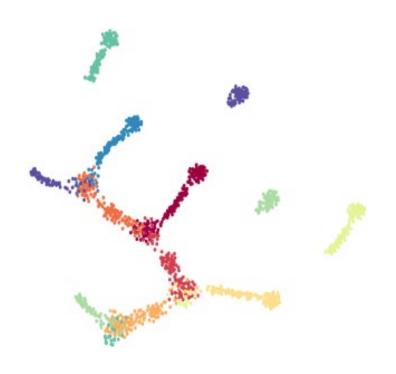
#### Visualization with t-SNE





1400 points, 60 dimensions

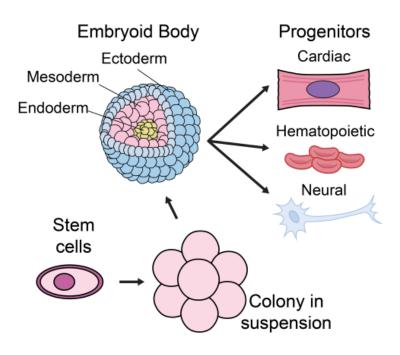
t-SNE



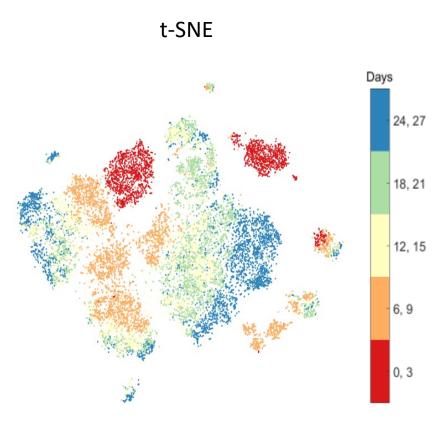
#### Visualization with t-SNE



Newly generated scRNA-seq data (27 days)



≈31k cells, more than 17k genes



#### Issues with t-SNE



- Sensitive to hyperparameters
- Cluster sizes in a t-SNE plot are meaningless
  - Distance adapts to density
- Distances between clusters are generally meaningless
  - Global relationships are not preserved
- Random noise may not look random
- See <a href="https://distill.pub/2016/misread-tsne/">https://distill.pub/2016/misread-tsne/</a> for more details

#### **UMAP**



- UMAP was proposed in 2018 to counter t-SNE's issues with preserving global structure
- It's also faster than t-SNE
- UMAP does do better at global structure, but not for the reasons the original authors claimed

#### UMAP vs. t-SNE



#### **T-SNE**

- Uses "perplexity" to determine similarity scale
- Normalizes similarities
- Averages similarities to symmetrize
- Uses KL-divergence as a loss function
- Initializes randomly
- Considers most pairwise similarities

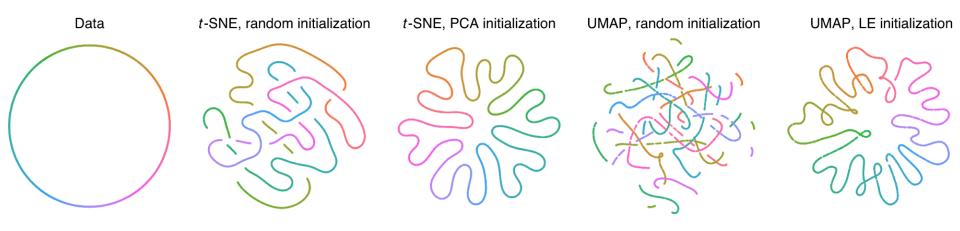
#### **UMAP**

- Uses k-nn distances to determine similarity scale
- Doesn't normalize
- Symmetrizes similarities differently
- Uses cross entropy as a loss function
- Initializes with Laplacian Eigenmaps
- Uses noise contrastive estimation (NCE)

#### Importance of Initialization



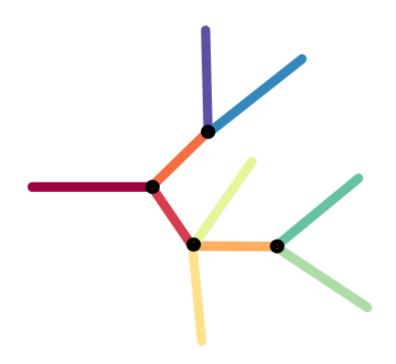
- The main reason UMAP does better at global relationships than t-SNE is the Laplacian eigenmaps initialization
- Initializing t-SNE with PCA gives near identical results (Kobak and Linderman, 2021)



## UMAP still fails at global structure

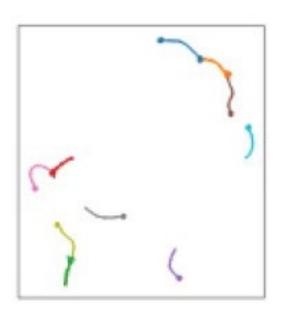


**Artificial Tree Data** 



1400 points, 60 dimensions

UMAP



#### Why is UMAP faster?

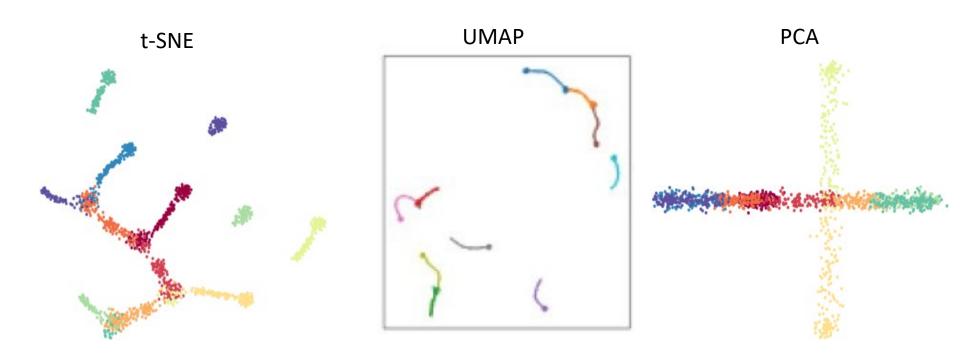


- Loss function is faster to compute and differentiate (no similarity normalizations)
- Biggest gain is in NCE
- NCE compares points to only a few real points as well as a bunch of "fake" points
  - Allows for faster computations
  - Used in a lot of neural network settings to speed up training as part of "self-supervised learning"
  - However, it changes the effective loss function (Damrich and Hamprecht, 2021)

