



University of Wisconsin  
**SCHOOL OF MEDICINE  
AND PUBLIC HEALTH**

# Clustering

Moo K. Chung

University of Wisconsin-Madison

[www.stat.wisc.edu/~mchung](http://www.stat.wisc.edu/~mchung) →

[github.com/laplcebeltrami](https://github.com/laplcebeltrami)

Motivating  
problem

# Magnetic resonance imaging (MRI)

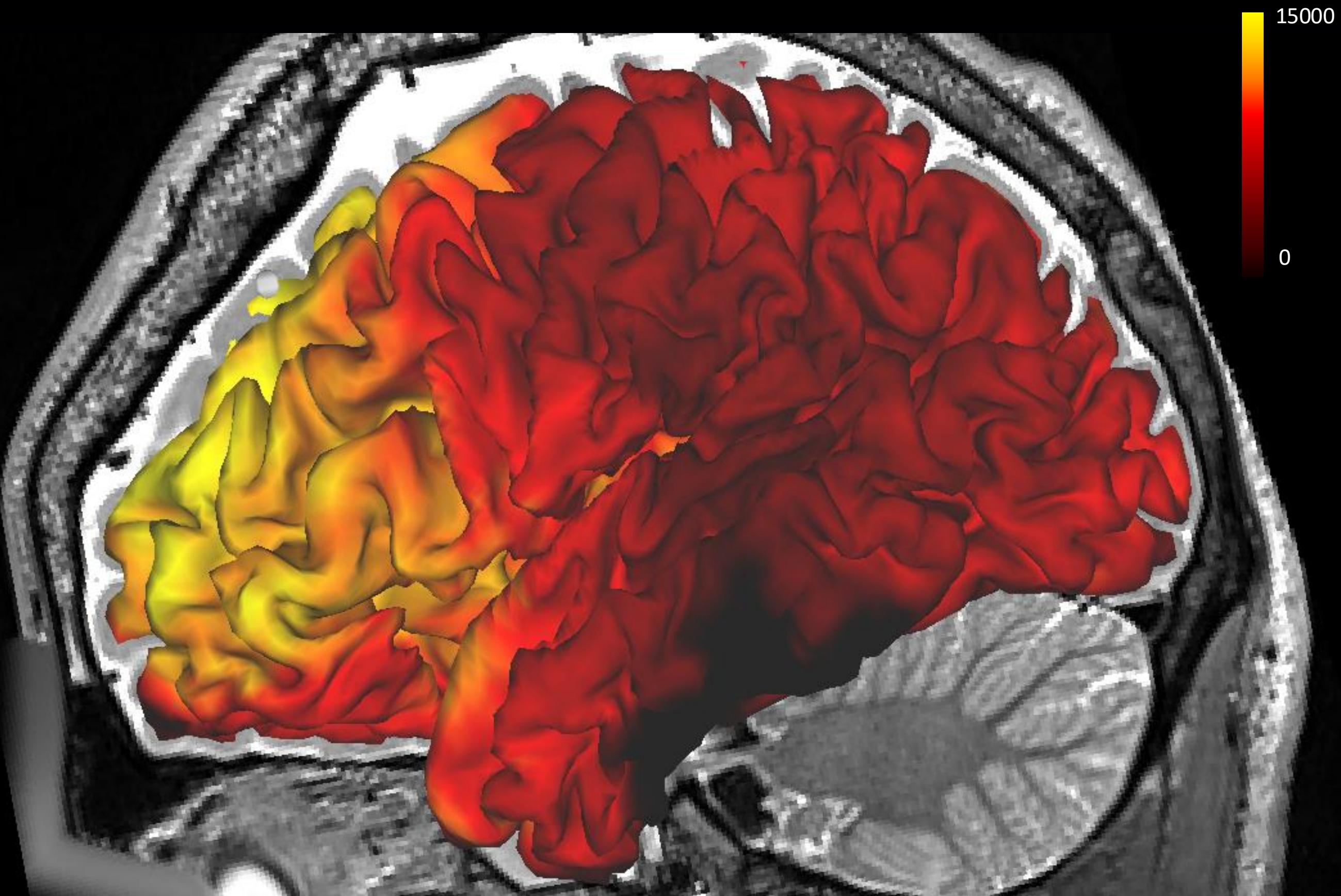


3T GE Discovery X750  
Waisman Brain Imaging Laboratory  
University of Wisconsin-Madison



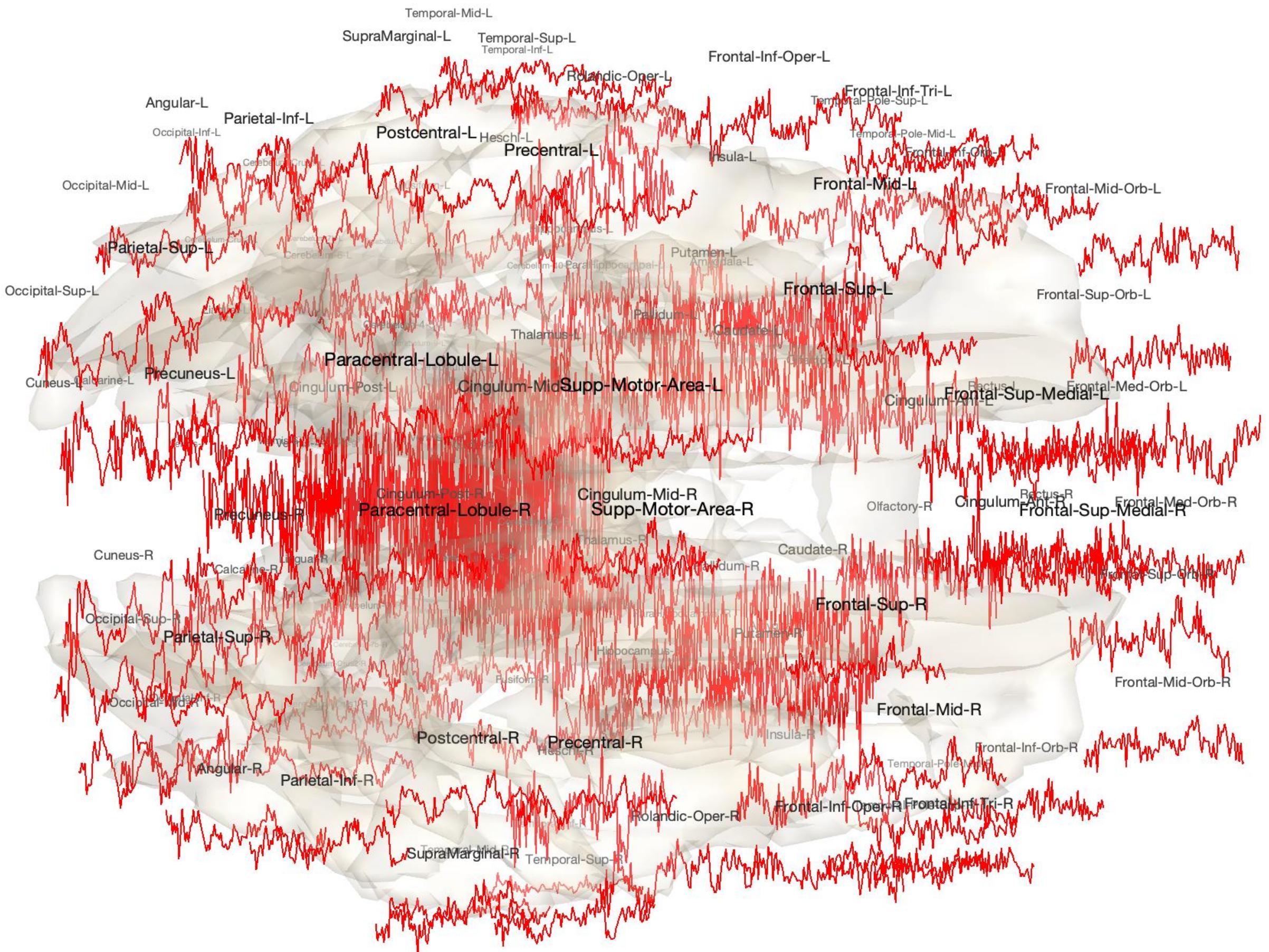
3T GE Discovery MR750  
Center for Imaging Research  
Medical College of  
Wisconsin, Milwaukee, WI

rs-fMRI (every 30 second)



