# Analytic (Harmonic) Geometry of subsets

*Learning "Dual" geometries of Databases/Matrices and Tensors.*

*Harmonic Analysis and Geometric Measure theory .*
*CIRM Luminy October 2017 . "Happy birthday Guy David"*

*R. Coifman ,*

*Department of Mathematics ,*
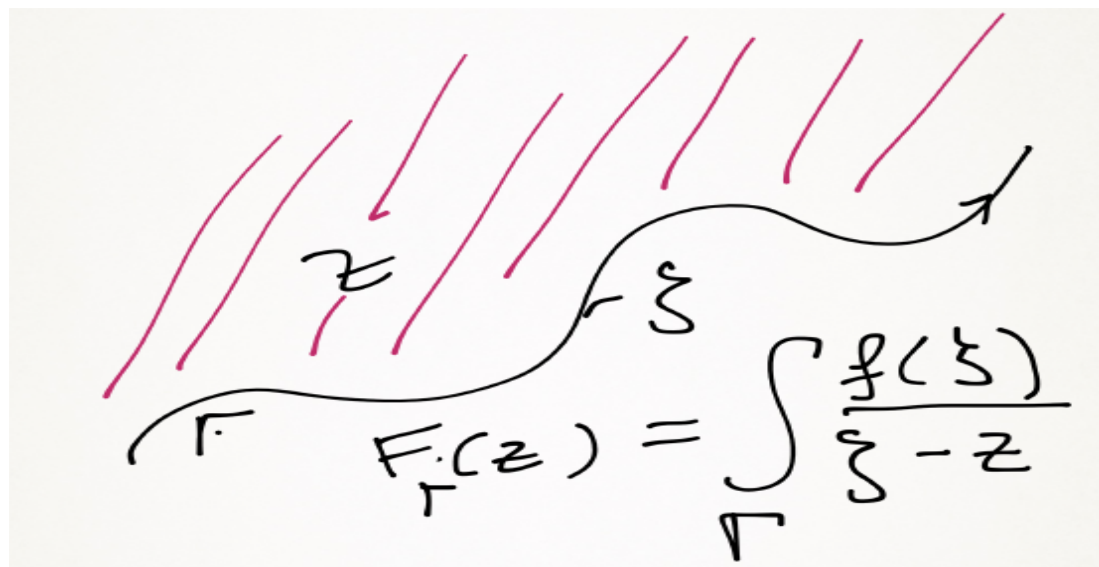*and program of Applied Mathematics*
*Yale University*

The geometries of a subset of Euclidean space can be characterized by the properties of operators defined on functions on the subset . Through the restriction of classes of functions to the subset .
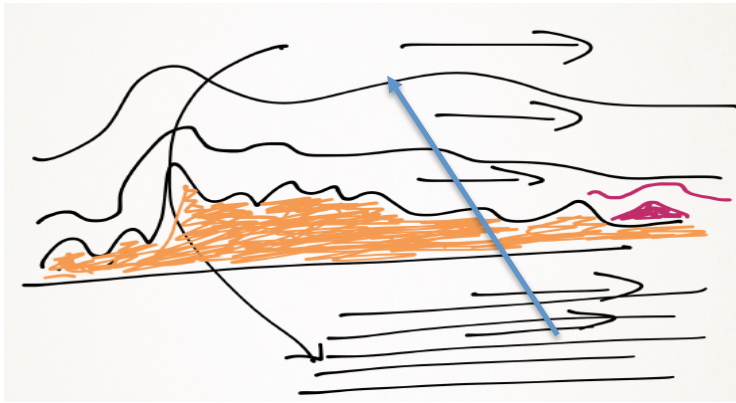
The most elementary example are the coordinate functions ,

More interesting examples are the electrostatic or acoustic, potential fields generated by charges supported on the set ,

The relationship between the geometry of the boundary of a domain and properties of the space of solutions to boundary value problems has a long history .

Here we focus on learning the geometry from the properties of extensions of functions on the set as solutions to some PDE. For example, the Cauchy transform, which is an oblique projection on the space of holomorphic functions on one side of a curve



$$F(z) = \int_\Gamma \frac{f(\zeta)}{\zeta - z}$$

The Riemann mapping from the upper half plane to the region above the curve, maps the horizontal lines to flow lines above the curve (a river bed). The Cauchy transform is easily converted through the Kerzman-Stein formula into the Szego orthogonal projection which , itself yields the Riemann mapping.

$$C(f)(z) = \frac{1}{2\pi i} \int_{\Gamma} \frac{f(w)}{w-z} dw$$

$$S = C(I + C - C*)^{-1}$$

We also,connect the curve Geometry , and distance between two curves to the distance between two Cauchy or Szego operators.

$$let \ \Gamma_1 \ \Gamma_2 \ be \ two \ curves \ parametrized \ by \ arc \ length$$

with arguments $\alpha_1(s), \ \alpha_2(s)$

$C_1, C_2$ the corresponding Cauchy transforms , then
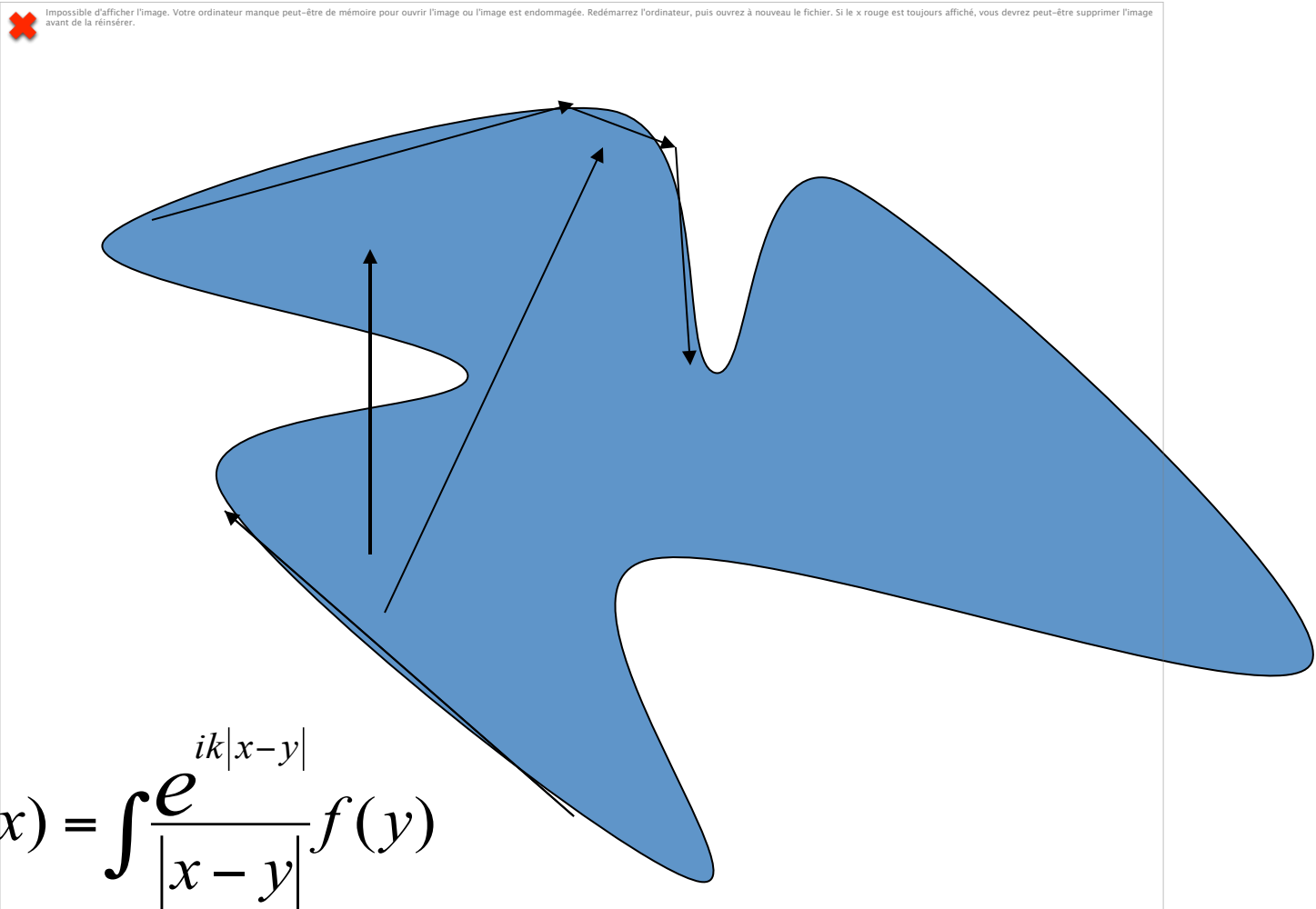
$$\left\| C_1 - C_2 \right\| \approx \left\| \alpha_1(s) - \alpha_2(s) \right\|_{BMO}$$

