

Scalable Functional Connectivity and Functional Clustering for fMRI

Orhan Firat, Ph.D. Student
Department of Computer Engineering,
Middle East Technical University
orhan.firat@ceng.metu.edu.tr

Affiliation and Collaboration



Department of Computer Engineering
Neuro-ceng



Department of Psychology



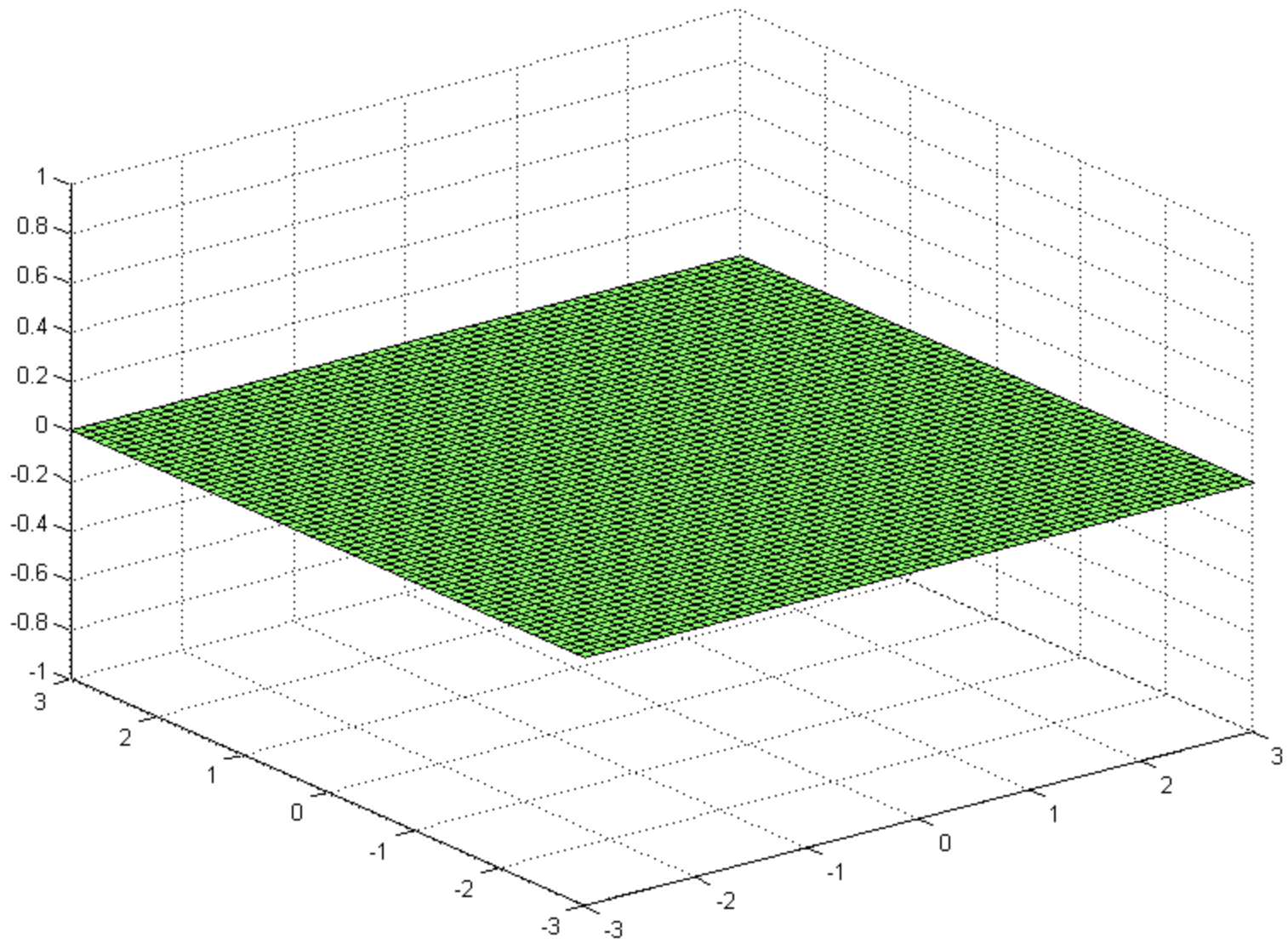
Graduate School of Informatics

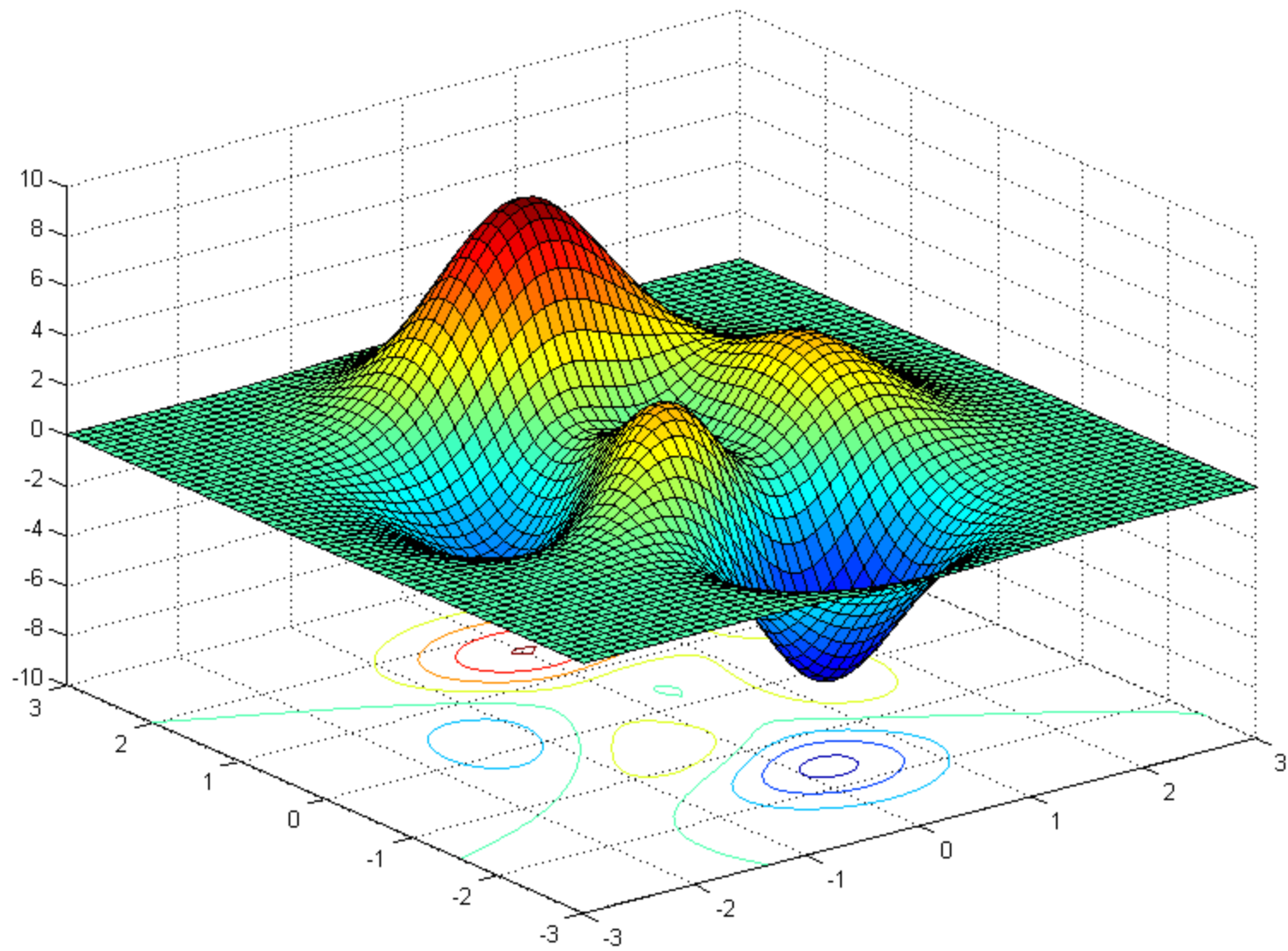


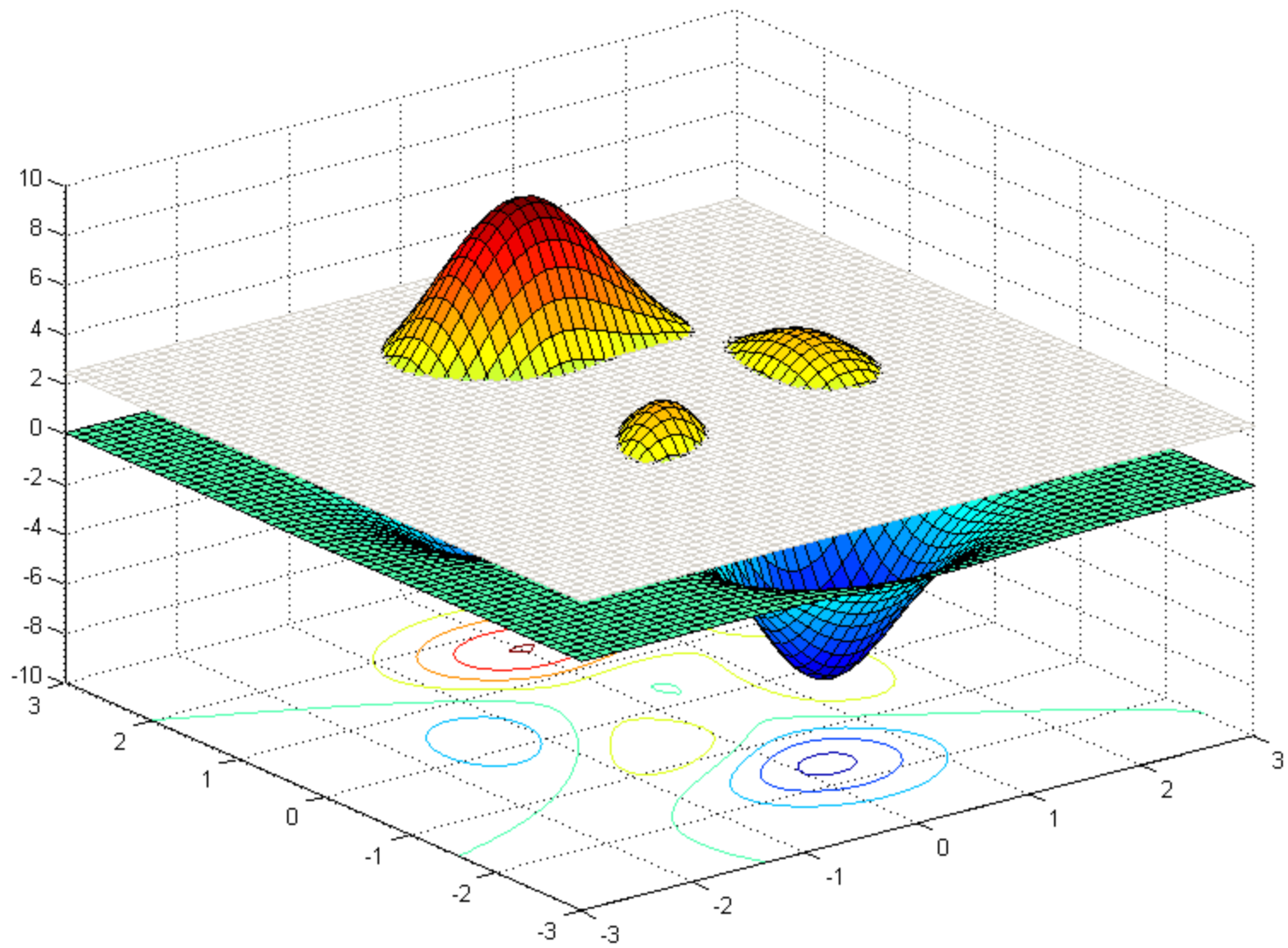
What do we do?