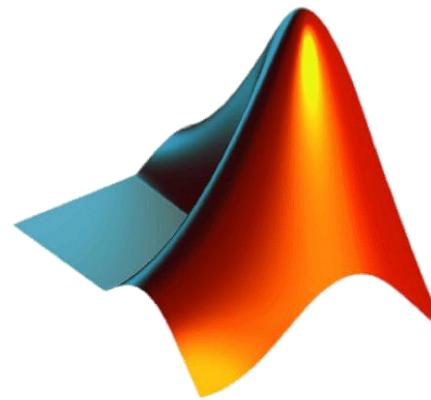
Time series averaged into 116 brain

1



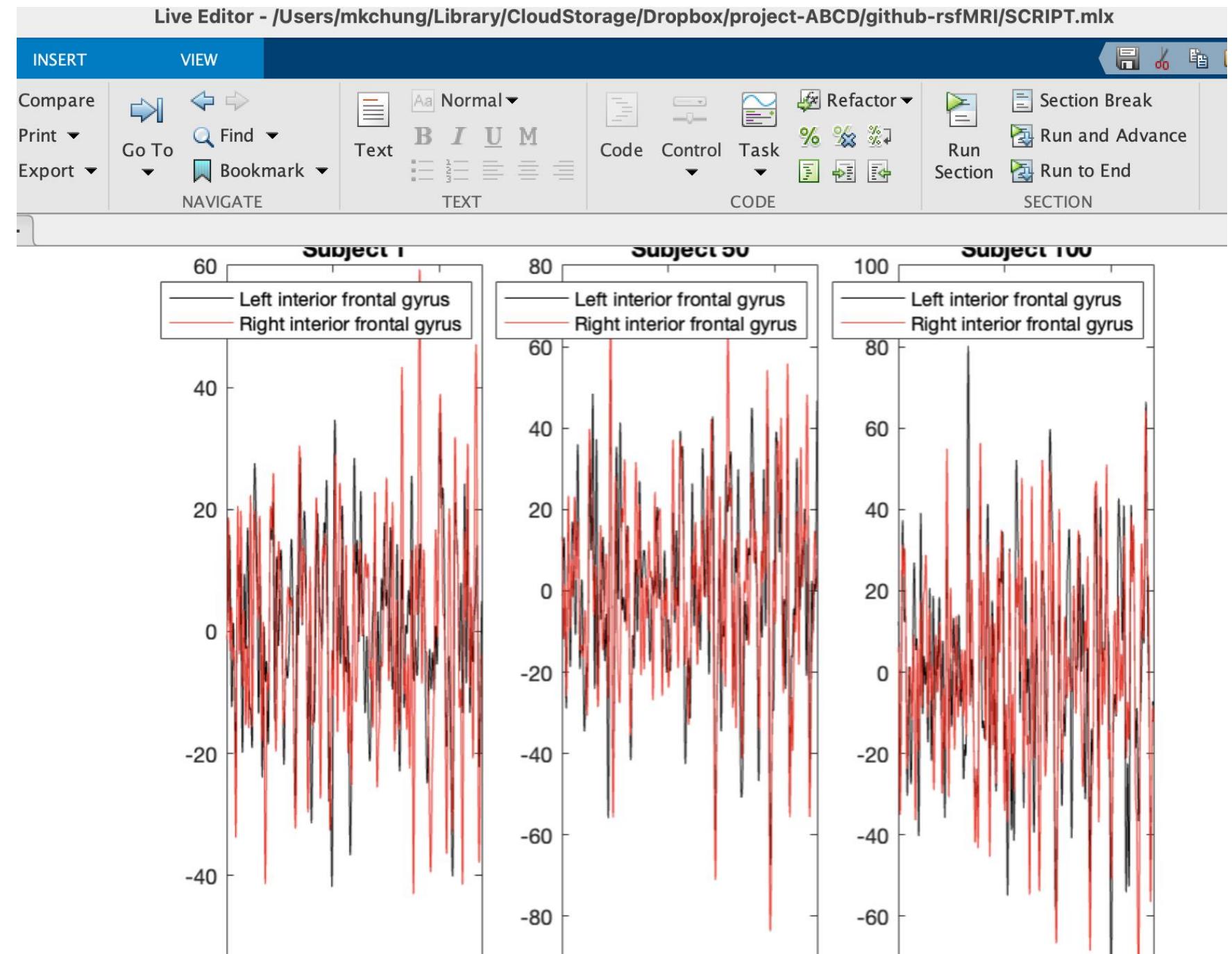
# Sample rs-fMRI time series data

## rsfMRI

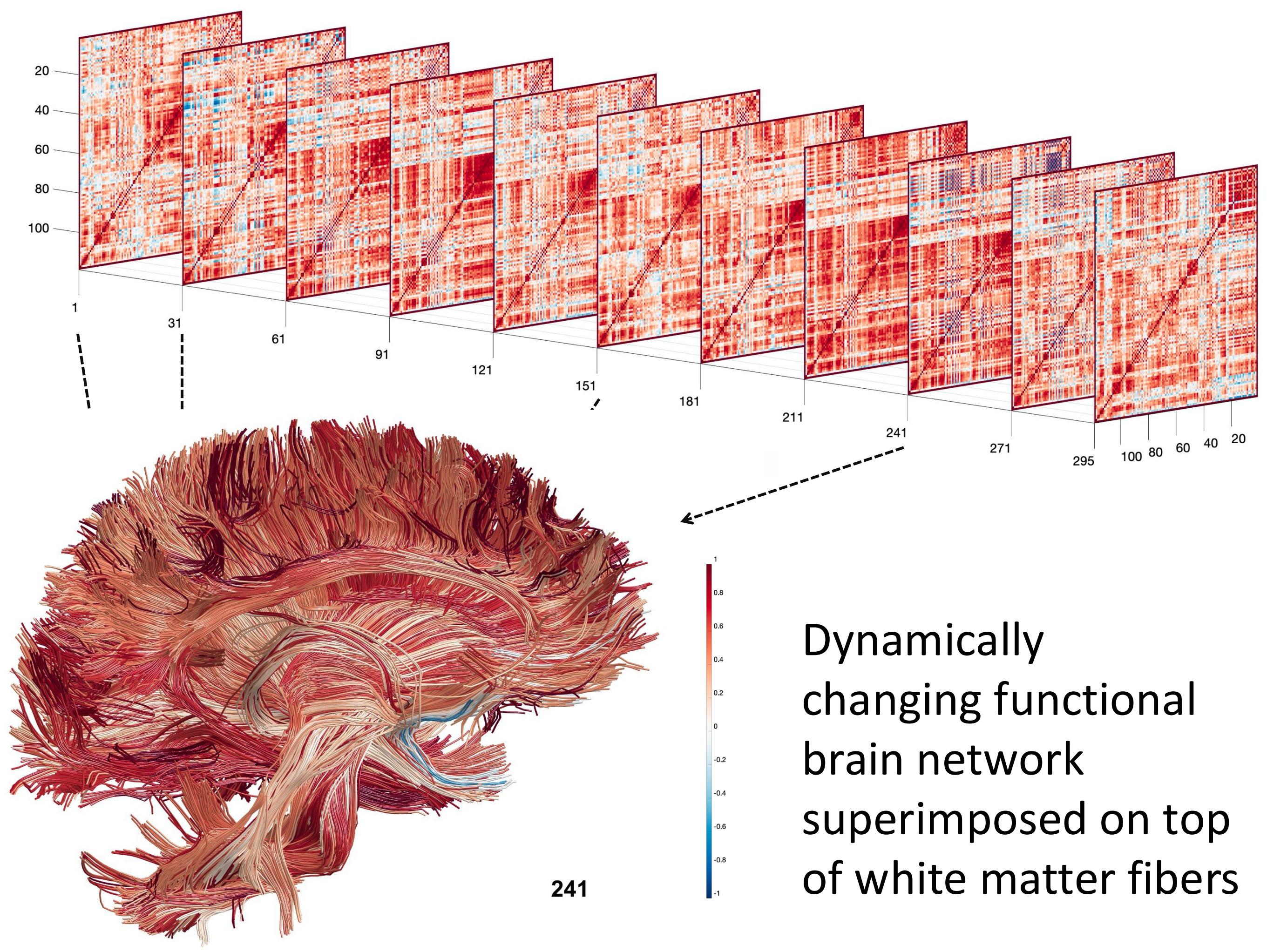