Similar results were obtained  in higher dimensions analyzing the restriction of Calderon Zygmund operators to subsets of $\mathbb{R}^n$ Guy David  , Steven Semmes , characterized the sets  on which these operators are bounded on $L^2$ .

These results and related detailed geometries of Guy David sets were also developed by Peter Jones and Raanan Schul , characterizing them by Carleson like conditions ( beta numbers and their variants).

Acoustic scattering off objects requires detailed effective
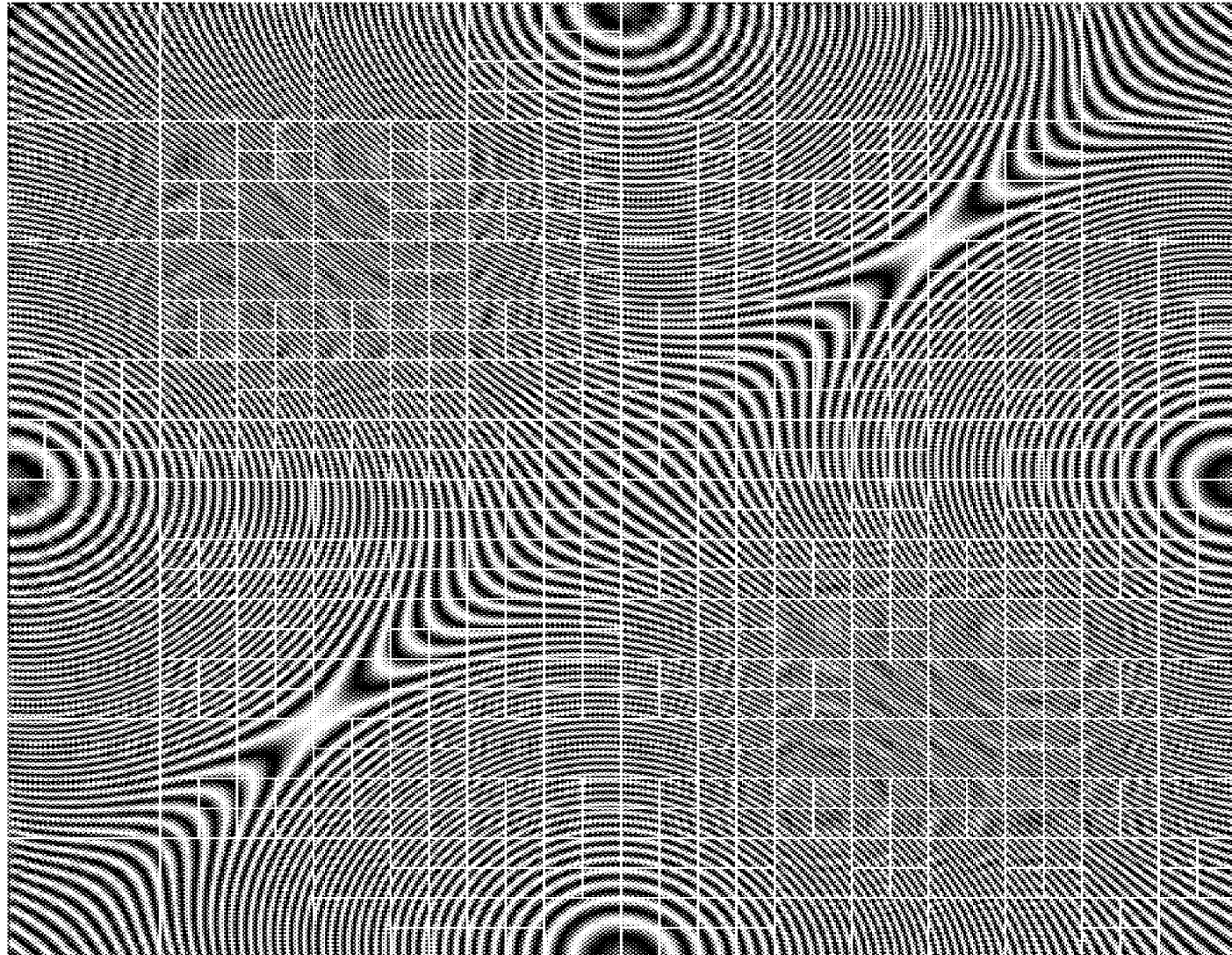Field  interactions between regions on boundary.

$$T(f)(x) = \int \frac{e^{ik|x-y|}}{|x-y|} f(y)$$

The first approximation is given by geometric optics , or Billiards and is
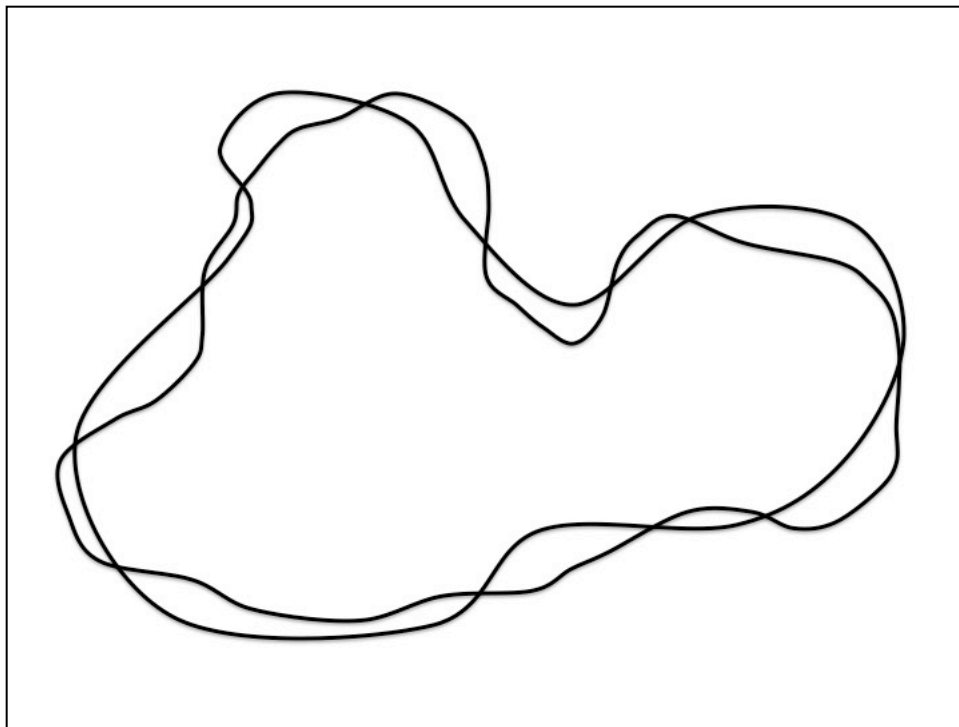obtained automatically through geometry learning.