#### Balancing local and global information



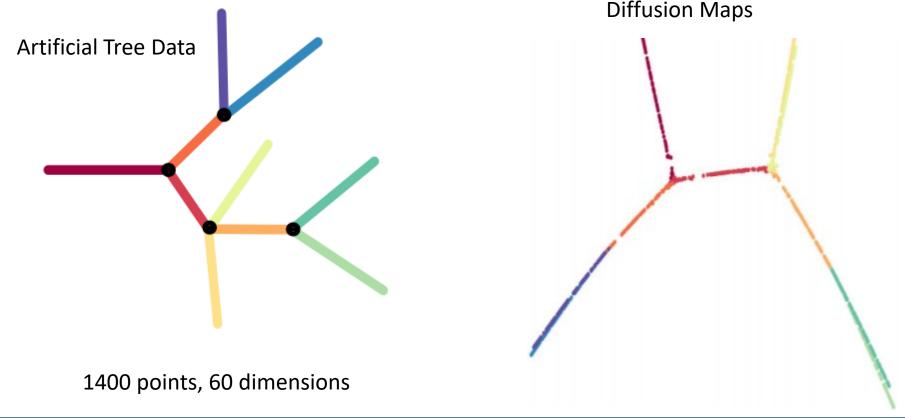
- T-SNE and UMAP are good for local structure
- PCA is good for global structure
- Diffusion maps (DM) does well at capturing both
- Idea: use DM to visualize data



#### Visualization with Diffusion Maps



- DM captures the information accurately and robustly, but not for visualization
- DM tends to place the structure in higher dimensions
- DM can also be unstable with boundaries and intersections



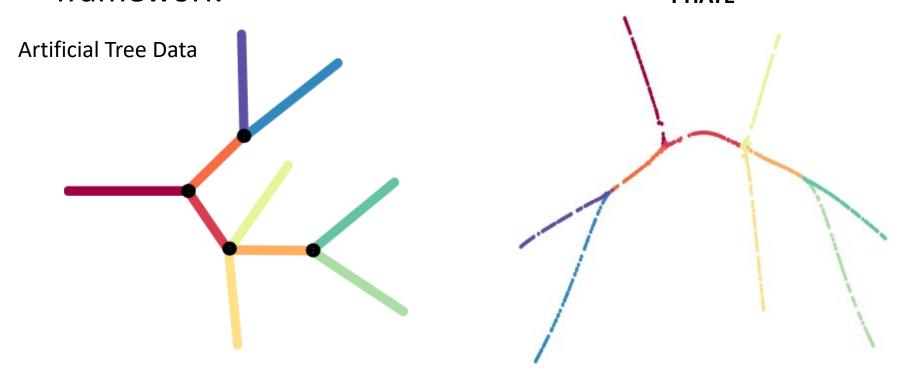
PHAIL (Potential of Heat-diffusion for Affinity-based Transition



Embedding)

#### Visualizing Structure and Transitions in High-Dimensional Biological Data (Nature Biotechnology, 2019)

 A method for visualizing data built off of the diffusion framework **PHATE** 

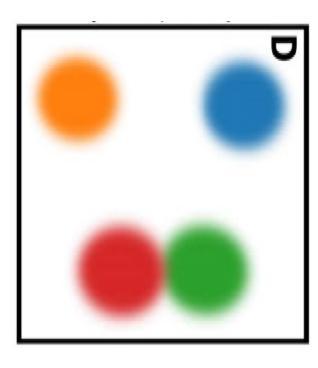


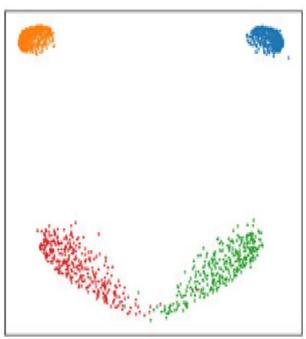
1400 points, 60 dimensions

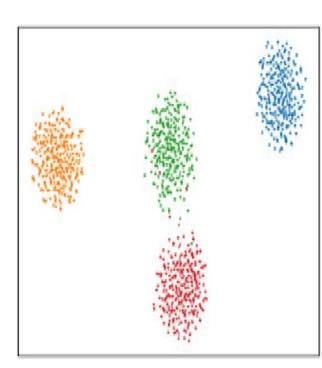
#### PHATE on GMM



Truth PHATE t-SNE

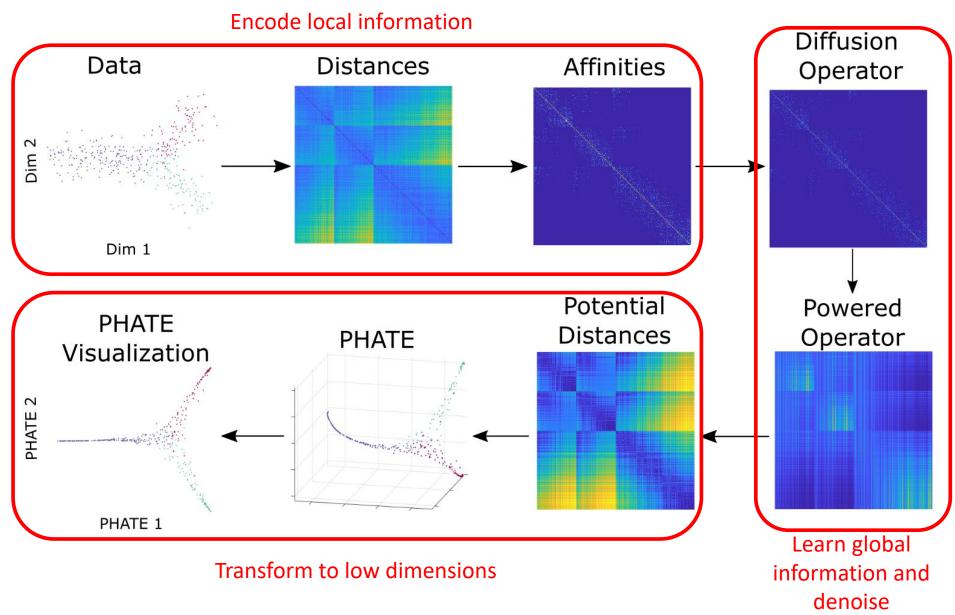






### The PHATE Algorithm





#### PHATE on Frey Faces



- Frey Faces dataset
  - 1965 frames taken from a video of Brendan Frey making various faces in front of a camera
  - Resolution: 20x28
- Frames are given out of order to PHATE w/o information about sequential ordering