MATLAB®

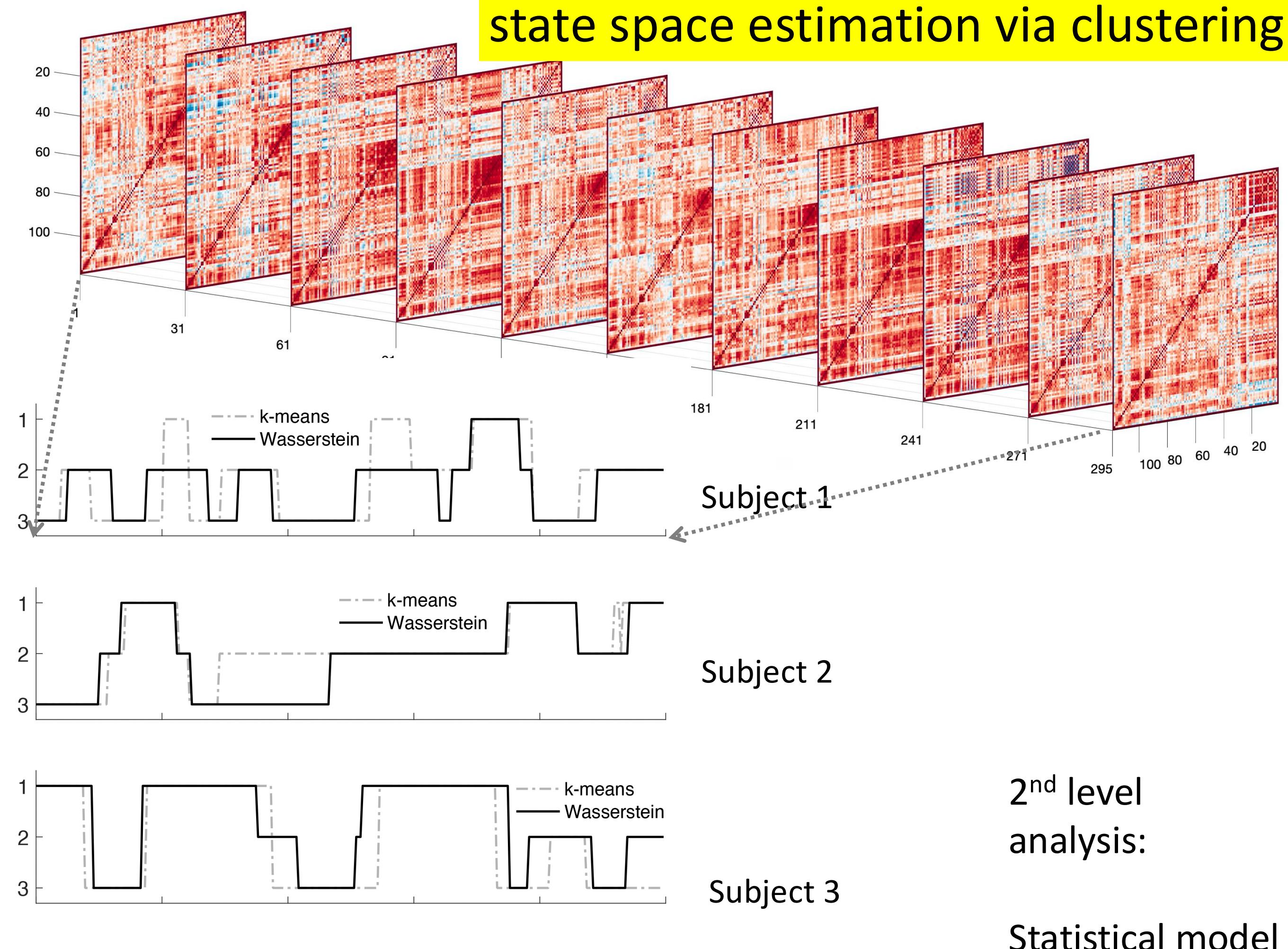
*Important  
biological  
questions are  
added*



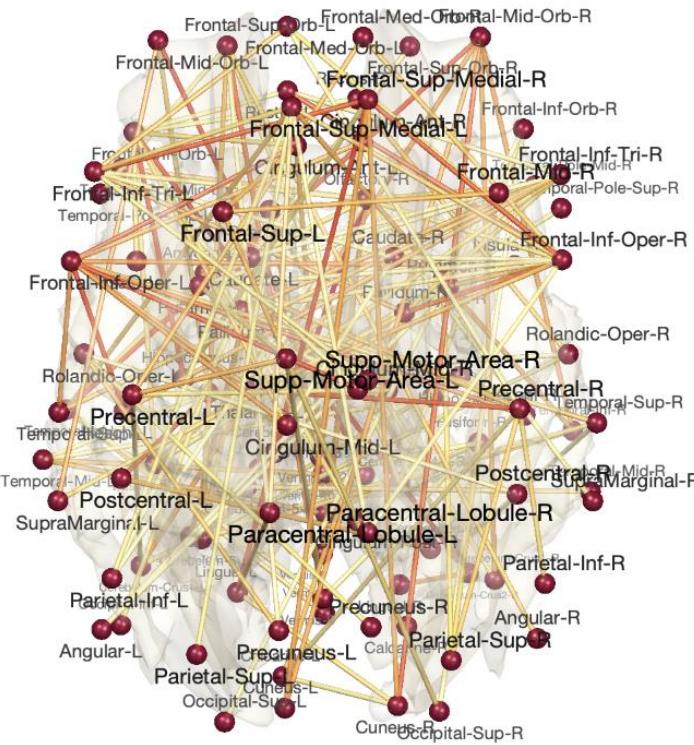
*Huang et al. 2020 Neuroscience Methods 331:108480*