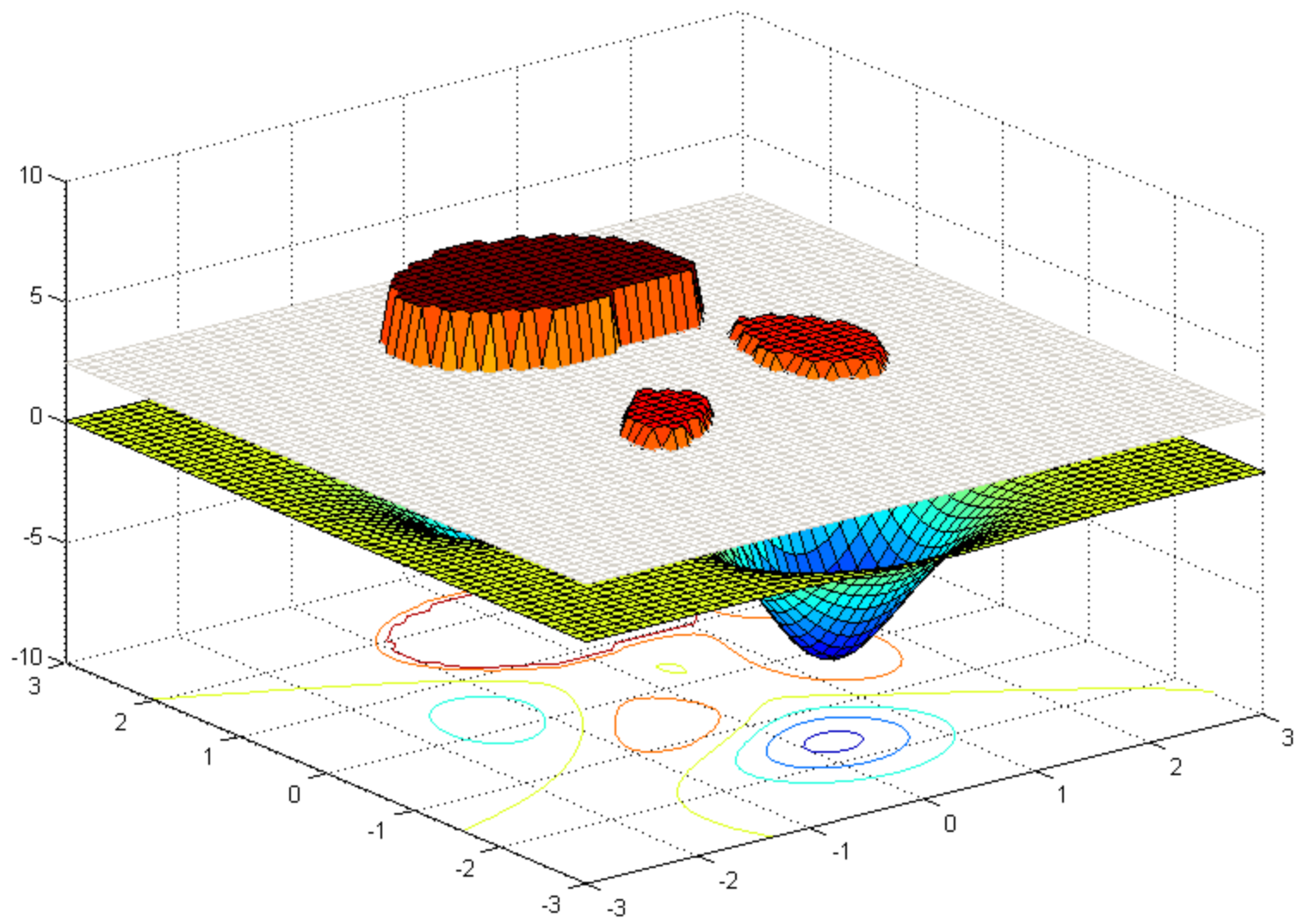
Computational Neuroscience

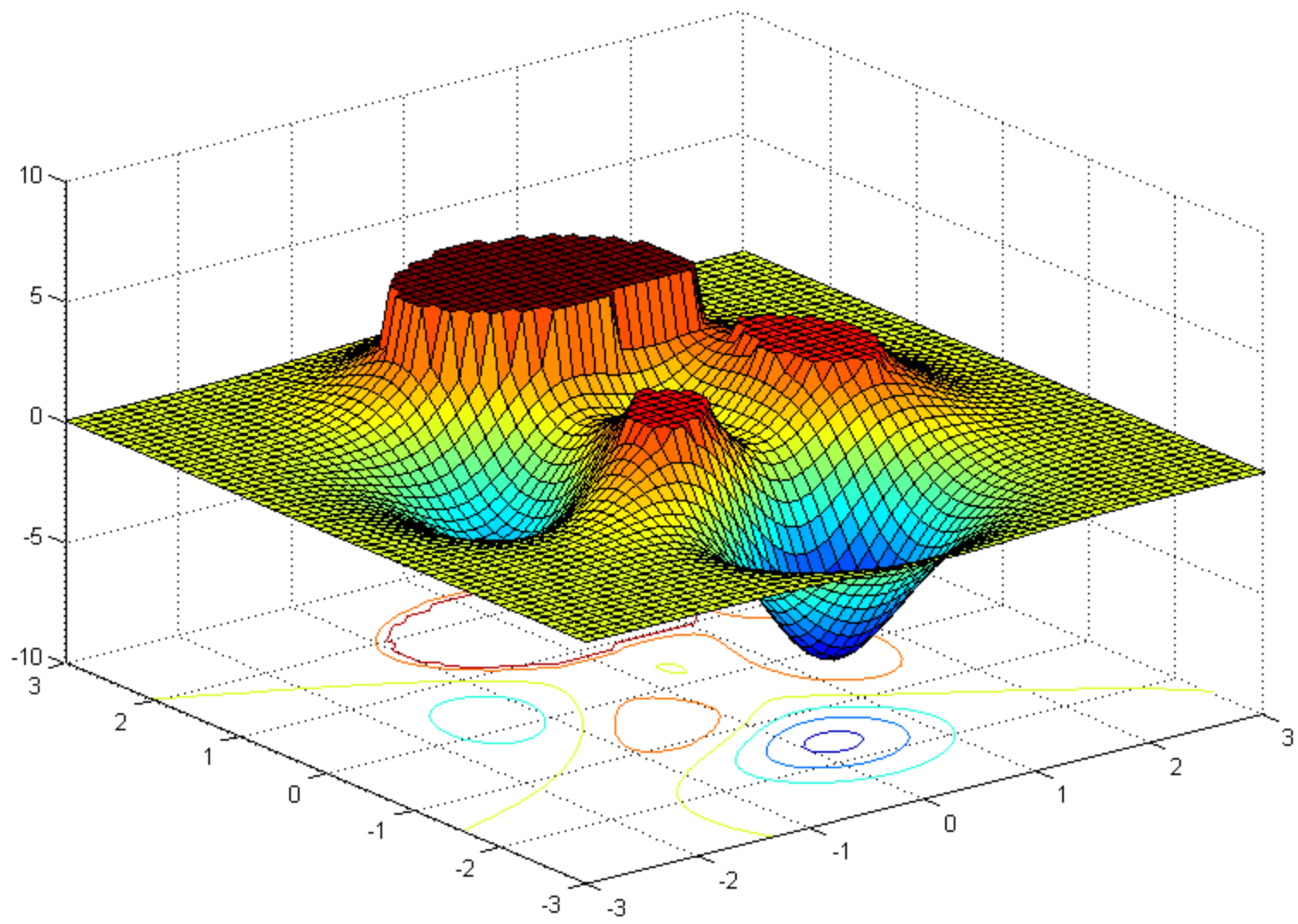
- Domain : fMRI
- Objective : Develop representative and informative models for cognitive processes
 - Pattern Analysis (MVPA), Machine Learning
 - Structure Learning
 - Structured Prediction
 - Mind Reading 😊