# Acoustic scattering matrix off an ellipse ,while dense ,the number of parameters (features) needed to describe it is small.each box encapsulates geometric optics interaction.

The same analysis
 could be otained
by local SVD
analysis to track
 rank of interactions
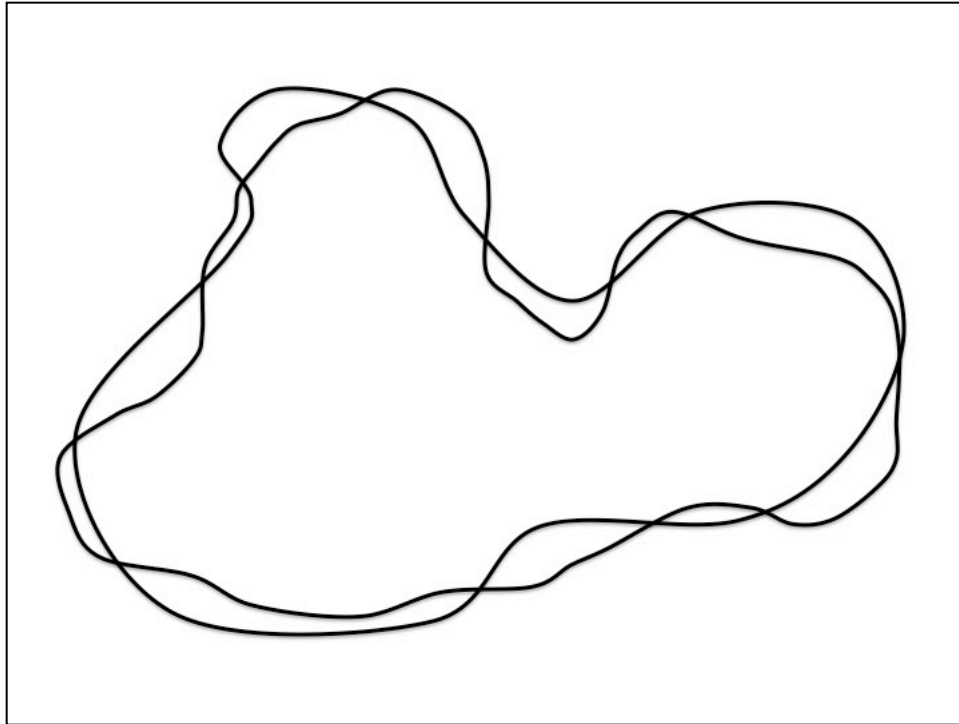The rank is one
at the geometric
 optic level

# Distances between subsets ?

- How to measure the distance between curves ?
- Viewing the points as defining a distribution, and taking a distance to be a dual norm for some function spaces
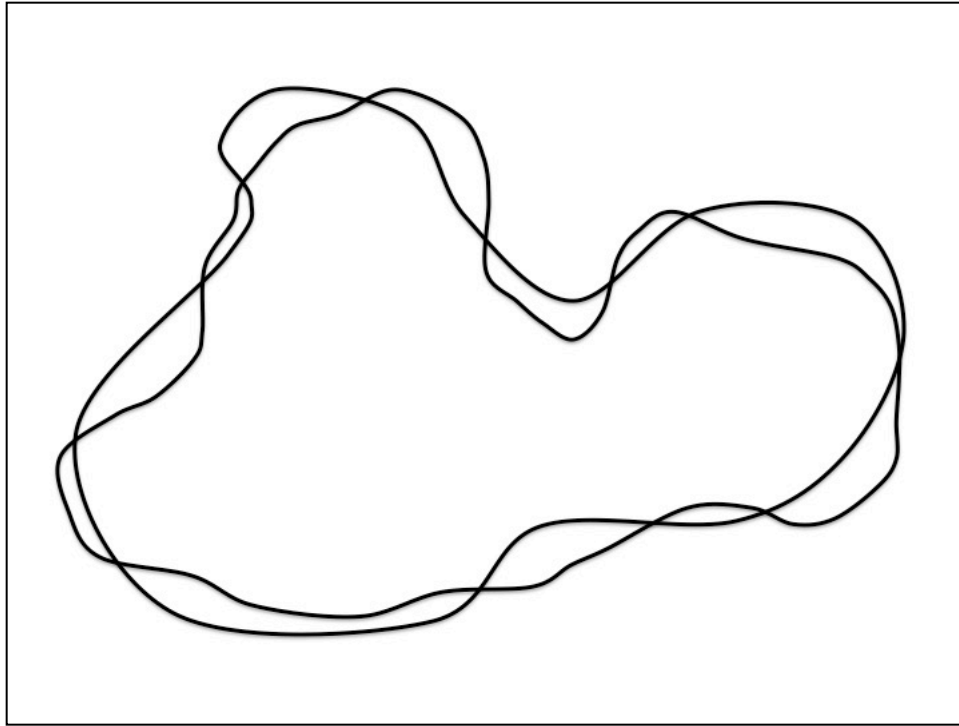
# Earth Mover's Distance

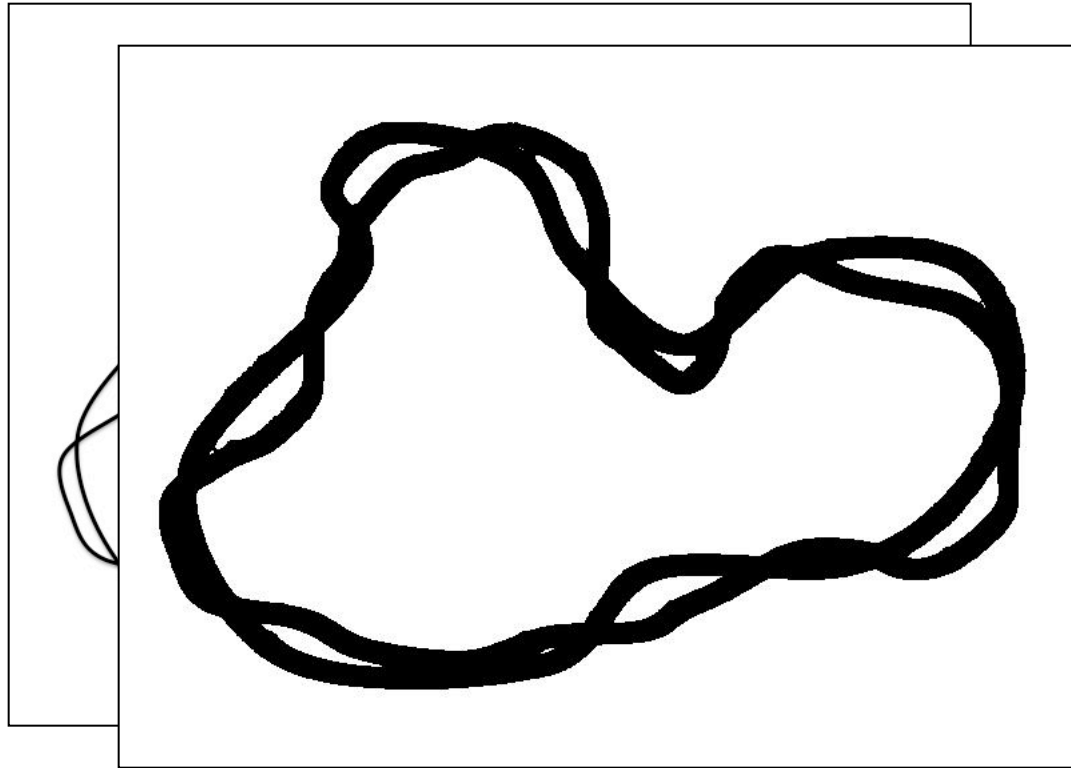- Didn't gain much popularity due to high computational complexity

Earth Mover's Distances or transportation distances measure the total cost of moving mass form one curve to the other, this optimization problem is by duality,equivalent to measuring the distance between the curves as being in the dual to Lipschitz or Holder .