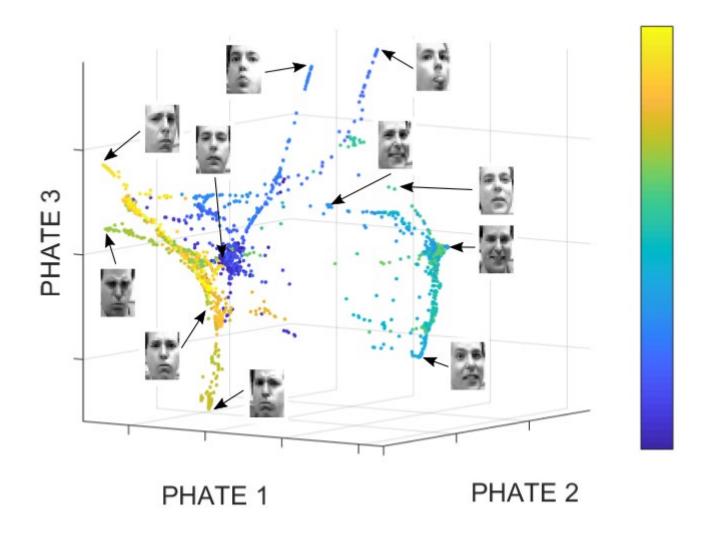






## PHATE on Frey Faces

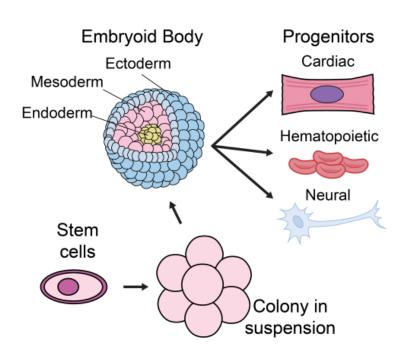




#### PHATE on EB scRNA-seq Data

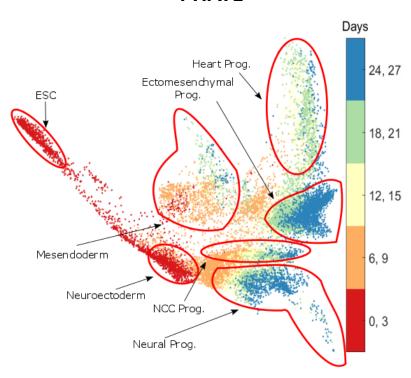


Newly generated scRNA-seq data (27 days)



≈31k cells, more than 17k genes

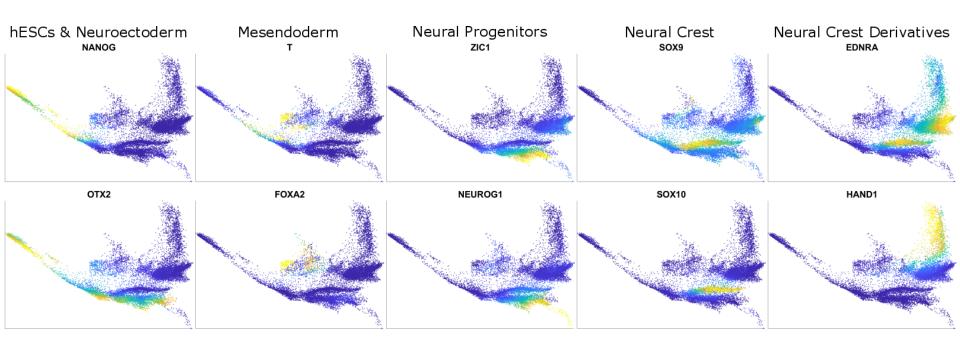
#### **PHATE**



#### Exploratory Data Analysis with PHATE



 Coloring the embedding by gene expression after MAGIC\* (van Dijk,...Moon et al., 2018) reveals lineages

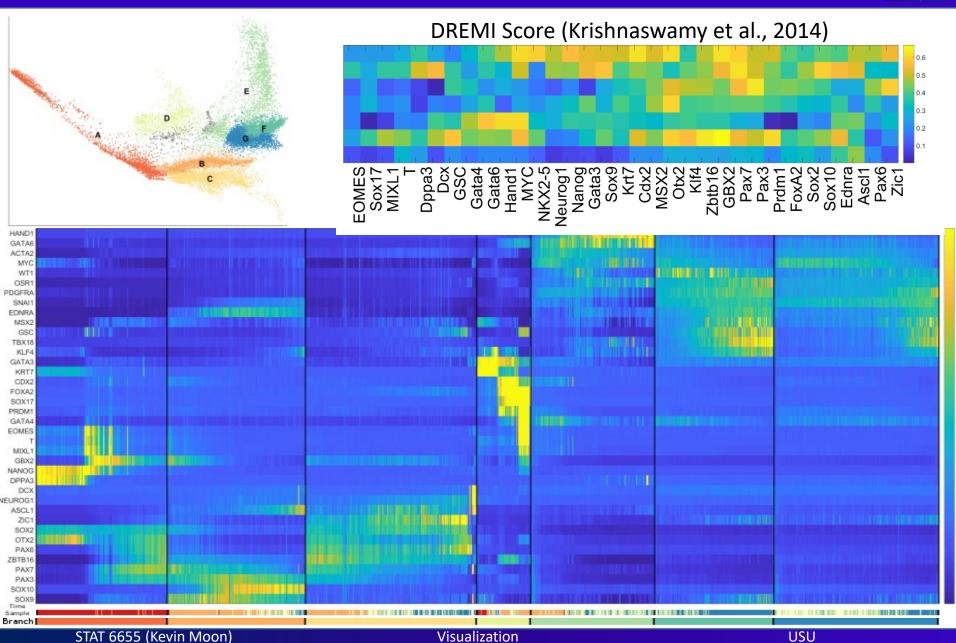


• Interactive web tool: krishnaswamylab.org/phatewebtool

\*Published in Cell

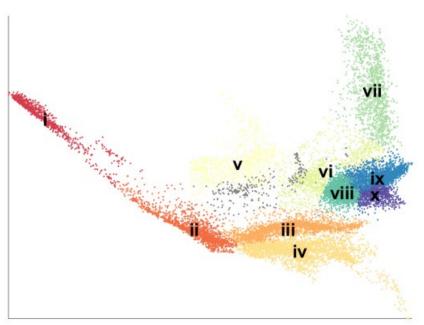
#### **Exploratory Data Analysis with PHATE**