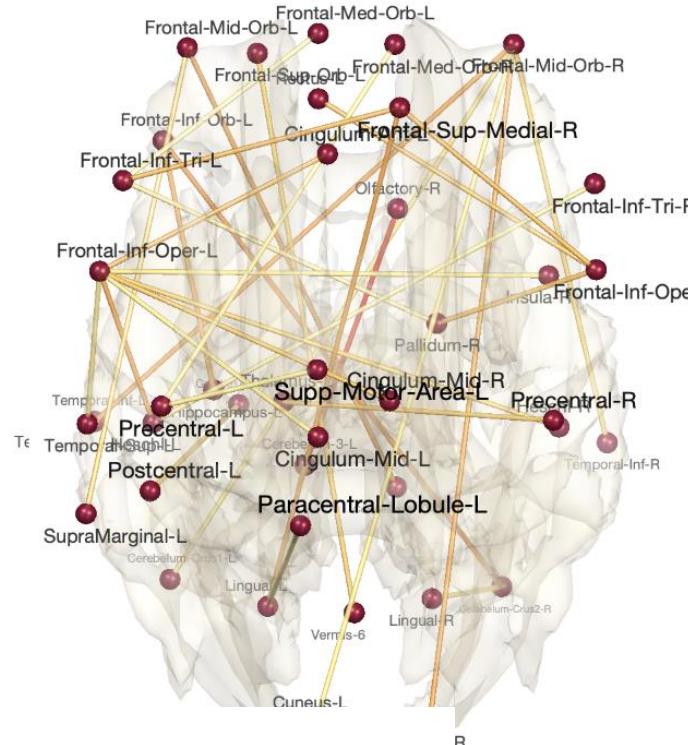
# state space estimation via clustering



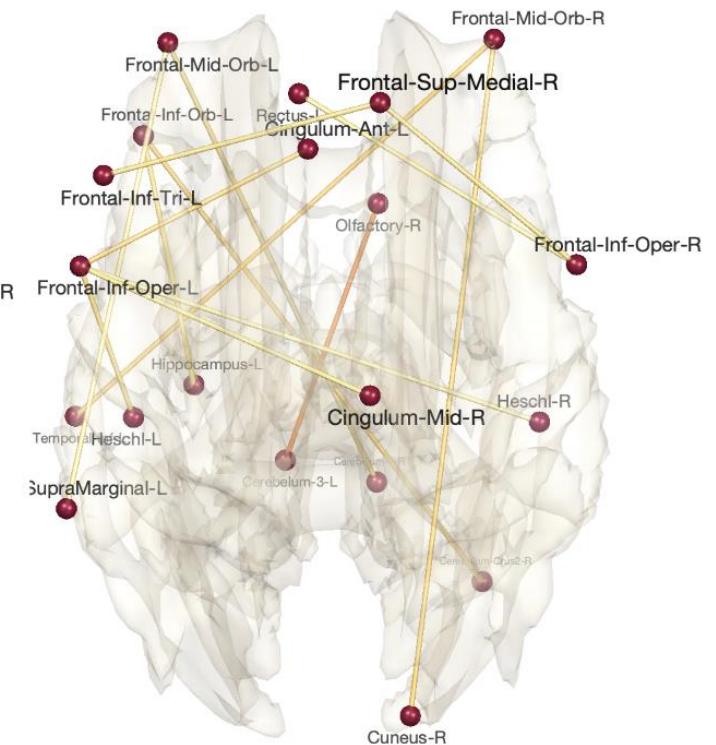
State 1



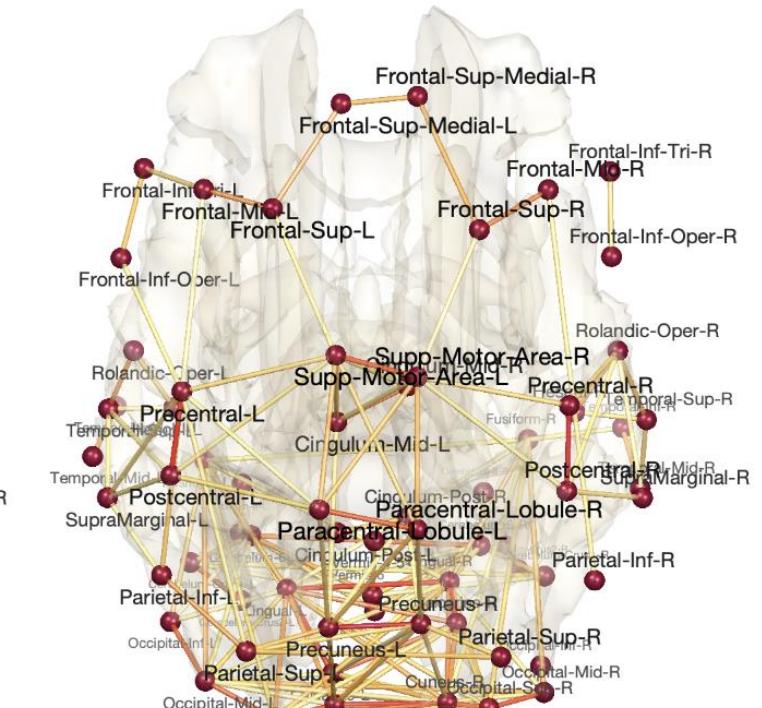
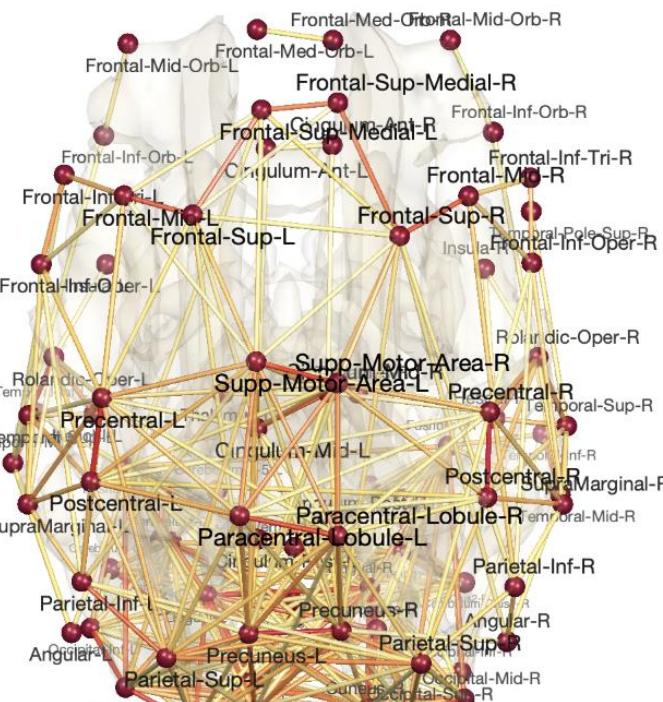
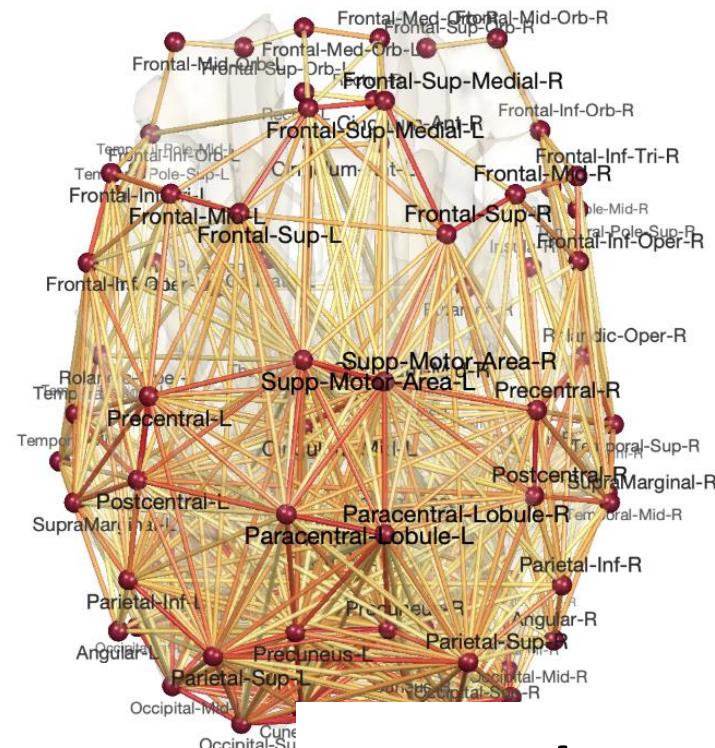
State 2