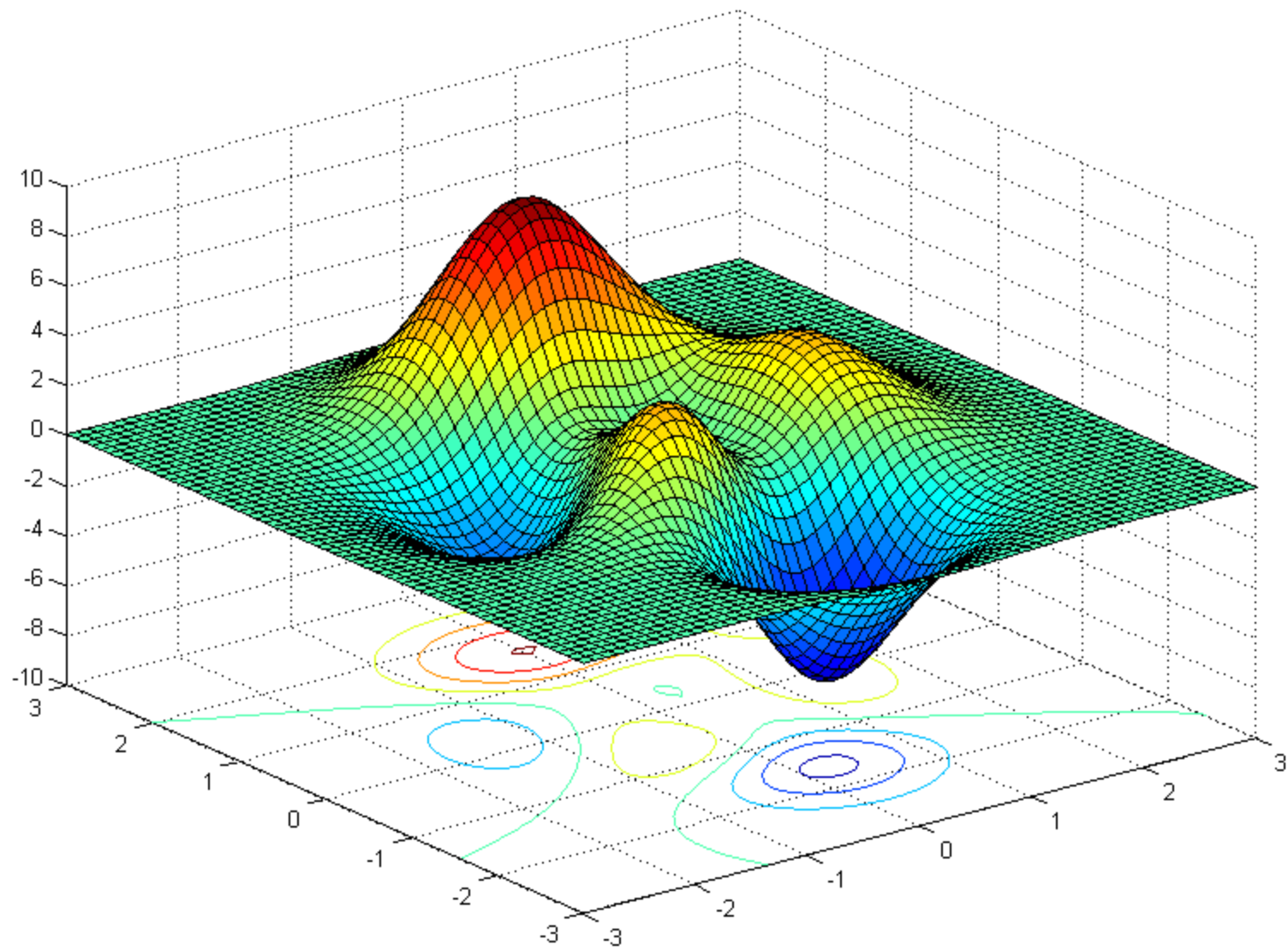












Pattern Analysis and Machine Learning

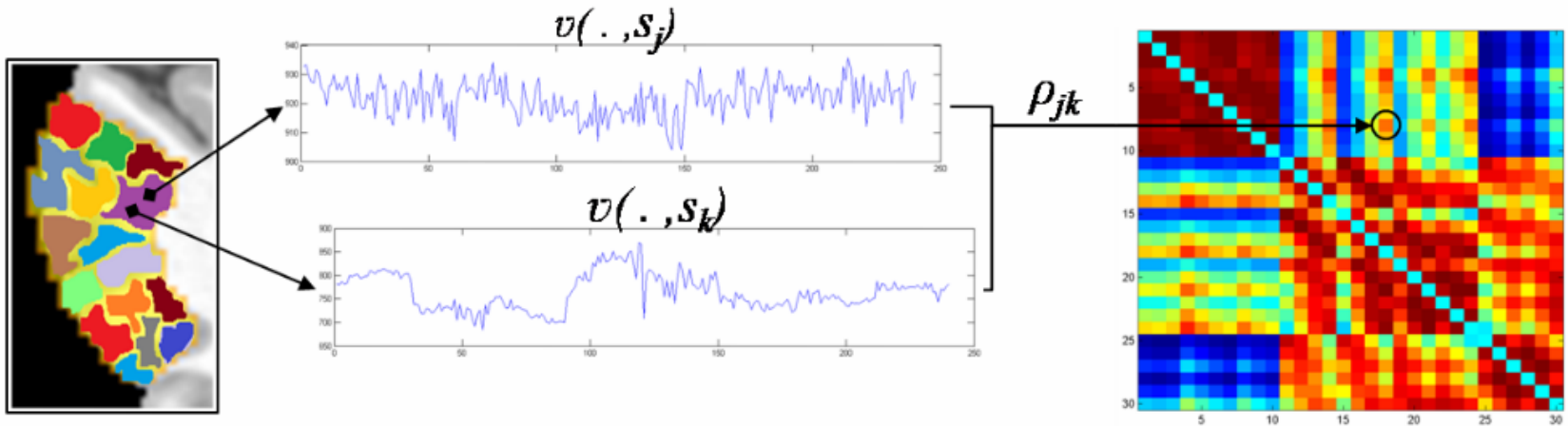
- More data → robust approaches for representation / discrimination
- Science
 - New learning paradigms
 - New models for cognitive process representation
- Engineering
 - Computational complexity; time-memory complexities.
 - Practical solutions for researchers

How do we do? (Engineering Part)

Case Studies:

1. Applying large scale functional connectivity on GPU
2. Clustering / parcellation functional connectivity matrices on hybrid architectures

Functional Connectivity



Functional Connectivity (FC) is the statistical dependence between remote neural elements or regions across time.

1. Find correlation between two voxels' time series using a correlation metric.
2. Construct correlation matrix (connectivity matrix) by using each pairs' correlation response.

Scalability of Functional Connections

- Connectivity matrices are expensive in voxel level.
- Considering functional relations of a voxel with all other voxels;
- 8142 voxels makes ~33M functional relations
- 82926 voxels makes ~3438M functional relations

Scalability of Functional Connections

- Suppose region of interest consists of M voxels
 - Let computation time of functional relations for a voxel pair $(v_i - v_j)$ takes t second
 - Computation time for functional connectivity matrix T_{fc} takes
$$T_{fc} = M^2 \times t \text{ seconds}$$
 - T_{fc} grows exponentially on M .

Ex: Suppose computation time of each functional connectivity metric takes ~ 1 sec

For 8142 voxels,

$$T_{fc} = (8142 \times 8142) / 2 - 8142 = 33.137.940 \text{ sec} \cong 9204 \text{ hours} \cong 383 \text{ days}$$

with 256-way perfect parallelization (in a cluster) $\cong 1,5$ days

Scalability Analysis

- Calculating even $\sim 8000 \times 8000$ sized correlation matrix is expensive where our domain spans $\sim 80,000$ voxels and their relations.
- Apply GPU parallelization for speed-up.
- Reformulate correlation metric for GPU.

Pearson Correlation Coefficient

- Cross-correlation of any two individual time-series x and y is given by

$$\rho_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}}$$

- This can also be written as:

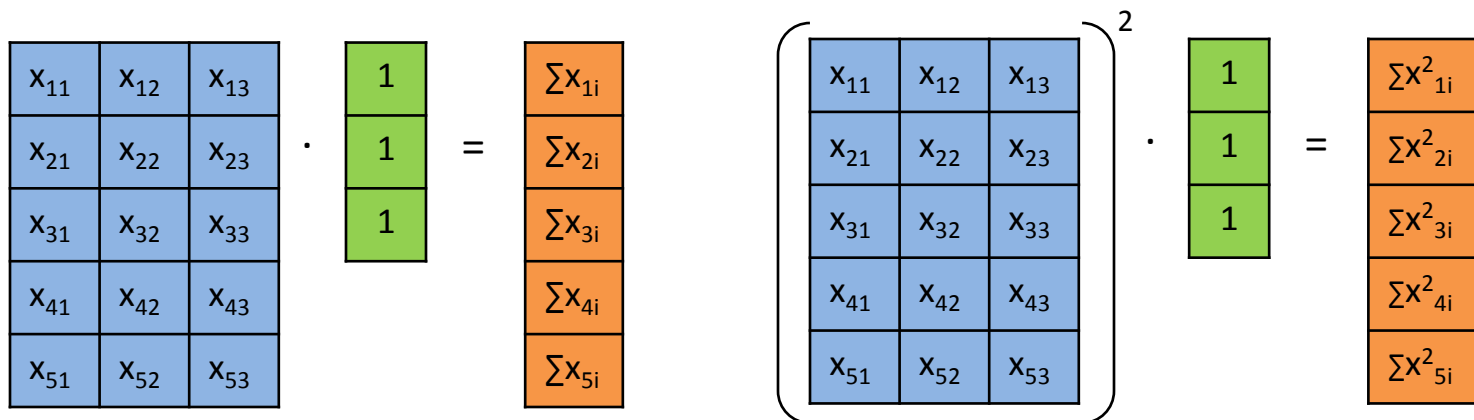
$$\rho_{xy} = \frac{\sum x_i y_i - n\bar{x}\bar{y}}{\sqrt{(\sum x_i^2 - n\bar{x}^2)(\sum y_i^2 - n\bar{y}^2)}}$$

- Design appropriate GPU architecture for reformulated correlation

$$\rho_{xy} = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sqrt{(\sum x_i^2 - n \bar{x}^2)(\sum y_i^2 - n \bar{y}^2)}}$$

$$\rho_{xy} = \frac{\sum x_i y_i - n(\sum x_i/n)(\sum y_i/n)}{\sqrt{(\sum x_i^2 - n(\sum x_i/n)^2)(\sum y_i^2 - n(\sum y_i/n)^2)}}$$

- Highlighted summations can be calculated with a single matrix-vector multiplication



$$\rho_{xy} = \frac{\sum x_i y_i - n(\sum x_i/n)(\sum y_i/n)}{\sqrt{(\sum x_i^2 - n(\sum x_i/n)^2)(\sum y_i^2 - n(\sum y_i/n)^2)}}$$

- Only highlighted summation needs to be calculated in a kernel as the last step.
- Other summation results are stored in texture memory.