- Efficient implementation via "filtering":
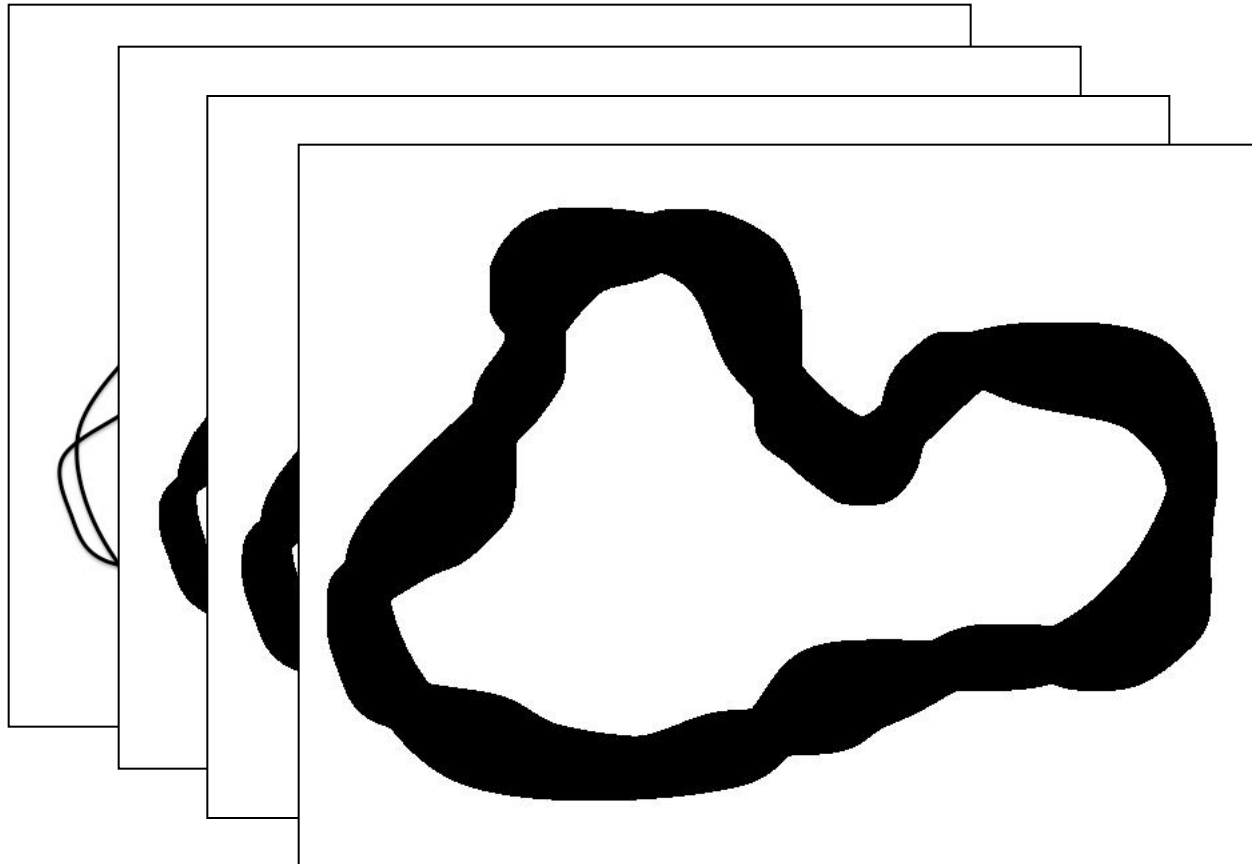  obtaining coarser and coarser views [Shirdhonkar & Jacobs, 08]

# Earth Mover's Distance

- Efficient implementation via "filtering": obtaining coarser and coarser views [Shirdhonkar & Jacobs, 08]

# Earth Mover's Distance

- Efficient implementation via "filtering": obtaining coarser and coarser views [Shirdhonkar & Jacobs, 08]

# Dual metrics and EMD

*C*onsider images $I_i$ to be sensed by correlation with a collection of sensors f, in a convex set B.

*W*e can define a distance $d_{B^*}(I_i, I_j) = \sup_{f \in B} \int_X f(x)(I_i(x) - I_j(x))dx$

If B is the unit ball in Holder classes we get the EMD distances ,
The point being that if B transforms nicely under certain distortions so does the dual metric.

The computation of the dual norm for standard classes of smoothness is linear in the number of samples. (Unlike the conventional EMD optimization or minimal distortion metrics)

This is applicable to general data sets , such as documents, or profiles .