State 3

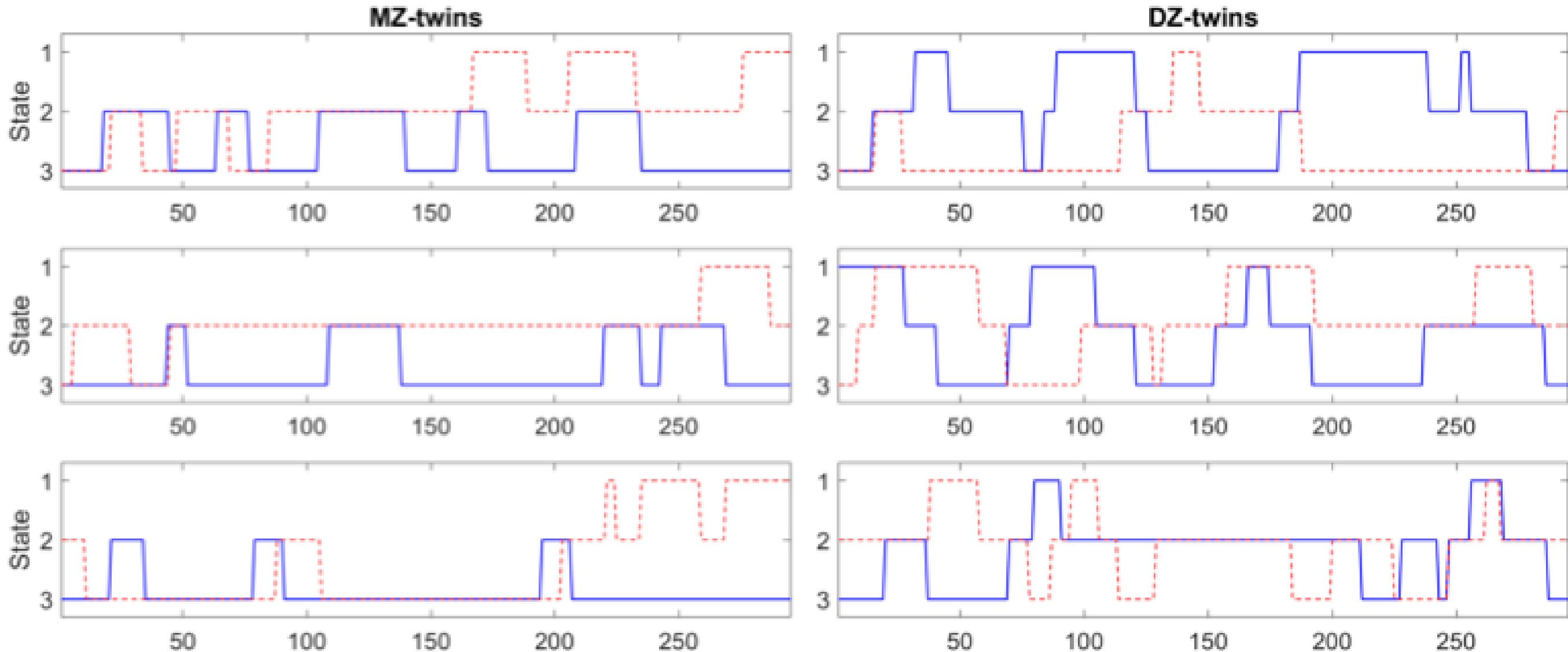


*k-means* *Sample mean in each state*



*Topological clustering* *Topological mean in each state*

# Is the state-change heritable?

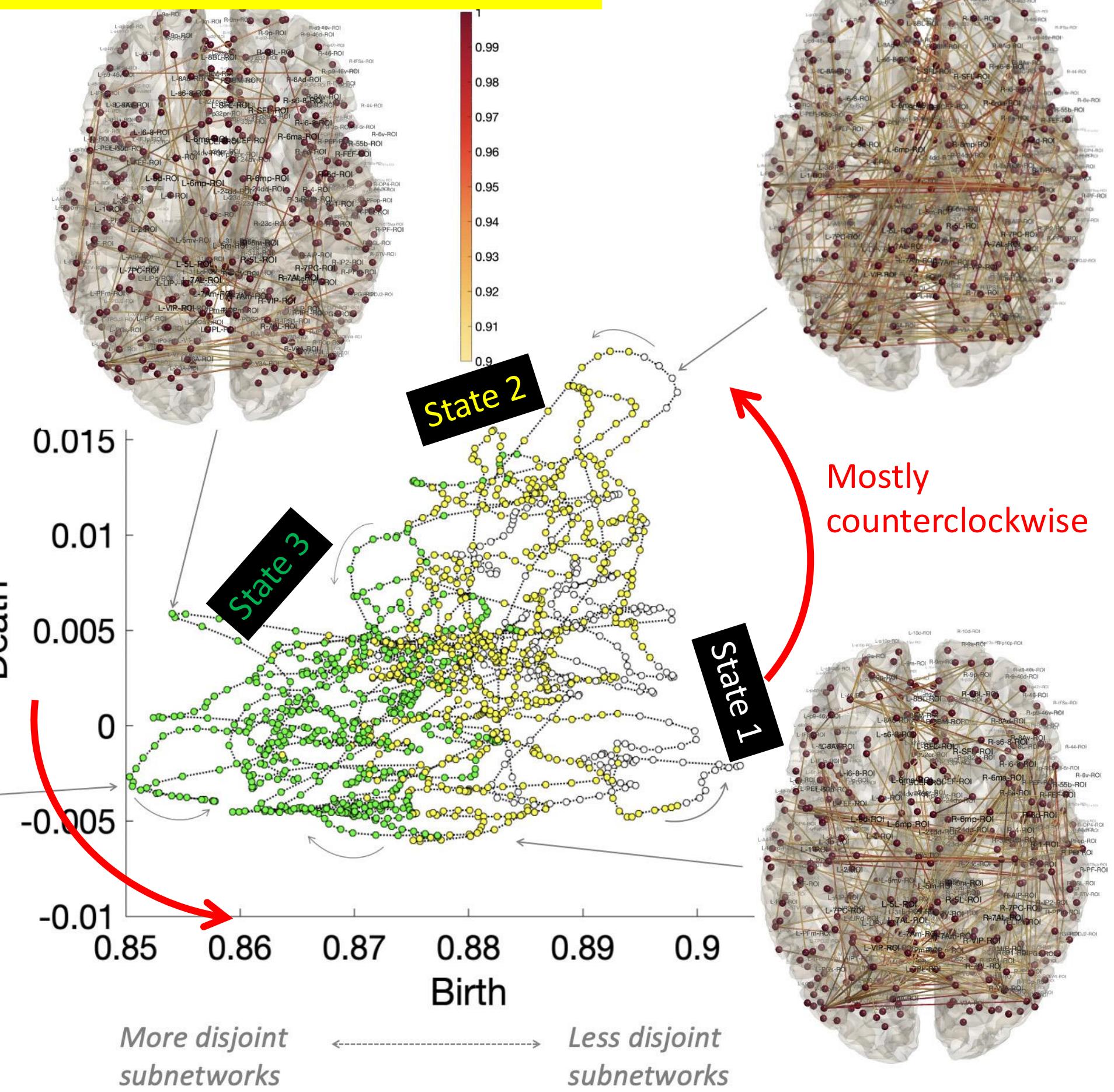


UW-Madison twin study (200 twin pairs)

# Interpreting clustering results is often needed

More cycles

less cycles



# Distances

# Wasserstein distance between networks

$$C_1 \cup C_2 \cup \dots \cup C_k = \{\mathcal{X}_1, \dots, \mathcal{X}_n\}, \quad C_i \cap C_j = \emptyset$$

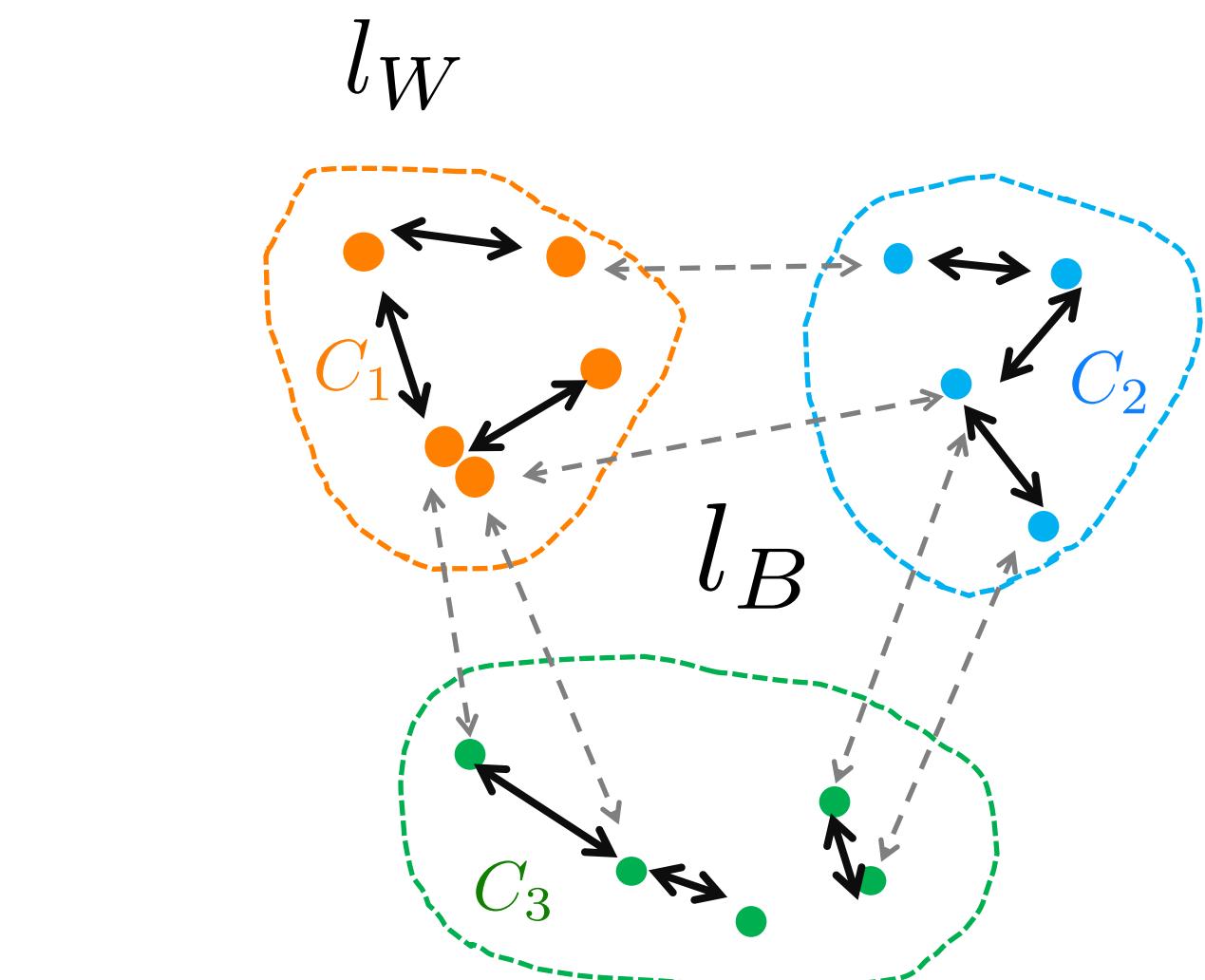
Between-group distance

$$l_B \propto \sum_{i \in C_1, j \in C_2} \mathcal{L}(\mathcal{X}_i, \mathcal{X}_j) \quad \leftarrow \text{----- 0D and 1D combined distances}$$

Within-group distance

$$l_W \propto \sum_k \sum_{i, j \in C_k} \mathcal{L}(\mathcal{X}_i, \mathcal{X}_j)$$