How to Handle Massive Data

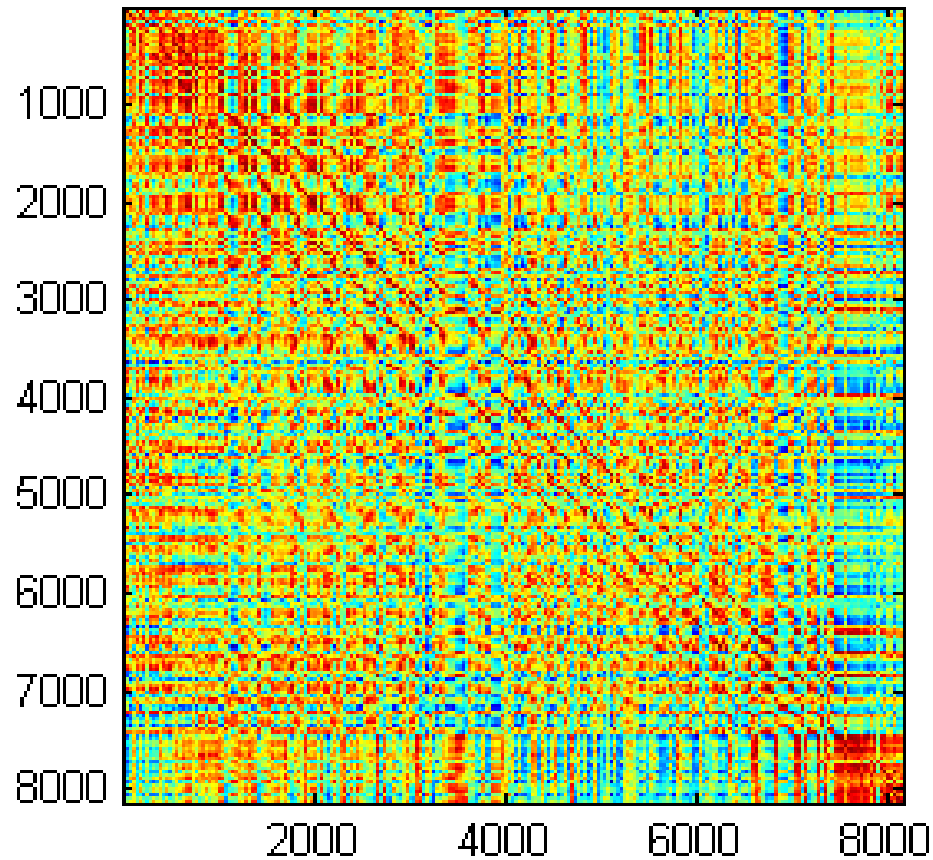
- Overlay CUDA-Grids on resulting symmetric correlation matrix (lower-diagonal).
- Determine chunk size and process each chunk sequentially.
- Transfer corresponding time-series data to texture memory.
- Process current chunk then load subsequent chunk (in column major order)

ρ_{11}	ρ_{12}				
ρ_{21}	ρ_{22}				
ρ_{31}	ρ_{32}	ρ_{33}	ρ_{34}		
ρ_{41}	ρ_{42}	ρ_{43}	ρ_{44}		
ρ_{51}	ρ_{52}	ρ_{53}	ρ_{54}	ρ_{55}	ρ_{56}
ρ_{61}	ρ_{62}	ρ_{63}	ρ_{64}	ρ_{65}	ρ_{66}

Overall Algorithm

1. Load data chunk to device memory
2. Calculate $\sum x_i$ with *cublasSgemv*,
 - save resulting vector on device memory
3. Square data matrix by using thrust library transform routine and calculate $\sum x_i^2$ with *cublasSgemv*.
 - save resulting vector on device memory
4. Load source and destination time series into texture memory.
5. Calculate pairwise correlation on GPU
 - save resulting matrix chunk on disk
6. Return to step 4

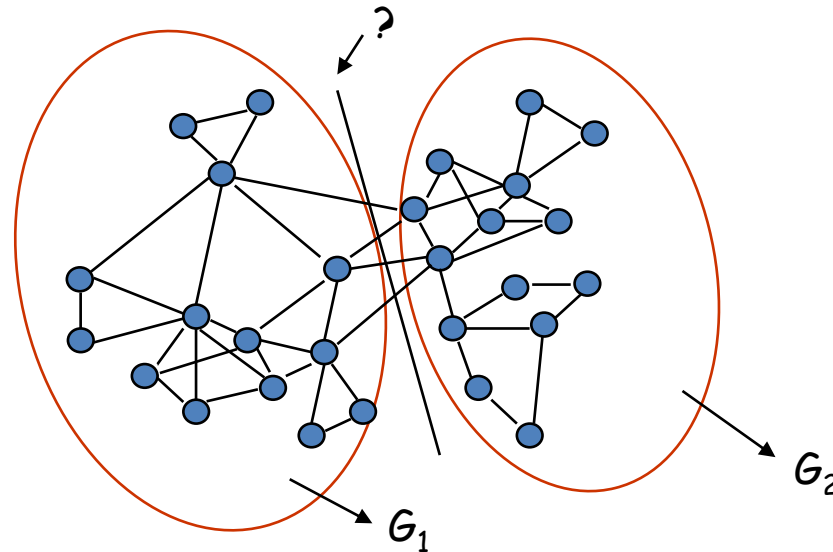
Functional Connectivity Matrix for 8000 voxels



Method	Conventional	Proposed
8142 voxels	~30 minutes	28 seconds
82926 voxels	~3 days	24 minutes

Case Study 2: Functional Parcellation / Clustering

- Consider a weighted Graph $G=(V, E, A)$
 - it is possible to partition G

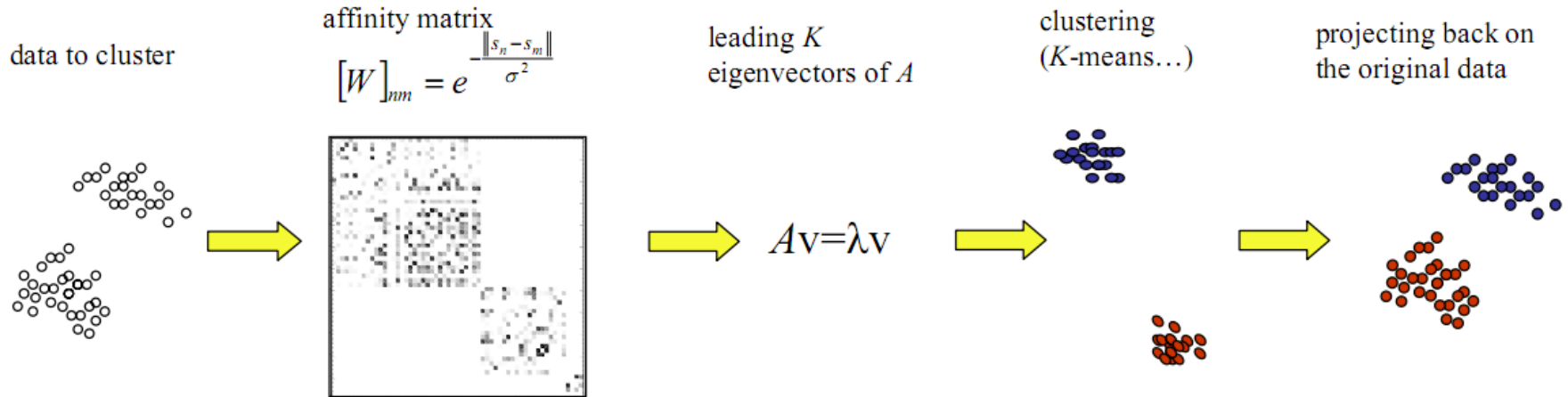


- into smaller components G_1 and G_2 where $G = G_1 \cup G_2$ with specific properties defined over A

How to Partition the Graph?

- Problem:
 - finding an optimal graph (normalized) cut is NP-hard
- Solution:
 - approximation & heuristics
- Approximation: spectral graph partitioning
 - Partitioning the graph by spectral analysis
 - (spectrum of a matrix is the set of its **eigenvalues**)
 - e.g. Spectral clustering
 - Based on Laplacian Matrix, or **Graph Laplacian**

Spectral Clustering

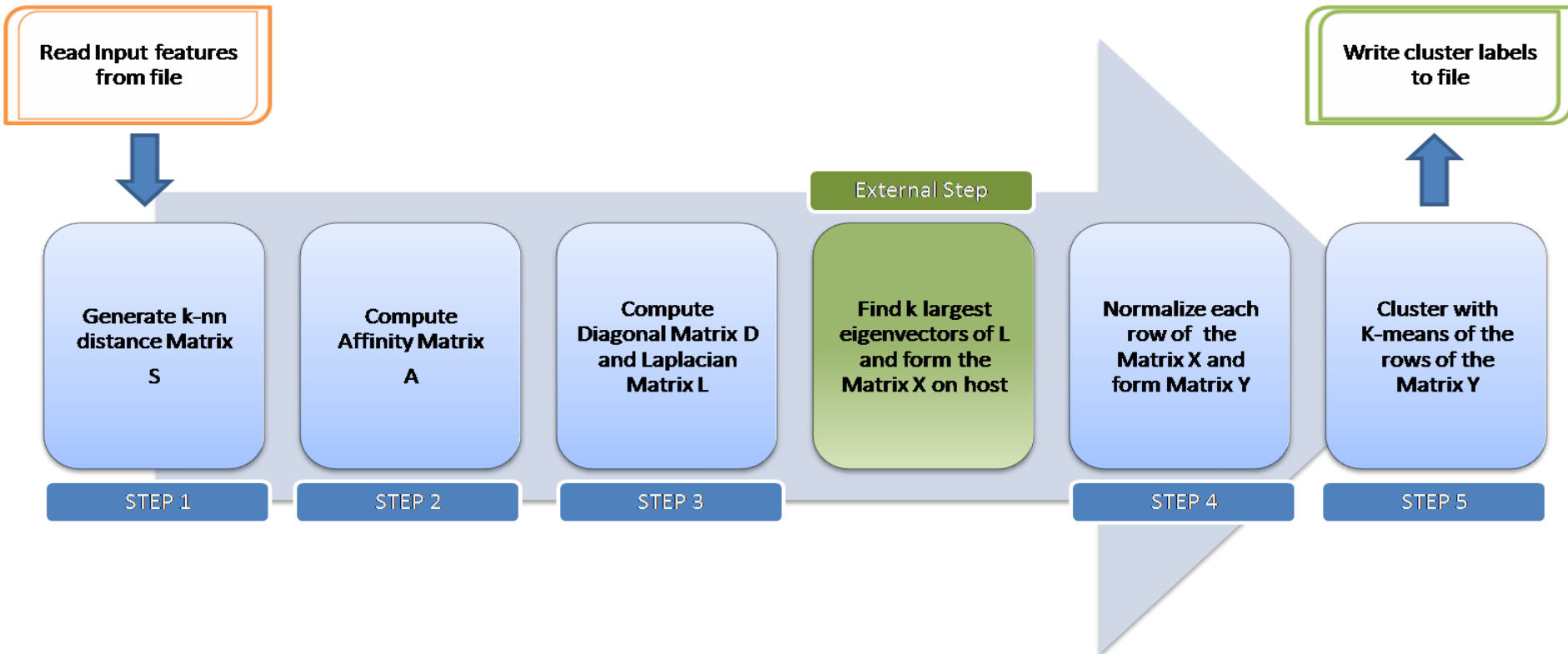


- SC is sensitive to the scaling parameter of the RBF kernel
- Main difference between algorithms is the definition of $A = \text{func}(W)$