Morever since dual norms are usually weighted combinations of $l^p$ norms at different scales, it is easy to adjust the weights to account for noisy conditions. (ie redefining smoothness).
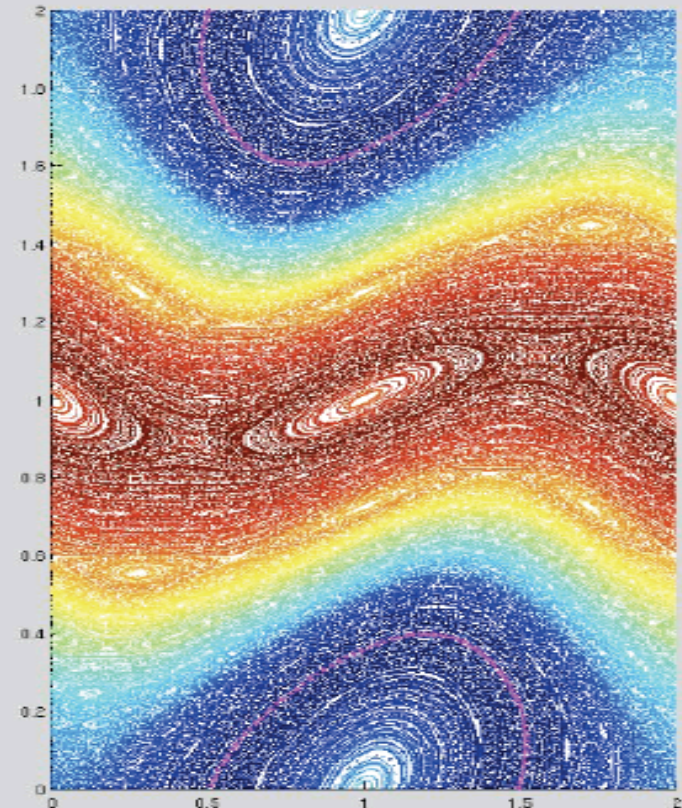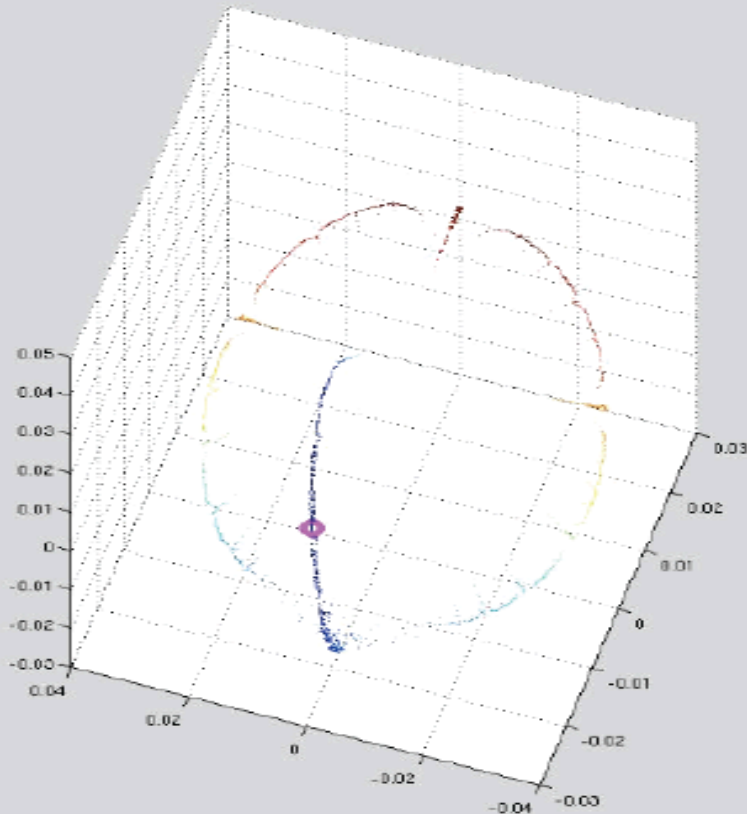
Diffusion ( or prolate function ) embedding of the graph of orbits of the standard map on the torus, each orbit is a measure , we use the earth moving distance to define distances between orbits and organize in a graph.

in particular the parameter alpha as well as the orbit can be easily retrieved from a few points.

$$p_{\ell+1} \triangleq p_\ell + \alpha \sin(\theta_\ell),$$
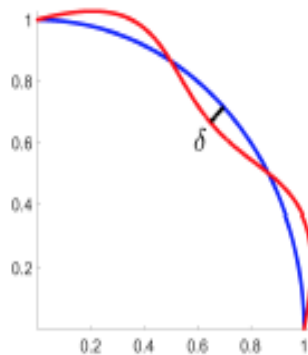
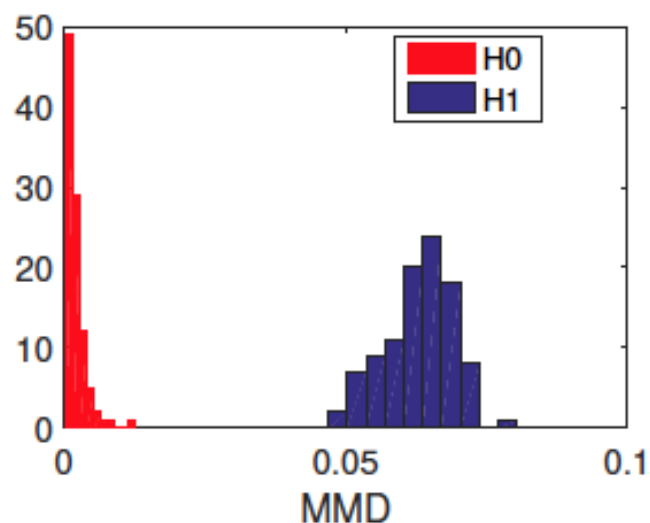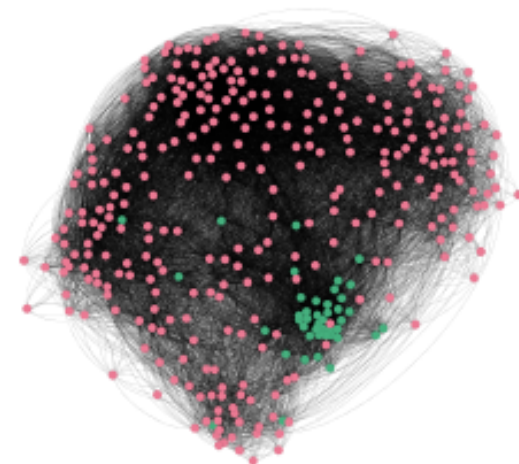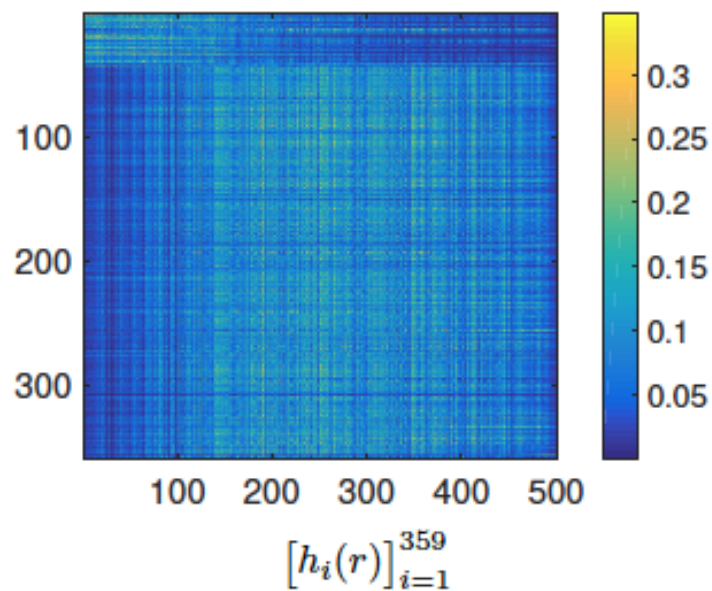$$\theta_{\ell+1} \triangleq \theta_\ell + p_{\ell+1},$$

A fundamental problem in statistical data analysis involves testing the hypothesis that two clouds of points in $\mathbb{R}^n$ are sampled from the same distribution. This is particularly painful when the distribution is unknown .
The ideas described here are useful for that purpose .

The basic Kolmogorov Smirnov test in one variable consists , in the comparison of the number of points in various bins ( intervals), ie our test functions are intervals having sufficiently many points to achieve accuracy.
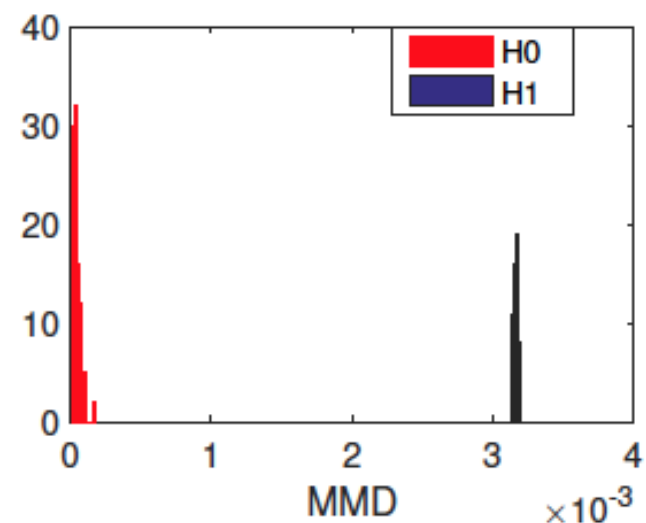
# Two-sample Statistics Based on Anisotropic Kernels

Xiuyuan Cheng[*1], Alexander Cloninger[*2], and Ronald R. Coifman[3]

Figure 9: Top: Unsupervised histograms and clustering of AML patients. Bottom: Two sample test between pool of healthy and unhealthy patients comparing isotropic and anisotropic gaussian.