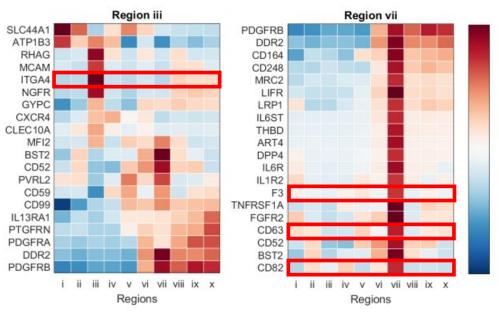


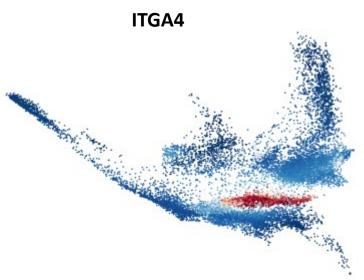


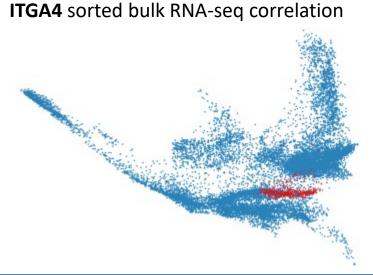
#### Discovering New Surface Markers for Sorting Populations





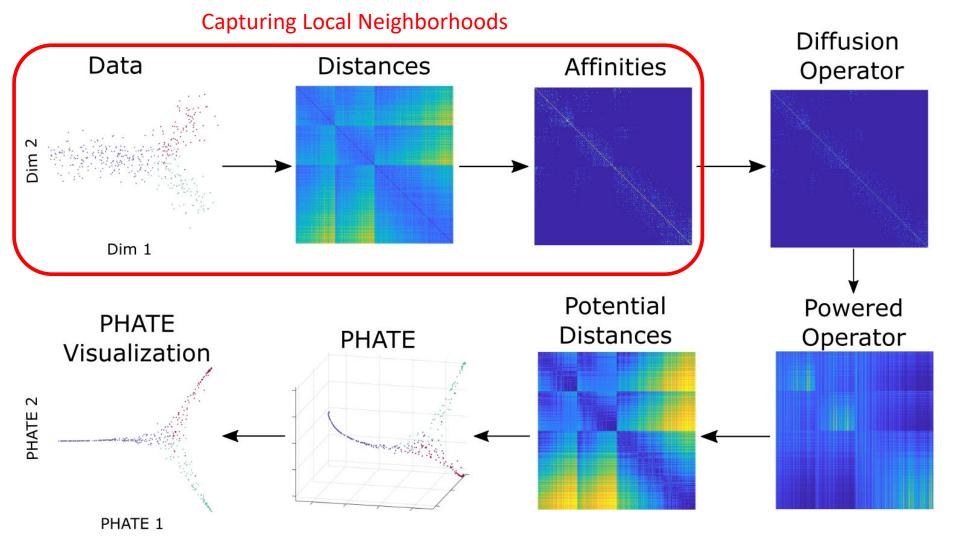






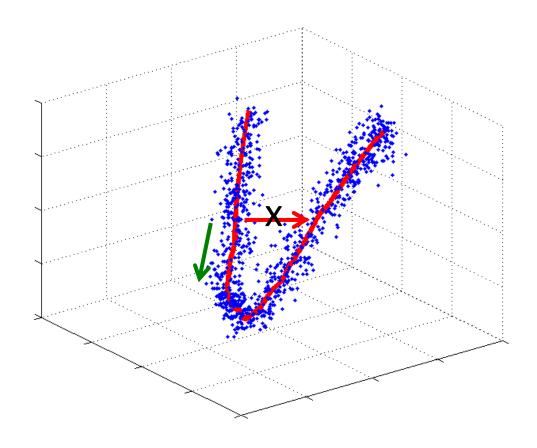
### The PHATE Algorithm





#### Capturing local neighborhoods





- Large Euclidean distance gives wrong global structure for visualizing
- Small Euclidean distances ok
  - Encodes local structure

#### From distances to affinities



- Step 1: Calculate all pairwise Euclidean distances
- Step 2: Convert the distances to pairwise affinities using a kernel function
  - E.g. the Gaussian kernel
  - Kernel function must be near zero for large distances and nonzero for small distances
  - Affinity: a measure of similarity between points
- Problem: many data have both dense and sparse regions
- A kernel bandwidth fixed for dense regions doesn't work well for sparse regions and vice versa
- Idea: use locally adaptive bandwidth

### Capturing Local Neighborhoods



• Adaptive 
$$\alpha$$
-decaying kernel  $\tilde{g}(x,y) = \exp\left(-\left(\frac{||x-y||}{\epsilon_x}\right)^{\alpha}\right), \rightarrow g(x,y) = \frac{\tilde{g}(x,y) + \tilde{g}(y,x)}{2}$ 

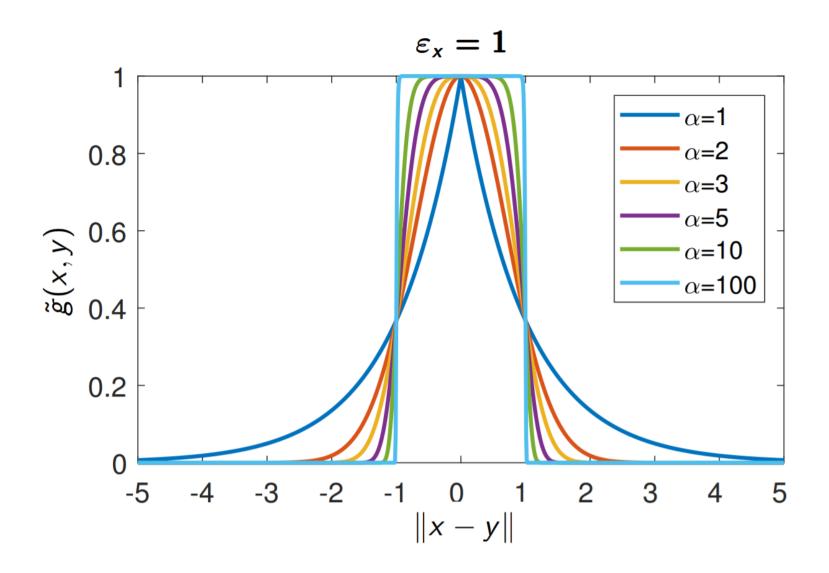
#### Where

- $\epsilon_x$  = distance from x to its kth nearest neighbor
- $\alpha$  controls the decay rate of  $\tilde{g}$