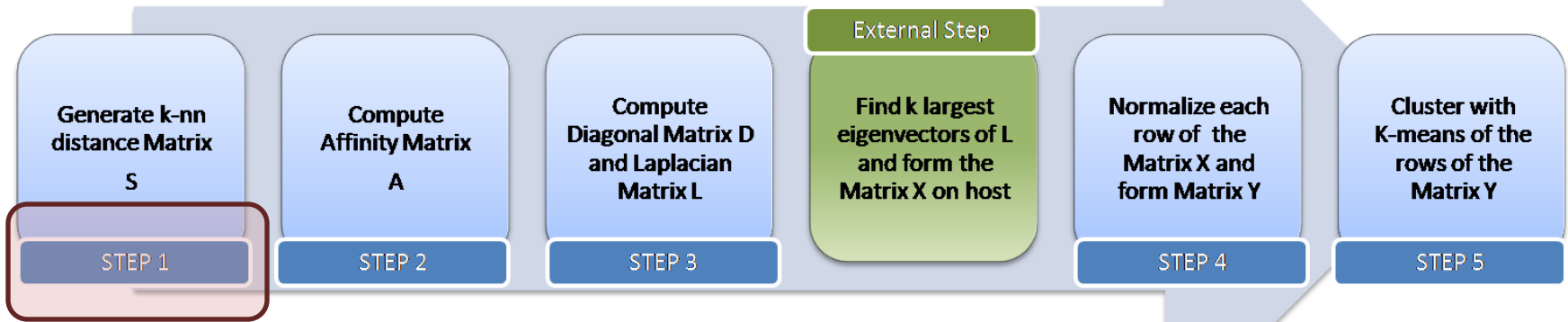
Major Operations and Parallelization of Spectral Clustering

- Sparse vector-vector multiplications
- Matrix-vector operations
- Matrix problems usually imply huge possibility to parallelization
 - Compute-intensive
 - Matrices are able to be divided by rows, columns or blocks

Algorithm Flow



- Step-wise analysis
 - Consists of comparisons with either a CPU implementation or GPU implementations



Distance matrix

Calculate distance from each point to others

$$d(\mathbf{q}, \mathbf{p}) = \sqrt{(q_1 - p_1)^2 + (q_2 - p_2)^2 + \dots + (q_n - p_n)^2}$$

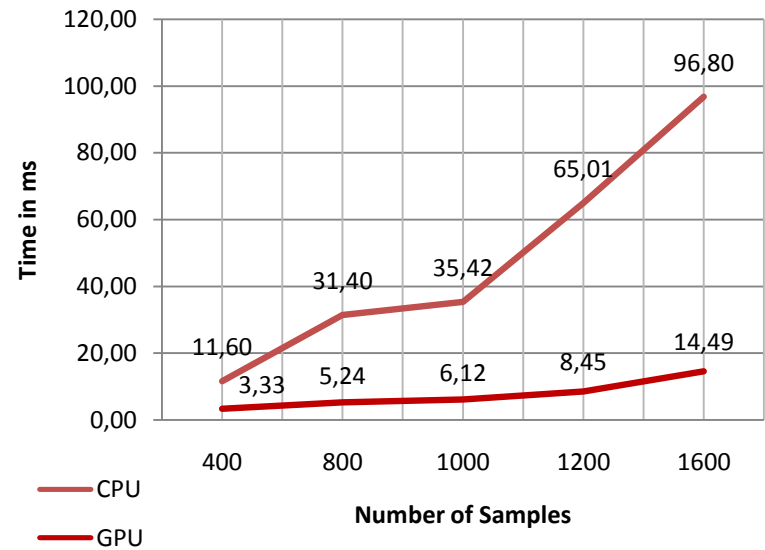
$$\|\mathbf{q} - \mathbf{p}\| = \sqrt{\|\mathbf{p}\|^2 + \|\mathbf{q}\|^2 - 2\mathbf{p} \cdot \mathbf{q}}$$

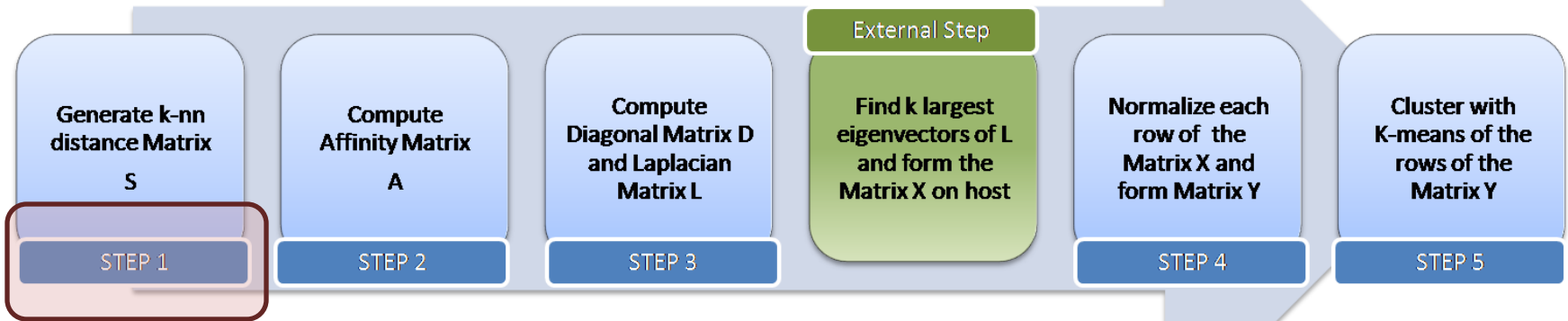
Multiply feature matrix with its transpose using CUBLAS

$$\begin{bmatrix} pp & pq & pq \\ qp & pp & pq \\ qp & qp & pp \end{bmatrix}$$

Launch 1 additional kernel to calculate distance matrix

CPU vs GPU



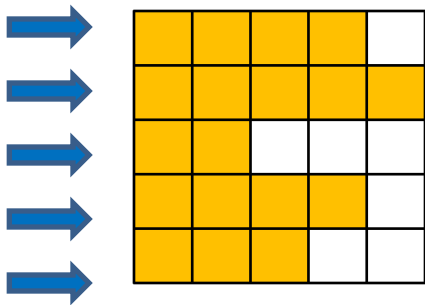


Find k-nearest neighbor

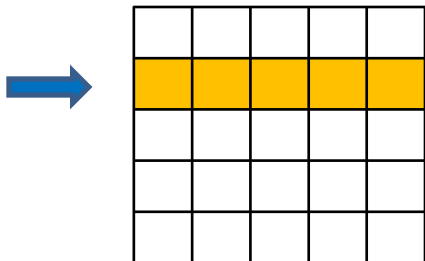
Extract closest points row-wise from its diagonal