One of the first applications of wavelet bases ,was the observation that CZ operators could be efficiently implemented in such bases .

Assume more generally that we have the matrix of potentials of a collection of sources located on a spiral, which are evaluated on a flat disk located away . We need to find a wavelet basis on each structure relative to its geometry. The full matrix is then expanded efficiently in the Tensor Wavelet basis .

Observe also that a matrix is usually given in garbled order.

We see that the basic Calderon Zygmund theory organizes both sets through the values of certain collection of functions on them .   The kernel establishes their relations.
 In the case of earth mover distance we view the distance between the sets through the eyes of Holder functions.

]We could also consider band limited functions, and match two sets by comparing their Fourier transforms .
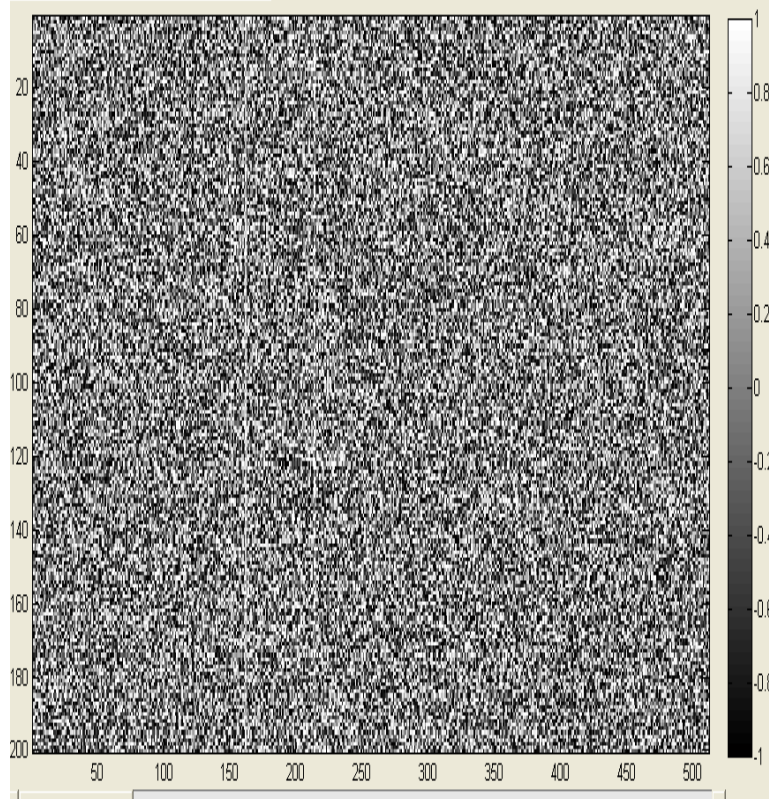 Or just organize a single set by looking at the kernel of the projection on band limited functions and compute its spectral decomposition when restricted to the set  ( as we did for the Cauchy transform).
The prolate functions so derived form a coordinate system on the set and enable its embedding into low dimensions.

Observe also that the kernel of band limited functions in high dimensions is numerically close to a Gaussian kernel.  Leading to prolates which are approximate eigenfunctions of a Laplace operator on the subset whenever it is an embedded manifold, leading to a diffusion geometry embedding.

By duality  we  can also organize  the" geometry" of the  eigenfunctions of the Laplace operator ,  this defines a Heisenberg geometry  of eigenfuncitons.

A permutation of the rows and columns of the matrix **sin(kx)**. On the left we recover the one dimensional geometry of x (which is oversampled ), while on the right we recover the one dimensional geometry
 of k .
More generally we can build a dual geometry of eigenvectors of Laplace Beltrami  operators on manifolds  { example SU(2)}

*No equations, no parameters, no variables: data, and the reconstruction of normal forms by learning informed observation geometries.*
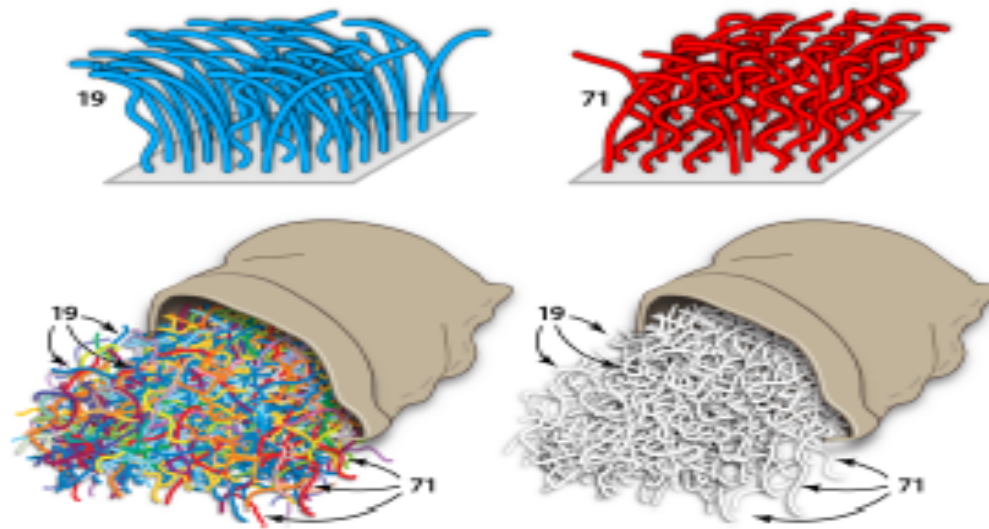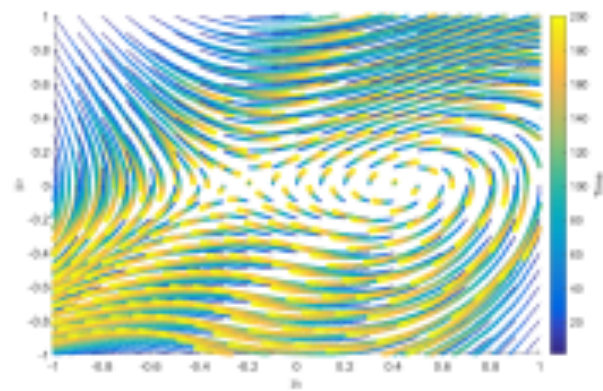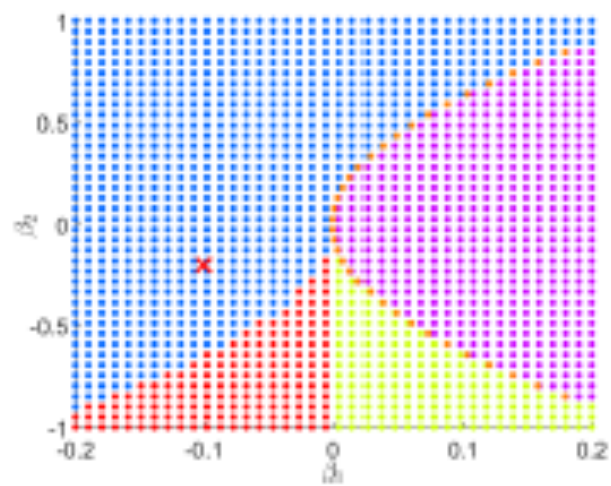
*Or Yair, Ronen Talmon, Ronald R. Coifman and Ioannis G. Kevrekidisc,*
PNAS 2017



**Fig. 2.** We observe an ensemble of short trajectories in what we call *a trial* from a fixed, yet unknown, dynamical regime. For example, one trial is labeled 19 and another trial is labeled 71. We store the short trajectories ("bag" them), and we only keep the label of the trial. Our goal is to empirically derive the dynamical regime information associated with each trajectory (labelled by the color in this case), by (a) deducing the state variables and the associated phase portraits; and then (b) organizing the phase-portraits of all the trials; so as to (c) derive (a tabulated form of) the evolution equation governing the dynamics.