Provides a robust notion of adaptive locality

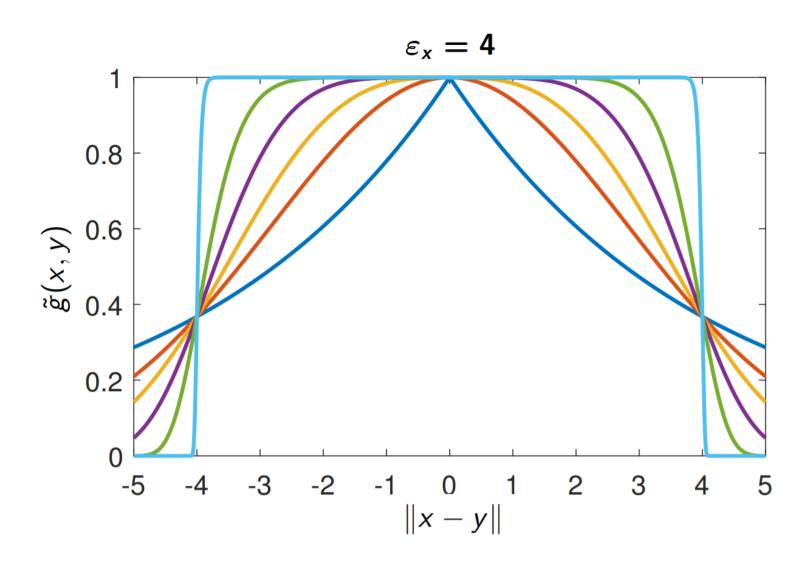
# Capturing Local Neighborhoods





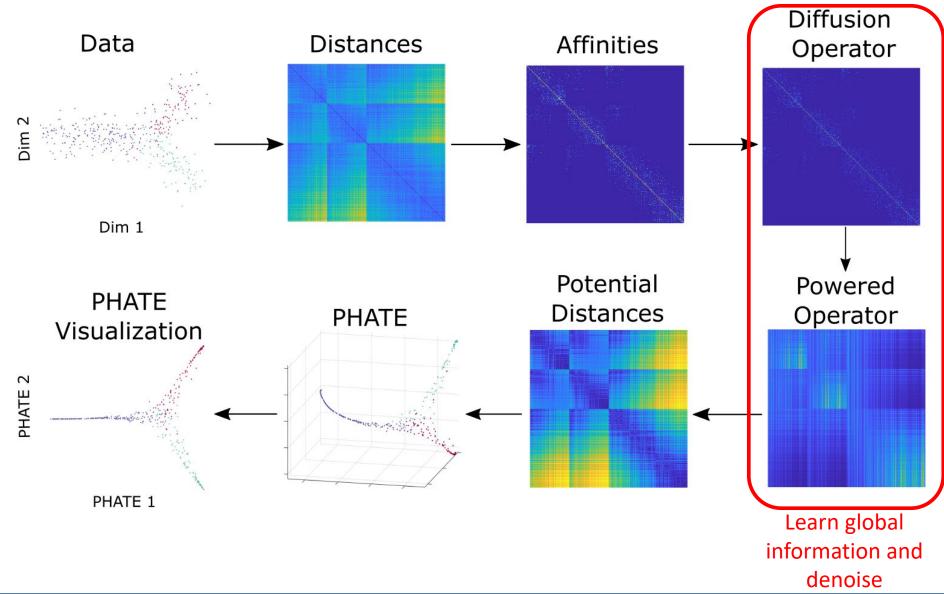
# Capturing Local Neighborhoods





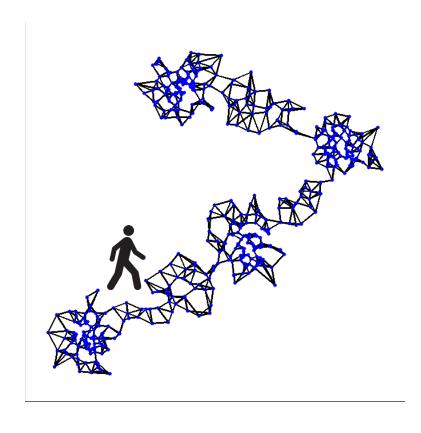
# The PHATE Algorithm





#### Diffusion Denoises and Recovers Global Structure





- Big Euclidean steps bad, likely to exit structure
- Small steps good, likely to stay within structure
- To learn global structure via small local steps we use random walk (diffusion)
  - Normalize affinity matrix to create Markov transition matrix (the <u>diffusion operator</u>) P (Coifman & Lafon, ACHA, 2006)
  - Power *P* by time step *t*

#### How much diffusion?



• Von Neumann entropy (von Neumann, 1932) of diffused operator  $P^t$ 

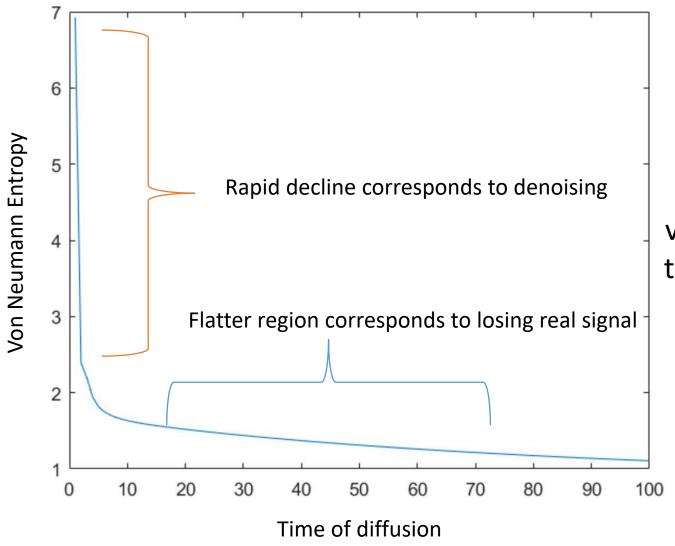
$$VNE(P^t) = -\sum_{j} \eta_j \log \eta_j$$
,  $\eta_j = \lambda_j^t / \left| |\lambda^t| \right|_1$ 