Sort rows of the distance matrix

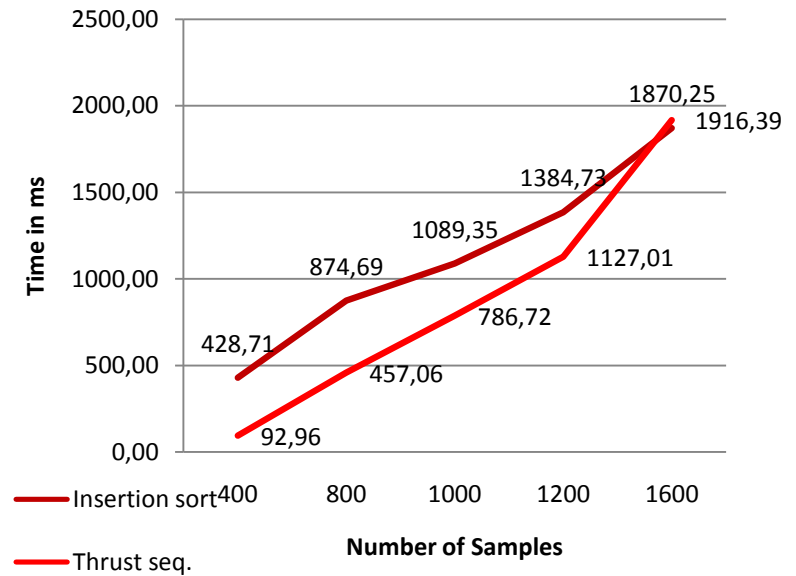
1. Insertion sort with kernel

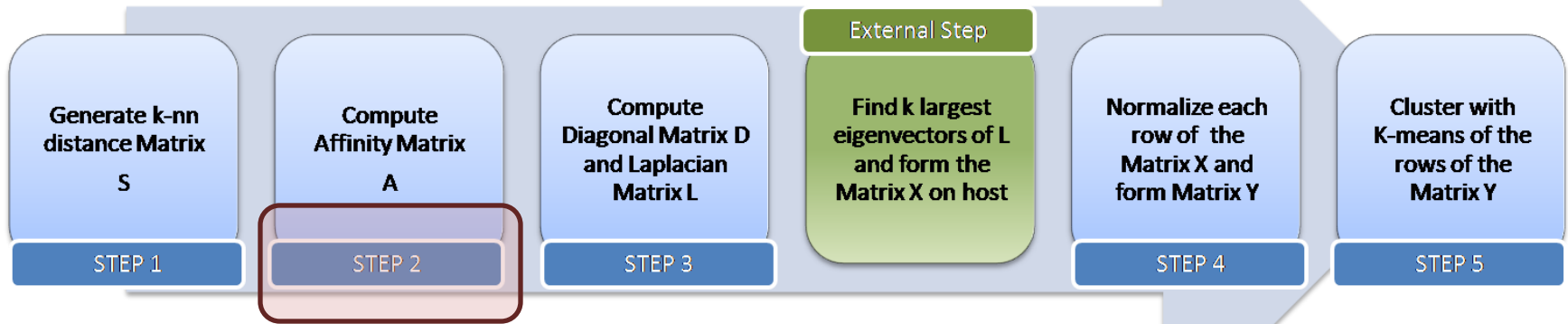


2. Sequential sort with thrust::sort



Insertion sort vs Thrust sort





Affinity (Similarity) Matrix

Convert distance matrix to a sparse similarity matrix

Use radial basis kernel in order to calculate similarities between data points

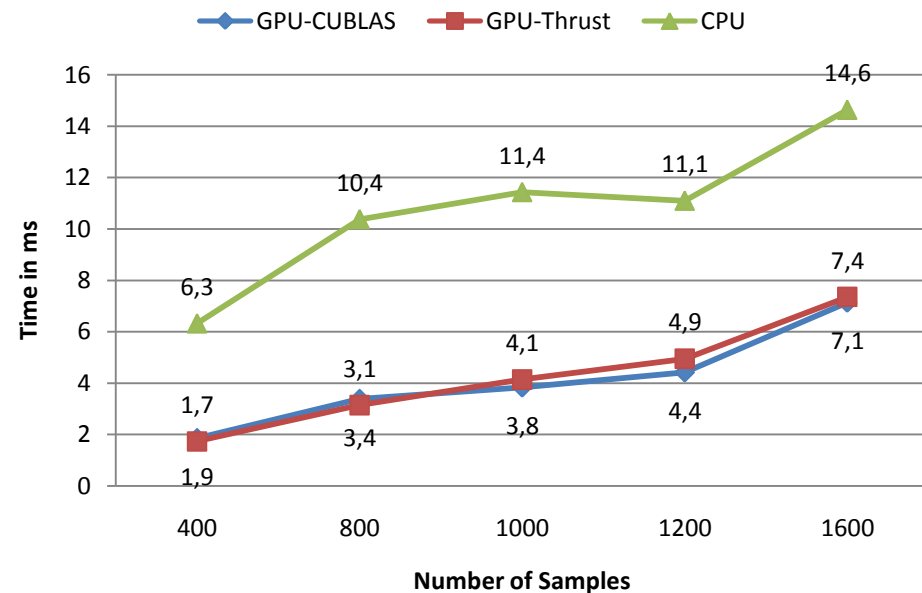
$$A_{i,j} = e^{-\|s_i - s_j\| / 2\sigma_i\sigma_j} \quad i \neq j$$

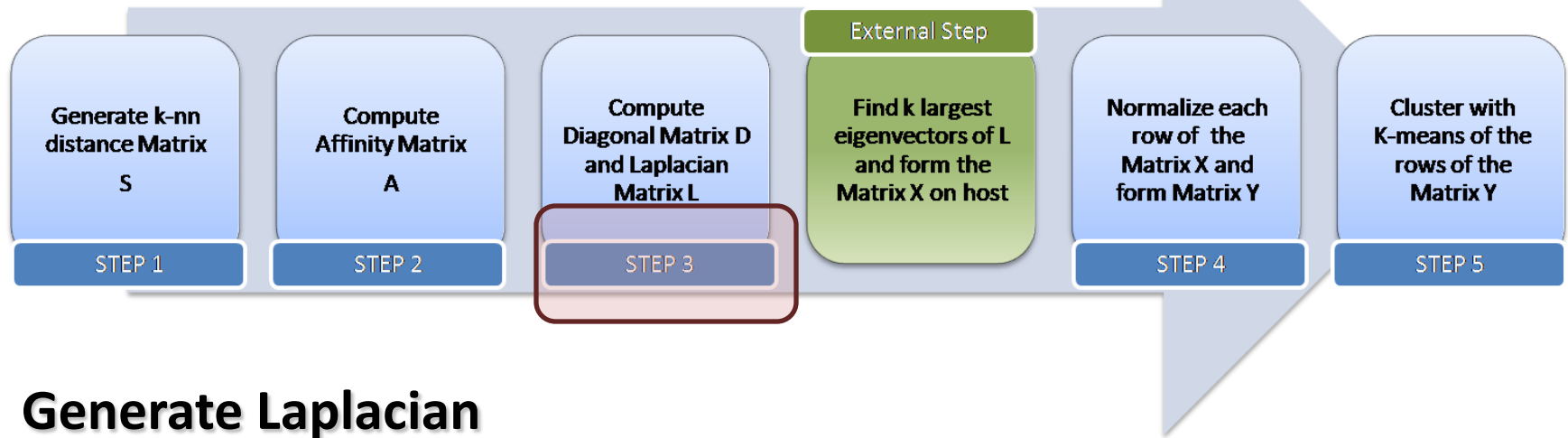
$$A_{i,j} = 0 \quad i = j$$

local scaling parameter σ_i calculated as the mean of row i in distance matrix S

Sum of each element in a row on GPU

1. sequential reduction using Thrust
2. Matrix-vector multiplication CUBLAS





Generate Laplacian

Compute diagonal matrix D as the **row sum** of Affinity Matrix A

Using CUBLAS *cublas_dgemv*

Multiply with Similarity matrix on left&right

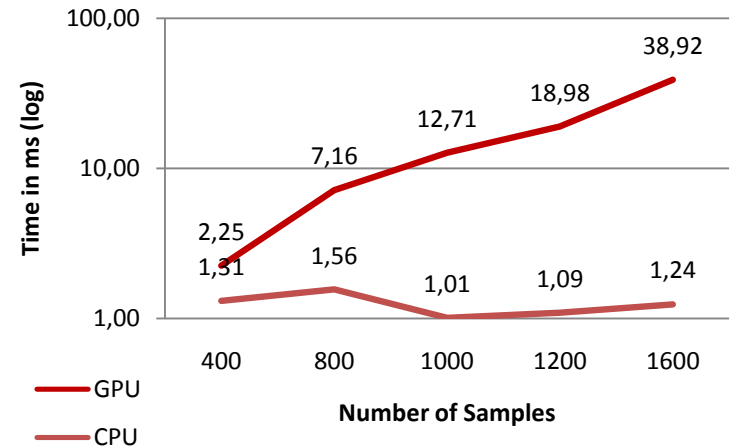
Using CUBLAS *cublas_dgemm*

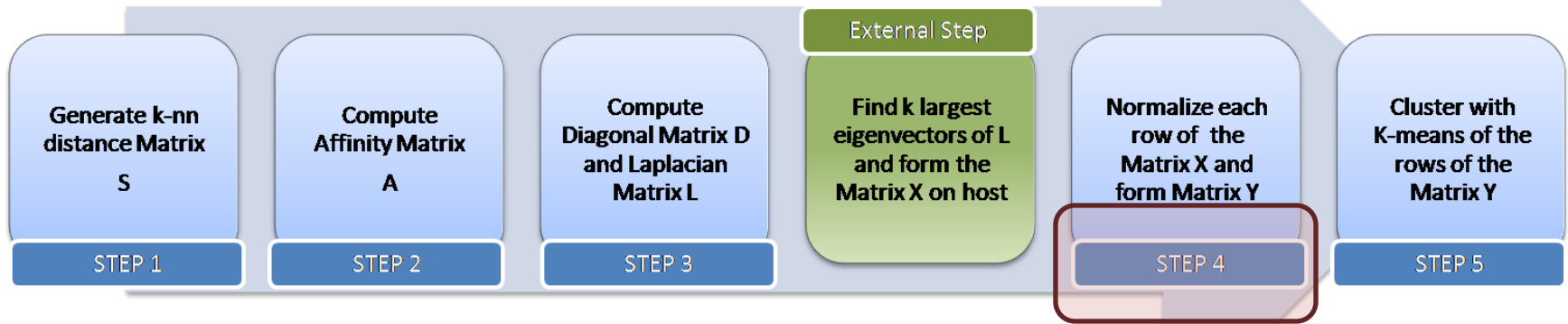
Launch 1 additional kernel to symmetrize

*CPU implementation outperforms GPU

Data transfer cost supresses the computation cost

Spectral Computations CPU vs GPU





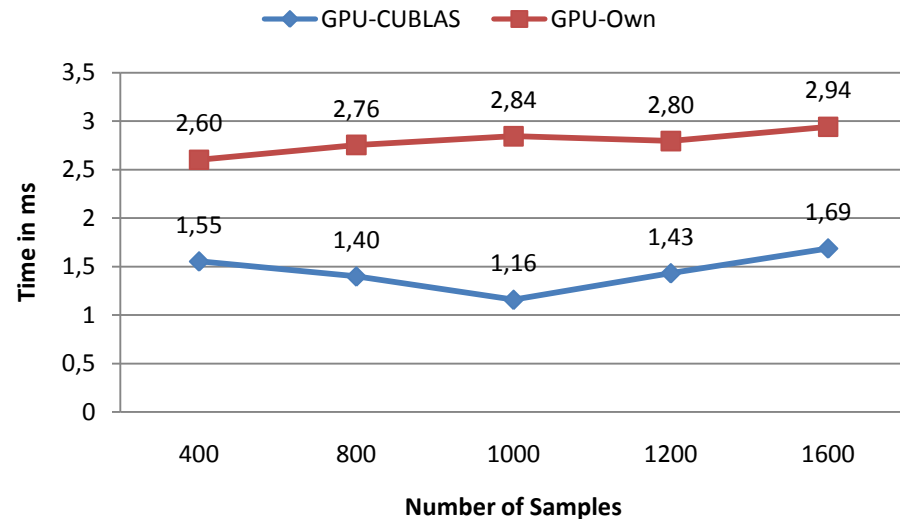
Pre-processing for K-means

Select k-largest Eigen values and their corresponding Eigen vectors form matrix X

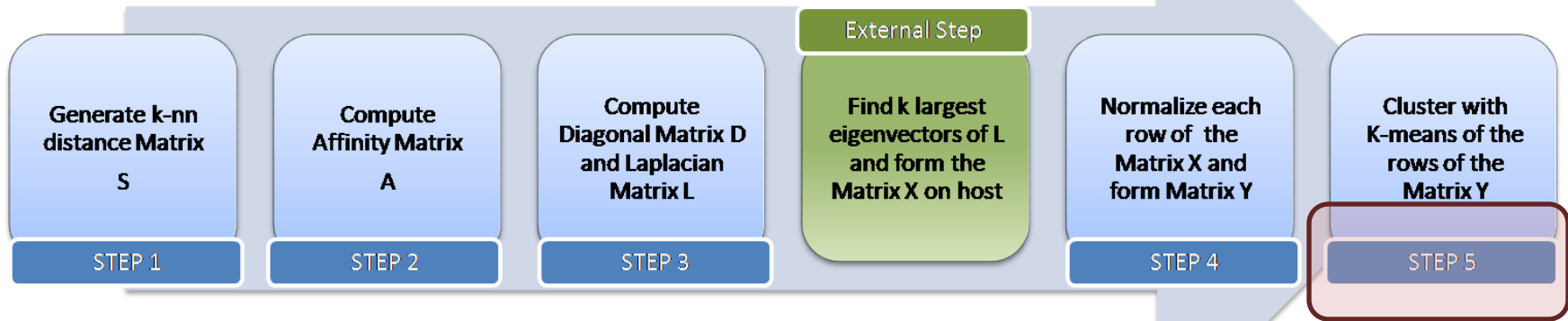
Using Thrust `thrust::sort_by_key`

Normalize rows of matrix X

1. Using CUBLAS `cublas_Sgemm`
2. Own implementation with kernels



*This step is a minor step and can be considered as a preprocessing step of K-means step.