**Fig. 4.** (left) Data-driven embedding of the parameters axis of the observations collected from the Bogdanov-Takens system (colored according to the true bifurcation map). Embeddings built from (a) state variable observations; and (b) observations through a nonlinear invertible function. (right) Data-driven embedding of the state variables axis (c) colored by the initial conditions of $x_1$, and (d) by the initial conditions of $x_2$.

# Coupled Pendulum

- A movie of the system for example:

# Coupled Pendulum

- Two normal modes:



$$\cos\left(\sqrt{\frac{g}{L}}\,t\right)$$

$$\cos\left(\sqrt{\frac{2k}{m}+\frac{g}{L}}\,t\right)$$

# Coupled Pendulum

- Time-varying spring constant:

# Coupled Pendulum

- "Scrambled" movie as input:

# Coupled Pendulum

- Result:

# "Empirical Physics"

$$\cdots \quad \cos\left(\sqrt{\frac{2k}{m} + \frac{g}{L}t}\right)$$

$$\text{--} \quad \cos\left(\sqrt{\frac{g}{L}t}\right)$$

# "Empirical Physics"

$$\cdots \quad \cos\left(\sqrt{\frac{2k}{m} + \frac{g}{L}} t\right)$$

$$\text{---} \quad \cos\left(\sqrt{\frac{g}{L}} t\right)$$

# Sensory-Motor Integration in the Mammalian Brain: experiment, dataanalysis and modeling
## A joint project with
**Gal Mishne , Jackie Schiller, Ronen Talmon and Ron Meir, Uri Dubin, Technion - IIT, Israel**



the setting of our tri-geometry analysis of the collected trial-based neuronal activity from the motor cortical region. These measurements were taken from a behaving mouse in a single day of experiments. The data is composed of 60 trials, where after the first 20 trials, the activity of the somatosensory region was silenced by pharmacogenic activation. A single trial consists of 12 seconds, during which 120 frames are measured. The recordings are taken from 121 neurons located in M1 cortex

# Experimental Setting

The Neuronal dynamics of mouse as it learns to accomplish a task, or suffers from a neural pathology



The data consists of a few hundred experiments, in which the mouse repeatedly performs a task, initially learning to pick up food when the bell rings, being adept at it, then is dis functional due an infection, and recovers later

Each experiment is a data base which when combined is a three dimensional array.