$$l_B + l_W = \sum_{i, j} \mathcal{L}(\mathcal{X}_i, \mathcal{X}_j)$$



# Clustering

Minimize the within cluster distance over all possible clustering configurations

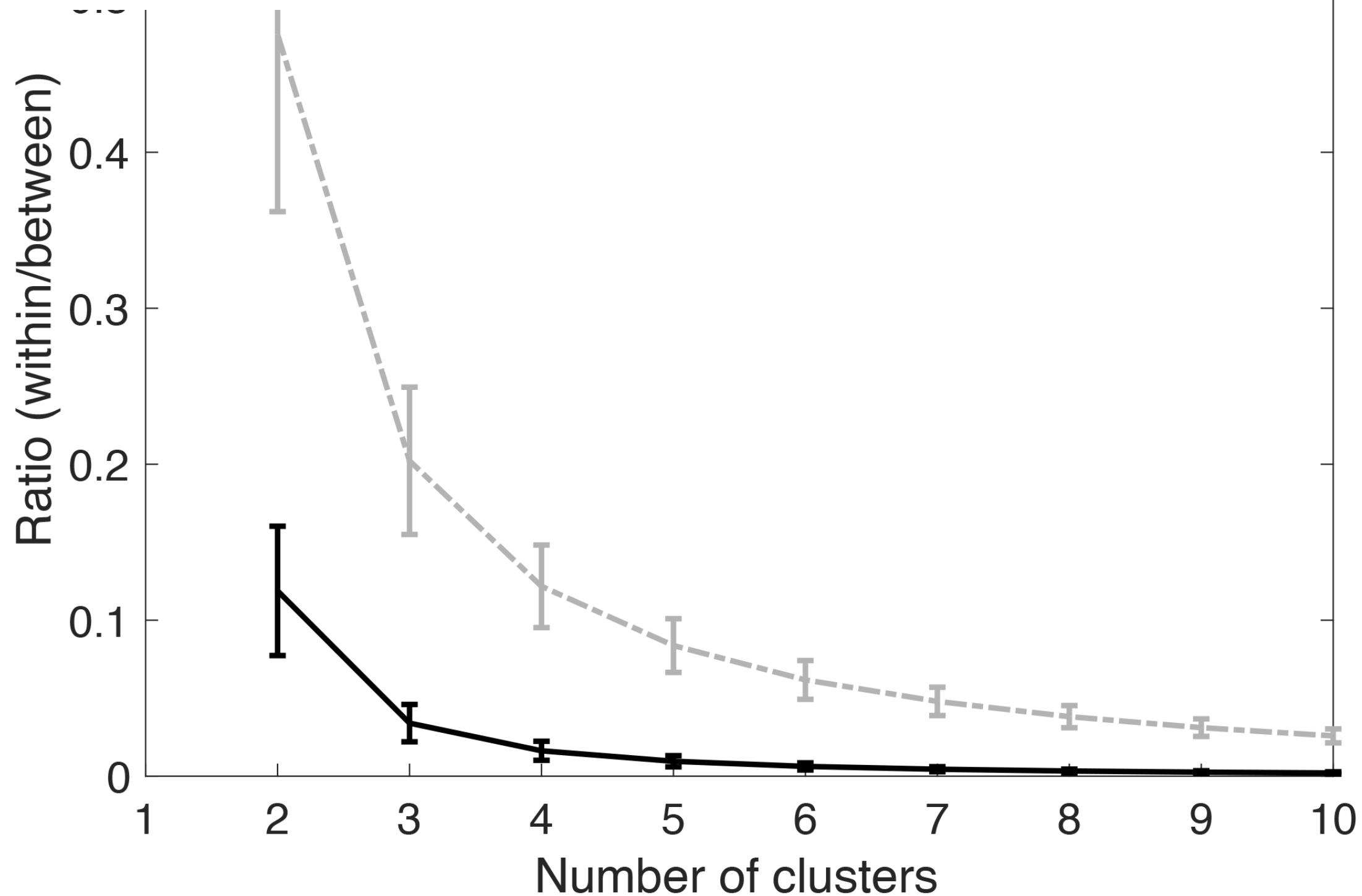
$$l_W \propto \sum_k \sum_{i,j \in C_k} \mathcal{L}(\mathcal{X}_i, \mathcal{X}_j)$$

Equivalently, maximize the between cluster distance

$$l_B = \sum_{i \in C_1, j \in C_2} \mathcal{L}(\mathcal{X}_i, \mathcal{X}_j)$$

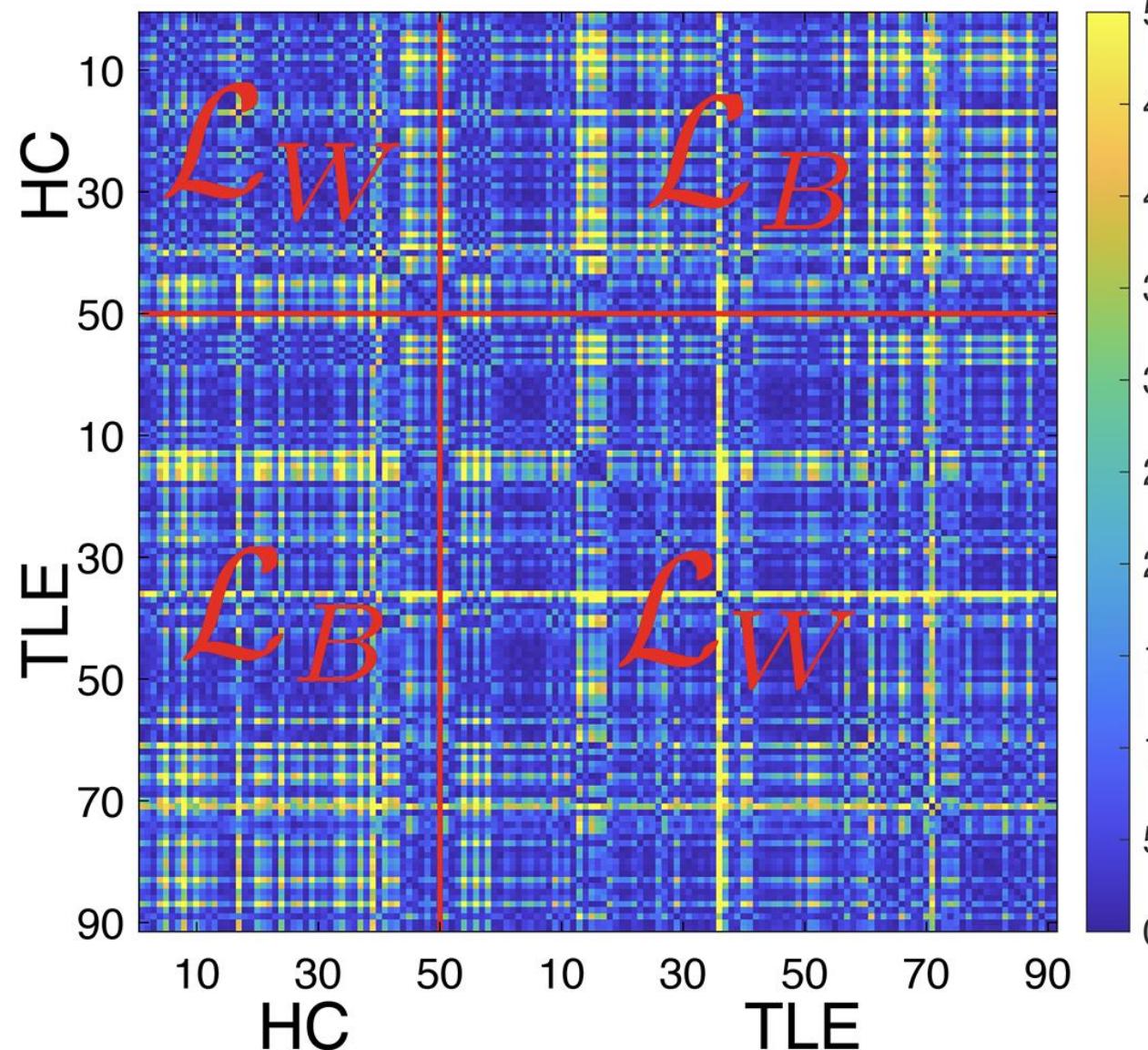
# Elbow method: Optimal number of clusters