K-means Clustering

Cluster rows of matrix Y with k-means

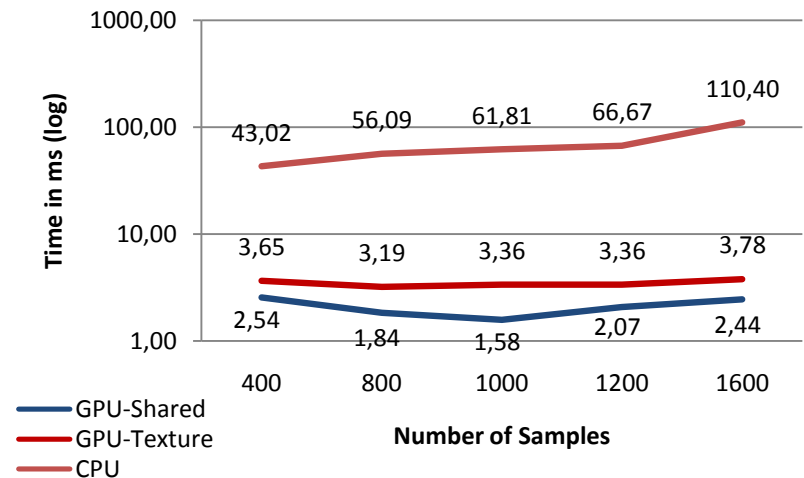
Two different utilization schemas

1. Shared memory reduction with texture
2. Dynamic block shared memory

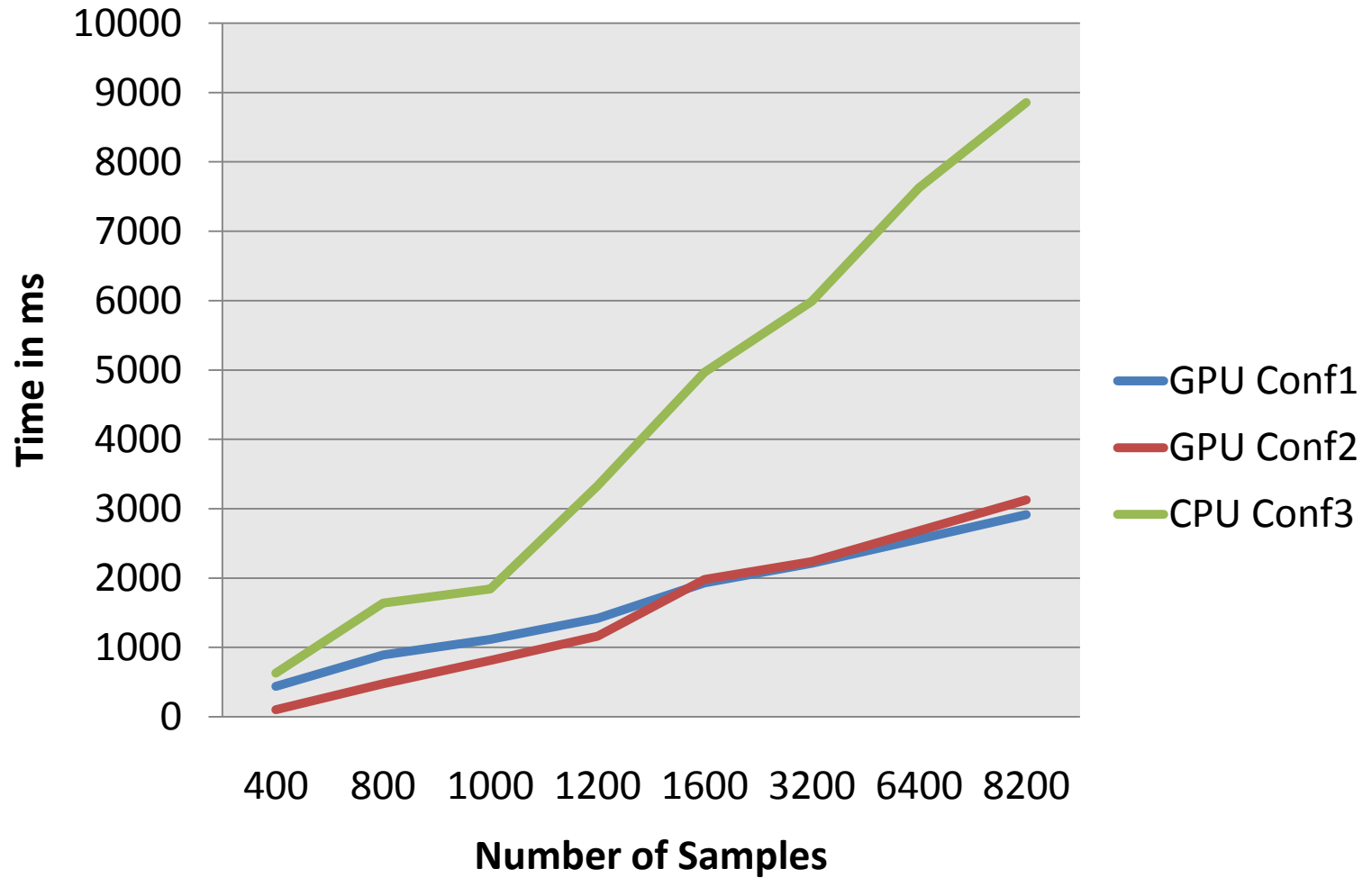
Both have pros. and cons.

1. Suffers low shared memory usage when number of cluster is low
2. Initiates consecutive loops to fill shared memory when number of cluster is high