Where  $\lambda^t = {\lambda_0^t, \lambda_1^t, ...}$  are the eigenvalues of  $P^t$ 

- VNE is a soft proxy of numerical rank
- Decays as diffusion time increases,  $\lim_{t\to\infty} VNE(P^t) = 0$
- Use rate of decay to choose time scale

#### How much diffusion?

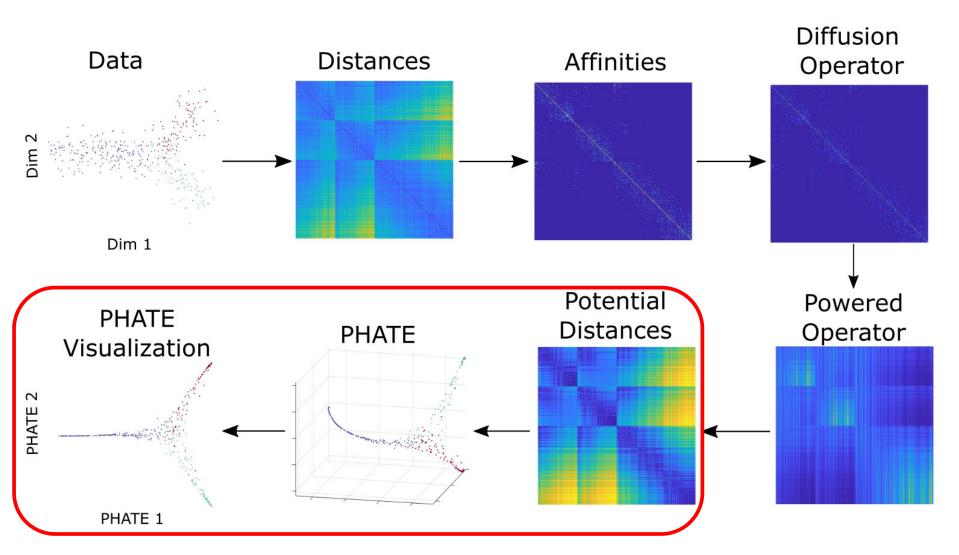




Embedding is visually robust in the flatter region

# The PHATE Algorithm



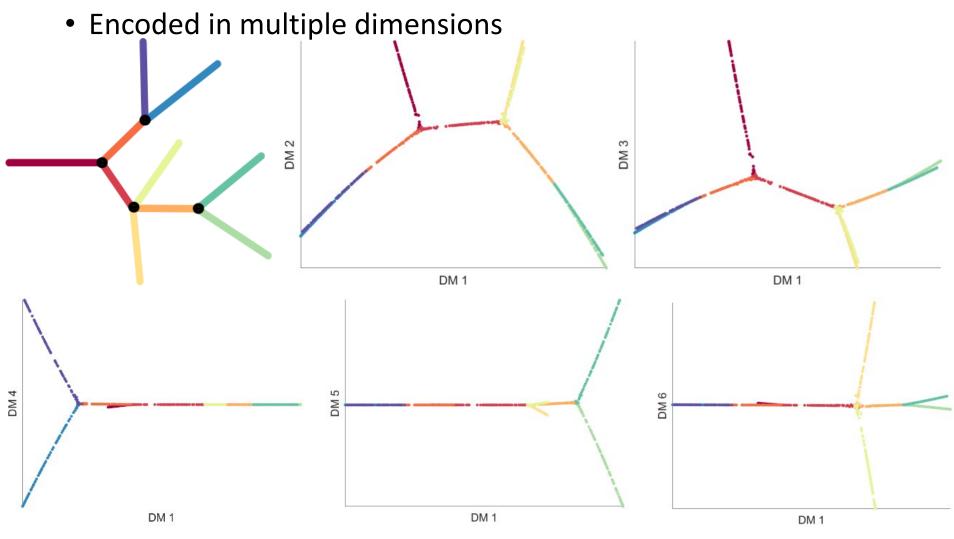


Transform to low dimensions

#### Potential Distances



Diffusion operator contains global and local structure

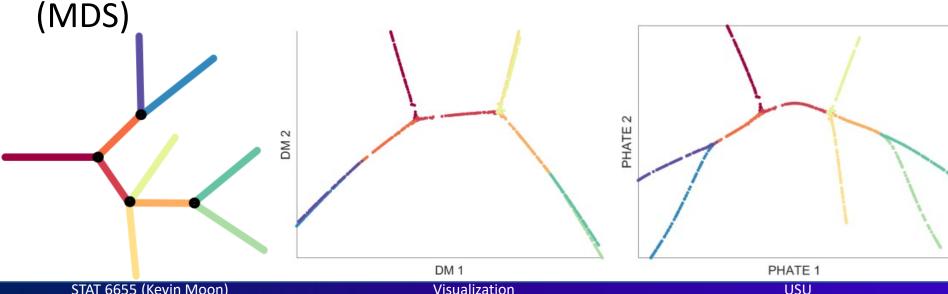


#### Potential Distances



- Diffusion operator contains global and local structure
  - Encoded in multiple dimensions
- To extract this information, we transform diffused probabilities using a potential transformation
  - Forms an information distance between diffused probabilities
  - Connected to heat potential

• Embed for visualization using multidimensional scaling



### Full PHATE Algorithm



**Input:** Data matrix X, neighborhood size k, locality scale  $\alpha$ , desired embedding dimension m (usually 2 or 3 for visualization)

**Output:** The PHATE embedding  $Y_m$ 

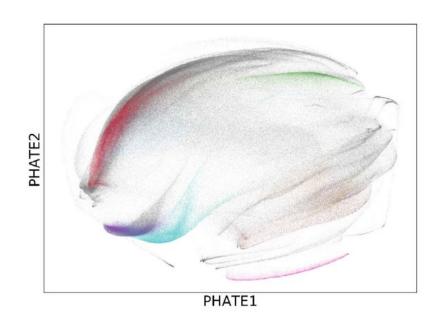
- 1:  $D \leftarrow$  compute pairwise distance matrix from X
- 2: Compute the k-nearest neighbor distance  $\varepsilon_k(x)$  for each column x of X
- 3:  $K_{k,\alpha} \leftarrow$  compute local affinity matrix from D and  $\varepsilon_k$  (see Eq. 3)
- 4:  $P \leftarrow$  normalize  $K_{k,\alpha}$  to form a Markov transition matrix (diffusion operator; see Eq. 2)
- 5:  $t \leftarrow$  compute time scale via Von Neumann Entropy (see Eq. 5)
- 6: Diffuse P for t time steps to obtain  $P^t$
- 7: Compute potential representations:  $U_t \leftarrow -\log(P^t)$
- 8:  $\mathfrak{V}^t \leftarrow$  compute potential distance matrix from  $U_t$  (see Eq. 6)
- 9:  $Y_{class} \leftarrow$  apply classical MDS to  $\mathfrak{V}^t$
- 10:  $Y_m \leftarrow$  apply metric MDS to  $\mathfrak{V}^t$  with  $Y_{class}$  as an initialization

Supplemental Table S1: Detailed steps in the PHATE algorithm.