$$\frac{1}{\phi} = \frac{l_W}{l_B}$$



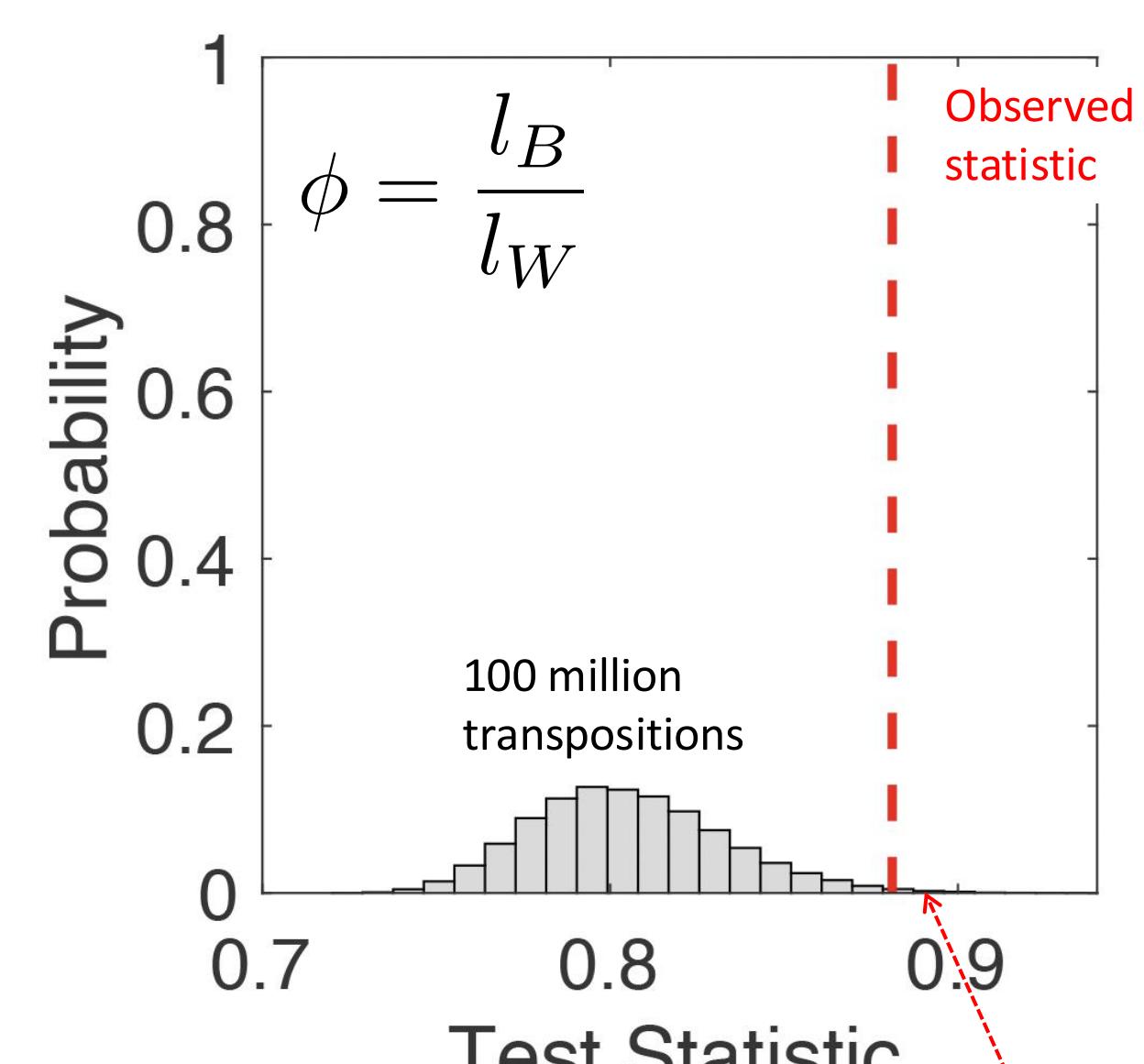
The within cluster variance *6 times* smaller

# Statistical inference on the equivalence of clusters



$$l_W \rightarrow l_W + \Delta(\text{tranposition})$$

$$l_B \rightarrow l_B + \Delta(\text{tranposition})$$



Songdechakraiwut and Chung  
2023

P-value 0.0086

# Mathematical equivalence of topological clustering and topological inference

There exists a monotonically decreasing function  $f$  satisfying

$$p\text{-value} = f(\text{clustering accuracy})$$



*Proof in Chung et al. 2023 NeuroImage*

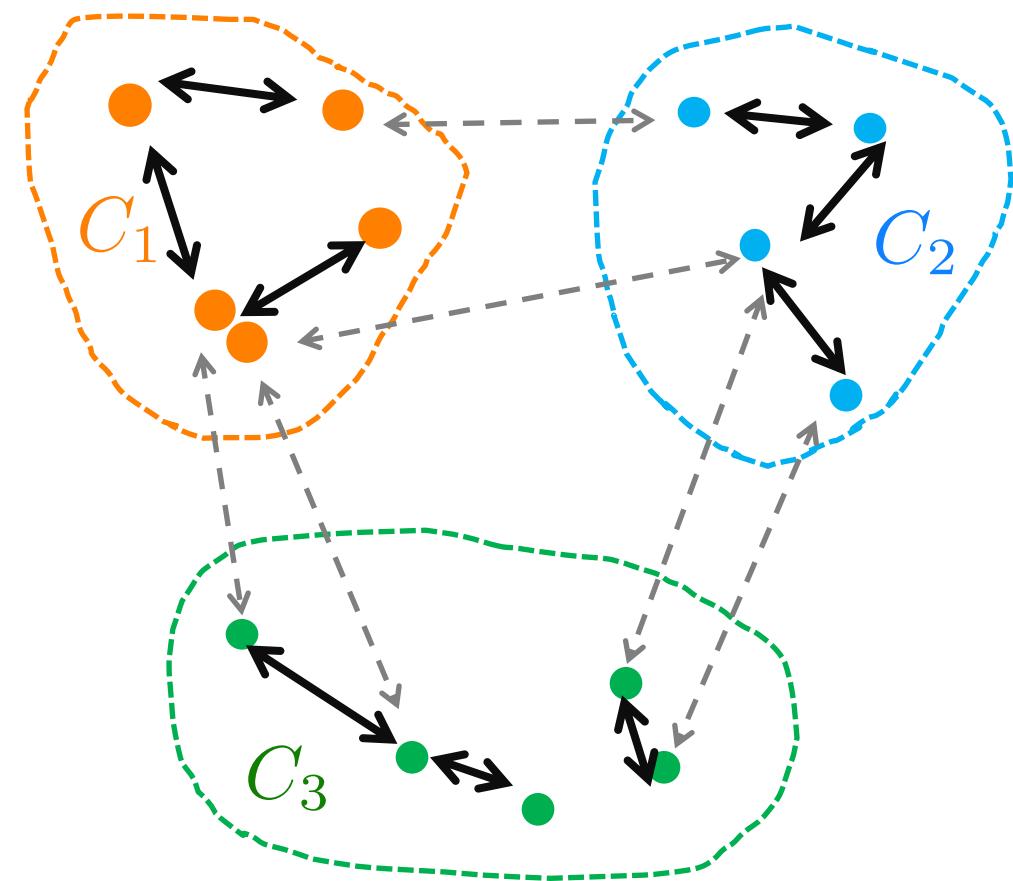
# Topological clustering = topological inference

$$\text{Within-group distance } d_W \propto \sum_k \sum_{i,j \in C_k} d_0(\mathcal{X}_i, \mathcal{X}_j) + d_1(\mathcal{X}_i, \mathcal{X}_j)$$

Clustering accuracy

$$\frac{dA}{dp} = \frac{dA}{d(d_W)} \frac{d(d_W)}{dp} \leq 0$$

*p-value*



There exists a **monotone decreasing** function  $f$  satisfying  $p = f(A)$

# Geometric clustering methods fail topological task