K-means on CPU vs GPU

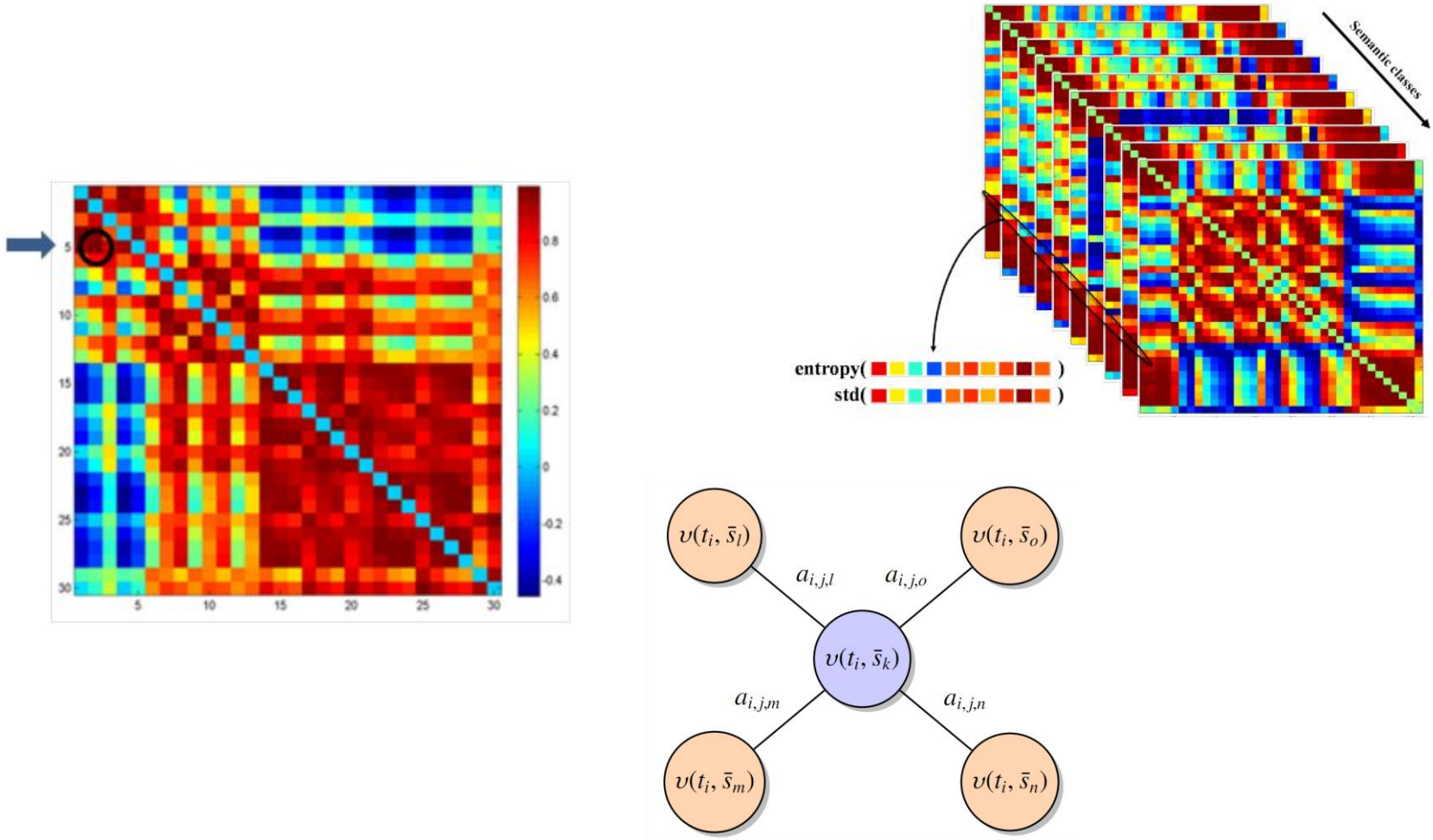


Overall Running Time

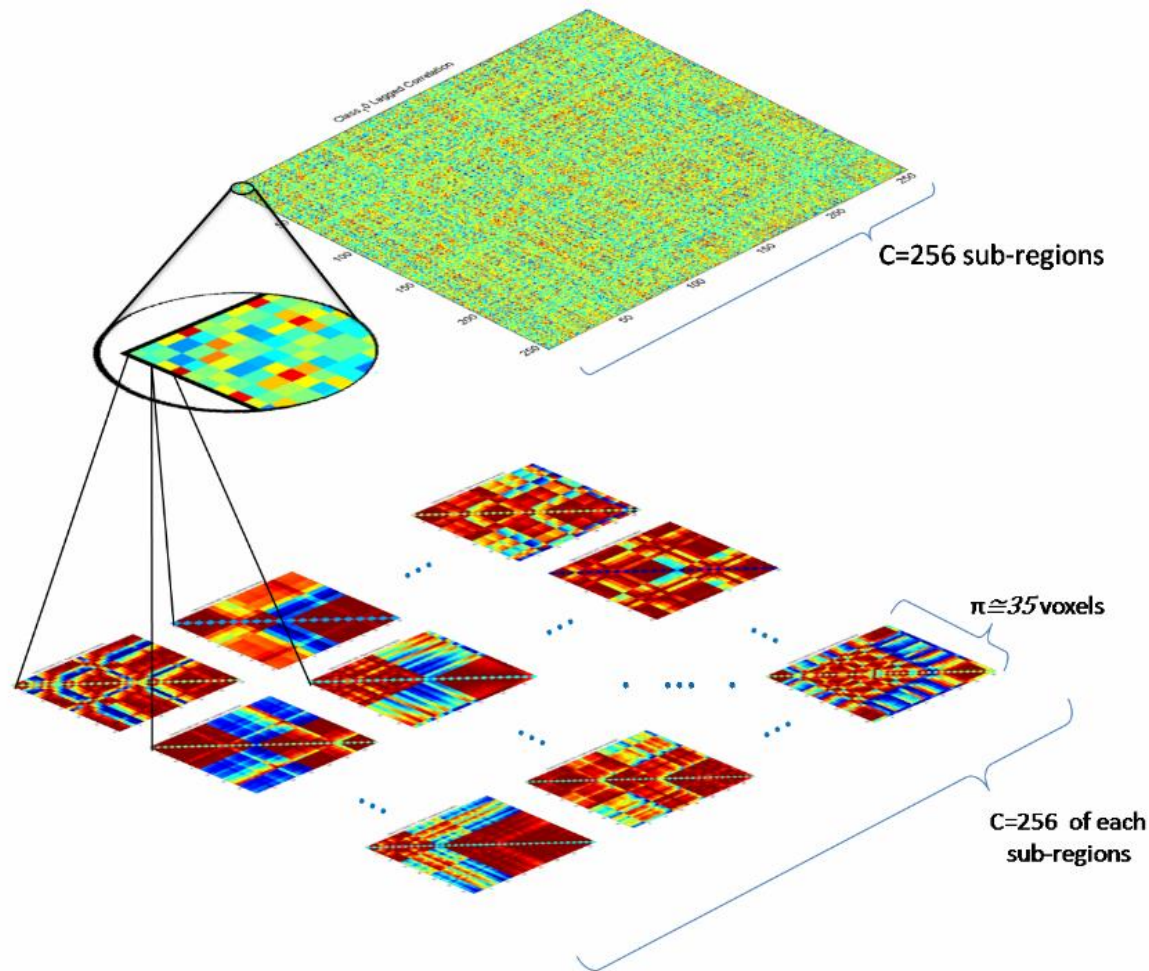


What we “actually” do?

What we “actually” do?



What we “actually” do?



- O. Firat, A. Temizel "Parallel Spectral Graph Partitioning on CUDA" GPU Technology Conference, San Jose, California, 2012
- O. Firat, M. Ozay, I. Onal, I. Oztekin, F. T. Yarman Vural, "Functional Mesh Learning for Pattern Analysis of Cognitive Processes", 12th IEEE International Conference on Cognitive Informatics and Cognitive Computing (ICCI*CC), 2013

What we “actually” do?

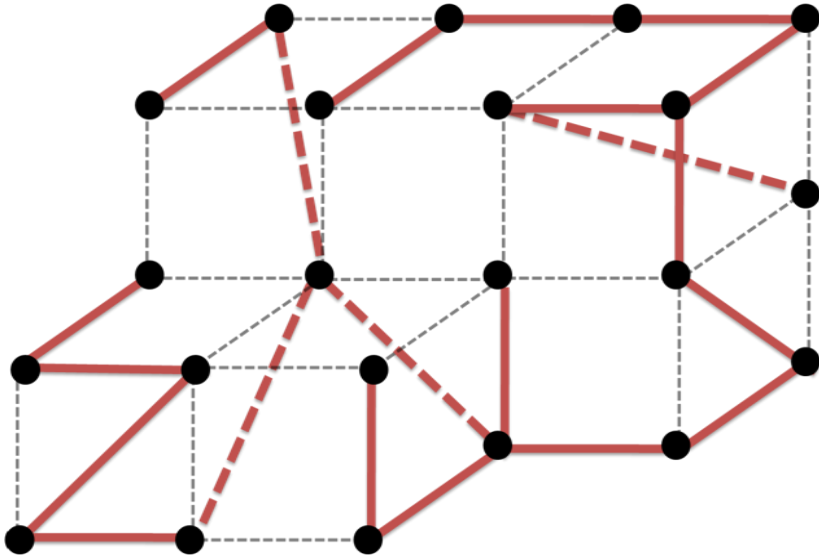


TABLE I. CLASSIFICATION ACCURACIES FOR EACH METHOD

Employed Method	Number of Features	Classifier Accuracy	
		<i>k</i> -NN	SVM
Classical MVPA Method*	8142	44.77%	39.75%
MST-F (6 neighbors)	8141	54.39%	58.16%
MST-F (18 neighbors)	8141	58.16%	59.83%
MST-F (26 neighbors)	8141	59.83%	61.09%

* Voxel intensities are directly fed to classifiers as features.

- O. Firat, M. Ozay, I. Onal, I. Oztekin, F. T. Yarman Vural, "Representation of Cognitive Processes Using the Minimum Spanning Tree of Local Meshes", 35th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBS), 2013.

What we “actually” do?

$$\varepsilon_{i,j,p}^2 = \left(v(t_i, \bar{s}_j) - \sum_{\bar{s}_k \in \eta_p(\bar{s}_j)} a_{i,j,k} v(t_i, \bar{s}_k) \right)^2. \quad (2)$$

By taking the average of squared errors for all time instants t_i and for all seed voxels $v(t_i, \bar{s}_j)$, we approximate the variance of the error for the mesh size p as follows :

$$E(\bar{\varepsilon}_p^2) \cong \frac{1}{N} \frac{1}{M} \sum_{i=1}^N \sum_{j=1}^M \varepsilon_{i,j,p}^2, \quad (3)$$

where $E(\cdot)$ is the expectation operator. This expected squared error is used to determine the Akaike's Final Prediction Error, FPE, in space-domain such that:

$$FPE_p = E(\bar{\varepsilon}_p^2) \left(\frac{M+p+1}{M-p-1} \right). \quad (4)$$

TABLE II. CLASSIFICATION ACCURACIES FOR ESTIMATED MESH SIZE AND CLASSICAL MVPA METHOD

Participants	Estimated optimum mesh size	<i>k</i> -NN Accuracy		
		Classical MVPA method	Classification with arc vectors using estimated mesh size	Average accuracy of the classifiers for $p \in [2-25]$
Participant 1	17	58%	66%	61%
Participant 2	23	58%	67%	61%
Participant 3	24	62%	60%	61%
Participant 4	25	53%	58%	57%
Participant 5	23	54%	59%	57%
Participant 6	16	53%	59%	57%
Participant 7	25	57%	56%	55%
Participant 8	17	57%	58%	57%

I. Onal, M. Ozay, O. Firat, I. Oztekin, F. T. Yarman Vural, "Analyzing the Information Distribution in the fMRI measurements by estimating the degree of locality", 35th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBS), 2013.

What we “actually” do?



TABLE III: Test accuracies of classifiers, * indicates Gaussian Kernel

	SVM		SVM*		k -NN		NB	
	Run1	Run2	Run1	Run2	Run1	Run2	Run1	Run2
MVPA [1]	53%	83%	53%	50%	53%	57%	53%	50%
LRF	60%	73%	60%	57%	50%	60%	63%	57%
FC-LRF	73%	90%	67%	83%	67%	70%	70%	57%

Where we come in handy?

- *libFCL* GUI
- Scalable Functional Connectivity Tools

The screenshot displays the *libFCL* GUI interface. At the top, a progress bar indicates the workflow steps: 1. Load Data, 2. Clustering, 3. Connectivity, and 4. FC. A dialog box titled *libFCL* shows a green checkmark and the message "Classification Successful".

The main interface is divided into several sections:

- Classification:** A "Select Classifier" dropdown menu is set to "K-Nearest Neighbors". A "Run" button is visible below the classifier options.
- Training Features:** A list of features is shown, including "A_TR_N_500_C_1_NONE_P_4_NN_LRF", "A_TR_N_500_C_8_PEAR_ALL_P_4_NN_POS", and "A_TR_N_500_C_8_PEAR_ALL_P_4_NN_NEG".
- Test Features:** A list of features is shown, including "A_TE_N_500_C_1_NONE_P_4_NN_LRF", "A_TE_N_500_C_8_PEAR_ALL_P_4_NN_POS", and "A_TE_N_500_C_8_PEAR_ALL_P_4_NN_NEG".
- Status Summary:** Displays classification metrics: "Classification completed for 2 classes", "Recall : 0.95833", "Precision : 0.96154", and "Fscore : 0.95826".
- Performance Results:** Includes a "Confusion Matrix" button, "ROC", "p-Test", and "Accuracy" buttons.

A "Confusion Matrix" window is open, showing a 2x2 matrix with the following values and percentages:

	1	2	
1	22 45.8%	0 0.0%	100%
2	2 4.2%	24 50.0%	92.3%
	91.7% 8.3%	100% 0.0%	95.8% 4.2%
	1	2	Target Class

The "Performance Results" section includes a "Confusion Matrix" button, which is highlighted with a blue circle. Below it are buttons for "ROC", "p-Test", and "Accuracy".

A note at the bottom of the "Confusion Matrix" window states: "User can also visualize the classification results on a confusion matrix by clicking the corresponding button in the 'Performance Results' section."

neuro.ceng.metu.edu.tr

orhan.firat@ceng.metu.edu.tr