# Scalability of PHATE



- Storing and performing operations on the diffusion matrix can be difficult for large samples
- Can reduce computation by diffusing through "landmarks" (compressed diffusion, Gigante et al., 2018)
  - Landmarks are chosen by clustering
  - Obtain an embedding of all points by projection

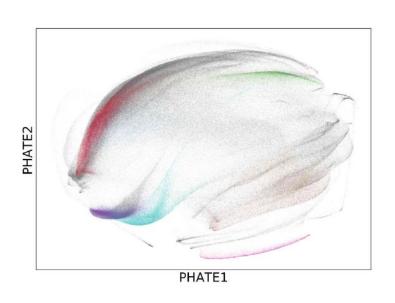


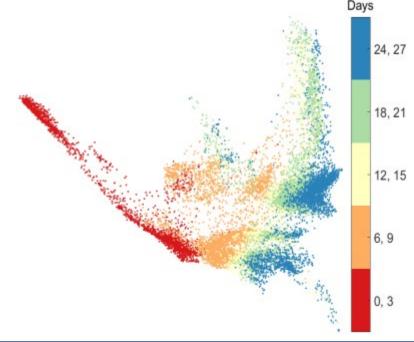
PHATE applied to the 10x megacell mouse brain data (>1 million cells)

#### Noisy Visualizations



- The visualizations are still pretty noisy when the data are noisy
- Can we just diffuse more to denoise more?
- Increasing the amount of diffusion can oversmooth
  - Loss of true signal
- Can we do better than this?

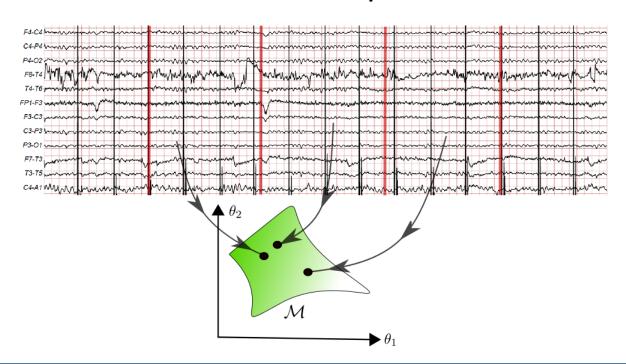




# Dynamical Systems Approach



- Let's add some more structural assumptions
- Let's assume that the data are generated by a dynamical system (e.g. EEG measurements)
  - I.e., we now have a time component to the data
- Goal: learn a low-dimensional representation of the data



# Diffusion with Dynamical Systems



State-space formalism:

$$\begin{aligned} \boldsymbol{x}_t &= \boldsymbol{y}_t(\boldsymbol{\theta}_t) + \boldsymbol{\xi}_t \\ d\theta_t^i &= a^i (\theta_t^i) dt + dw_t^i, & i &= 1, \dots, m \end{aligned}$$

- $x_t$  is the observed time series while  $\theta_t$  represents the unobserved states that drive the process
- $\xi_t$  is a stationary process independent of  $y_t$  (noise)
- $m{x}_t$  is a corrupted version of a clean process  $m{y}_t$  that is driven by  $m{ heta}_t$
- The unknown drift functions  $a^i$  are independent of  $\theta^j$  when  $j \neq i$ 
  - $\Rightarrow$  we assume local independence between  $\theta_t^i$  and  $\theta_t^j$  when  $j \neq i$
- $w_t^i$  are independent white noise
- Can we recover  $\boldsymbol{\theta}_t$  from  $\boldsymbol{x}_t$ ?

# Diffusion with Dynamical Systems



State-space formalism:

$$\begin{aligned} \boldsymbol{x}_t &= \boldsymbol{y}_t(\boldsymbol{\theta}_t) + \boldsymbol{\xi}_t \\ d\boldsymbol{\theta}_t^i &= a^i (\boldsymbol{\theta}_t^i) dt + d\boldsymbol{w}_t^i, & i = 1, \dots, m \end{aligned}$$

- Can view  $\mathbf{y}_t$  as being drawn from the conditional pdf  $p(\mathbf{y}|\boldsymbol{\theta})$
- Important insight 1 (Talmon & Coifman, 2015): the pdf  $p(x|\theta)$  is a linear transformation of  $p(y|\theta)$
- New feature space: histograms  $m{h}_t$  taken within a time window centered at  $m{x}_t$
- Important insight 2: the expected value of the histograms, e.g.  $\mathbb{E}[h_t]$ , is a linear transformation of  $p(x|\theta)$

## Diffusion with Dynamical Systems



- Important insight 1 (Talmon & Coifman, 2015): the pdf  $p(x|\theta)$  is a linear transformation of  $p(y|\theta)$
- Important insight 2: the expected value of the histograms, e.g.  $\mathbb{E}[\boldsymbol{h}_t]$ , is a linear transformation of  $p(\boldsymbol{x}|\boldsymbol{\theta})$
- Important insight 3: the Mahalanobis distance is invariant under linear transformations
- Therefore, the following distance is noise resilient (Talmon & Coifman, 2015):

$$D^{2}(\boldsymbol{x}_{t},\boldsymbol{x}_{s}) = (\mathbb{E}[\boldsymbol{h}_{t}] - \mathbb{E}[\boldsymbol{h}_{s}])^{T}C_{t}^{-1}(\mathbb{E}[\boldsymbol{h}_{t}] - \mathbb{E}[\boldsymbol{h}_{s}])$$

• I.e., under certain assumptions,

$$\|\boldsymbol{\theta}_t - \boldsymbol{\theta}_S\|^2 \approx D^2(\boldsymbol{x}_t, \boldsymbol{x}_S)$$

## Dynamical Information Geometry (Duque et al., 2019)



Compute the Mahalanobis between expected values:

$$d^{2}(\boldsymbol{z}_{t},\boldsymbol{z}_{s})=(\mathbb{E}(\boldsymbol{h}_{t})-\mathbb{E}(\boldsymbol{h}_{s}))^{T}(\boldsymbol{C}_{t}+\boldsymbol{C}_{s})^{-1}(\mathbb{E}(\boldsymbol{h}_{t})-\mathbb{E}(\boldsymbol{h}_{s})),$$

- ullet Expected values and covariance matrices computed in a time window of length  $L_2$
- Input these distances into diffusion maps to get EIG (Talmon and Coifman, 2015)
  - A noise-resilient dimensionality reduction method
- Input these distances into PHATE to get DIG
  - A noise resilient visualization method

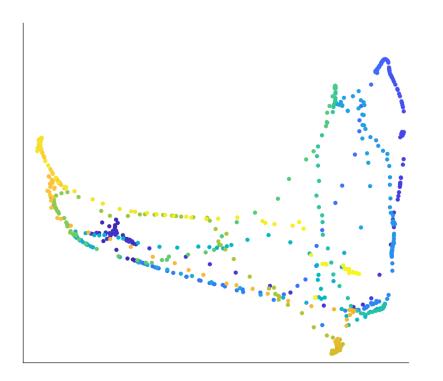
#### **EEG** Results: Information Distances

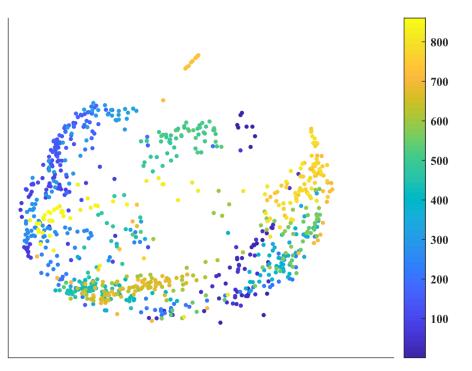


- Applied DIG to EEG sleep data (Terzano et al., 2002; Goldberger et al., 2000)
- Visualizations colored by time

DIG (Mahalanobis Distance)

An alternative geodesic distance





#### Information Distances



 We also explored other information distances (applied to the diffusion operator) besides the potential distance

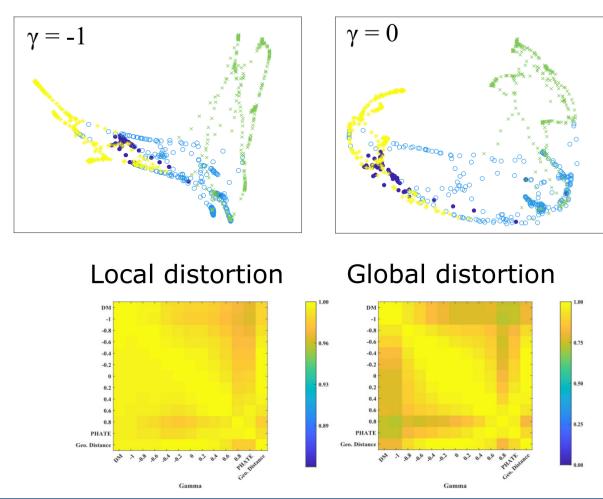
$$D_{\gamma,t}^{2}(\boldsymbol{x}_{i},\boldsymbol{x}_{j}) = \begin{cases} \sum\limits_{k=1}^{N} \frac{(\log P_{ki}^{t} - \log P_{kj}^{t})^{2}}{\phi_{0}(k)}, & \gamma = 1 \\ \sum\limits_{k=1}^{N} \frac{(P_{ki}^{t} - P_{kj}^{t})^{2}}{\phi_{0}(k)}, & \gamma = -1 \\ \sum\limits_{k=1}^{N} \frac{2((P_{ki}^{t})^{\frac{1-\gamma}{2}} - (P_{kj}^{t})^{\frac{1-\gamma}{2}})^{2}}{(1-\gamma)\phi_{0}(k)}, & -1 < \gamma < 1. \end{cases}$$

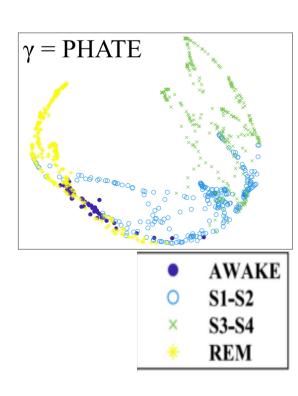
 Top corresponds to potential distance (PHATE), middle to diffusion maps

# **EEG Results: Information Distances**



Visualizations colored by sleep stage

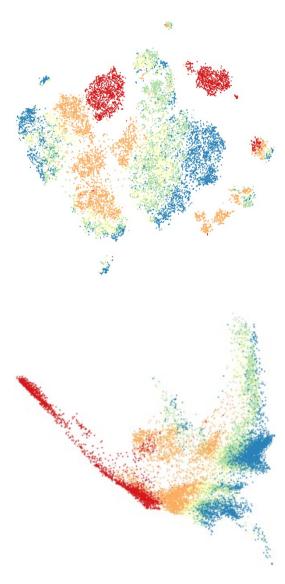




#### PHATE Summary



- Data have structure at different scales
  - Local branching structure
  - Global relationship between branches
- Existing visualization methods fail to account for all scales
- PHATE captures both global and local structure
- PHATE reveals new biology
- DIG extends this to better denoise the data



### Further reading



#### PHATE

- Paper: <a href="https://doi.org/10.1101/120378">https://doi.org/10.1101/120378</a>
- Code: <a href="https://github.com/KrishnaswamyLab/PHATE">https://github.com/KrishnaswamyLab/PHATE</a>

#### • DIG

- Paper: <a href="https://doi.org/10.1109/MLSP.2019.8918875">https://doi.org/10.1109/MLSP.2019.8918875</a>
- Code: https://github.com/KevinMoonLab/DIG

#### T-SNE

- https://lvdmaaten.github.io/tsne/
- https://distill.pub/2016/misread-tsne/

#### UMAP

- Damrich and Hamprecht (2021): https://openreview.net/pdf?id=DKRcikndMGC
- Kobak and Linderman (2021): https://www.nature.com/articles/s41587-020-00809-z