Trial 1 , t=4

Trial 1 , t=15

Trial 6 , t=15

Here we present the multiscale hierarchical organization of the 2D slices of all the neurons (Fig. 4 (right)) in a flexible tree. Same data as in previous figures. To demonstrate the organization obtained by the tree we highlight several interesting tree folders and present the 2D neuron slices under each folder of interest. We observe that indeed the neurons are grouped together according to similar properties. While this result is very preliminary and the full meaning of such an organization should be examined in depth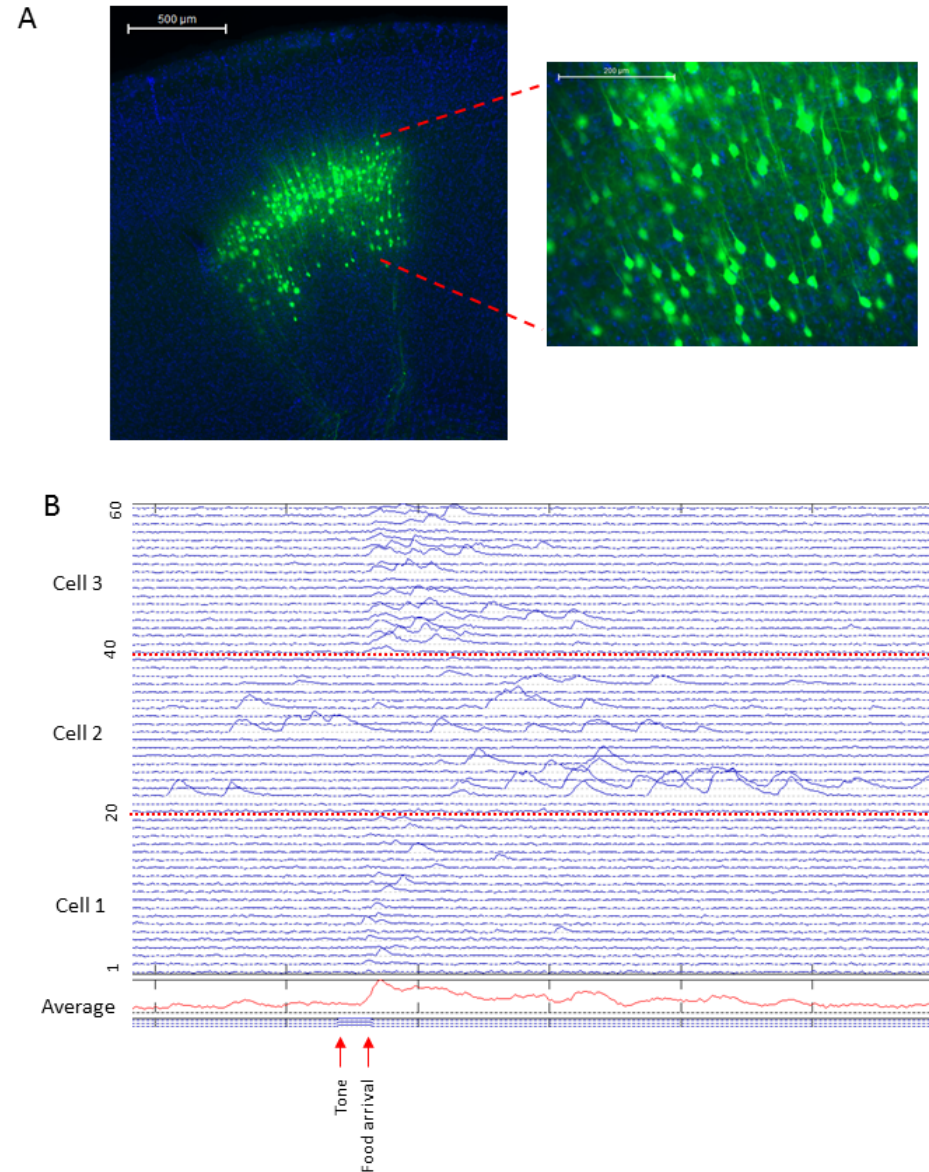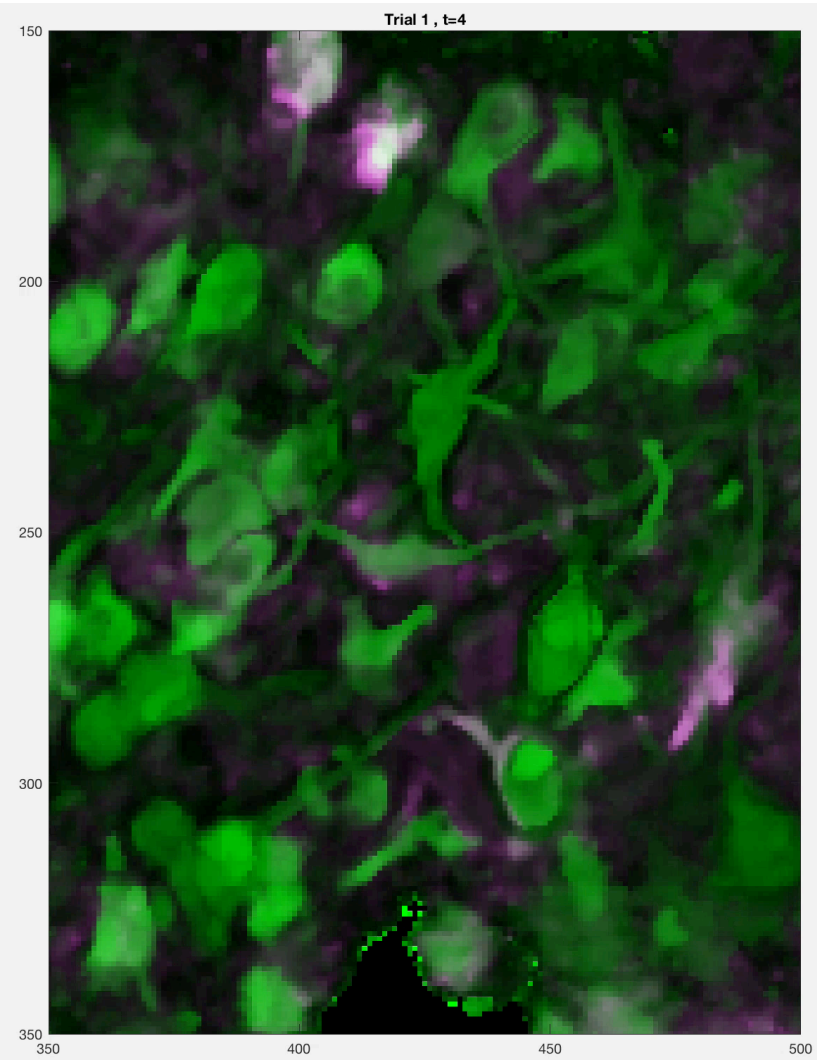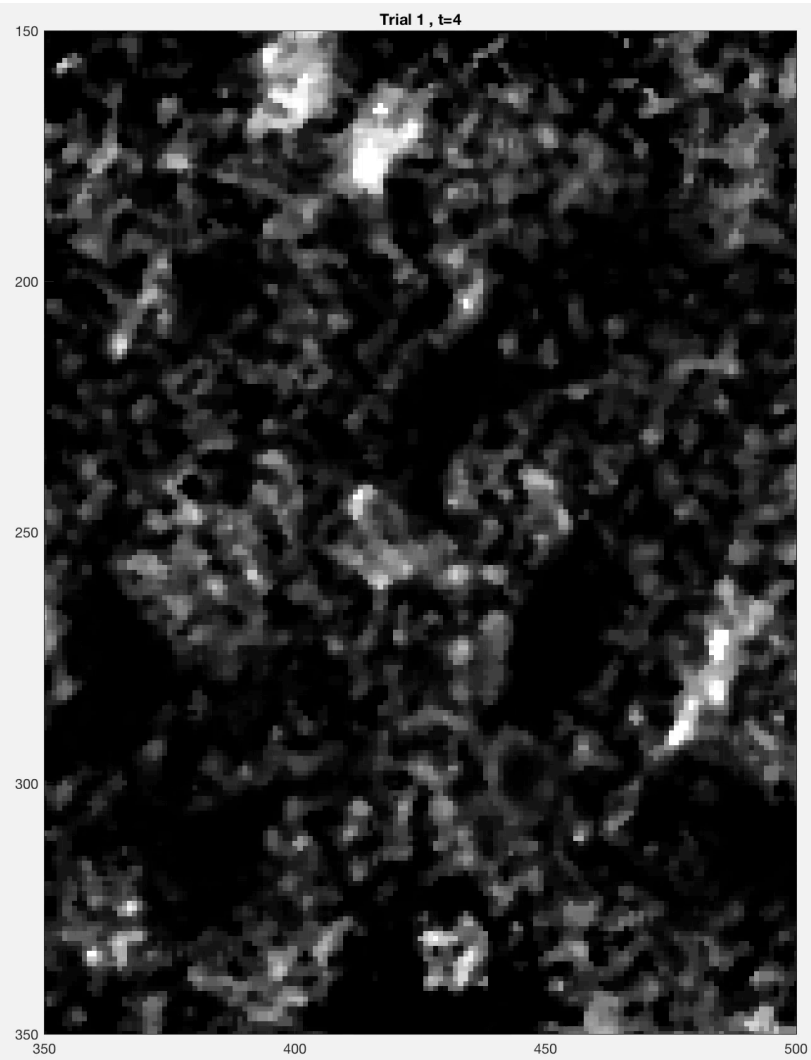, we can immediately detect few interesting features. For example, the orange folder (b) consists of neurons, which are active only during trials under the effect of somatosensory silencing (horizontal separation). The yellow folder (a) consists of neurons, which are active only at or after the tone (vertical separation), and mostly in trials under the effect of the silencing (horizontal separation). In contrast, the purple folder (c) consists of neurons, which are active after the tone but during trials without the silencing effect. Finally, the green folder (d) consists of neurons, which were silenced by the manipulation. We note that folders consisting of neurons, which are active only during a particular behavior (such as grabbing the food), were also identified, but are not displayed due to space constraints

# *Mathematical challenges for "emergent" data organization and knowledge building.*

*The following items are the traditional functions of the librarian, they now need to be performed, on a massive scale, for heterogeneous digital documents.*

- The initial tasks , are storage, retrieval, and relevance search.

- In view of massive amounts of continuously generated data , this needs to occur automatically, in a data agnostic way .

- In particular both "context" and "concept" need to emerge spontaneously from the data, once goals are defined.

- An analog of a geographic Data Atlas , in which heterogeneous documents are organized by context into different "geometries" of Knowledge.

*Methodologies for principled mathematical learning of empirical functions, or regression, in an environment where most measured parameters are irrelevant or only marginally relevant are needed.*

*The community of machine learning have come up with remarkable useful tools such as various deep learning neural nets , to enable filtering irrelevant data and to build useful data features for coding ,and classification of complex data .*

*Unfortunately not much theory is available to insure reliability or precision, more critically no explanatory model is associated with the "black box"*

*The challenge is to render these methods transparent , and explanatory.*

• Harmonic Analysis has over the last 60 years focused on the relationship between geometry , and appropriate representations, as a tool to understand and prove estimates on operators .  In particular kernels of operators restricted to subsets of Euclidean space have played a fundamental role in understanding the geometry and combinatorics of the set.

•We claim that these methodologies open the door to organization of matrices viewed as either databases, or as linear transformations.

*The challenge is to organize a database or a matrix without any a priori knowledge of its internal model, in particular can we find data anomalies, fill in missing entries build classifiers and in general build  data agnostic, analytic mathematics for processing any kind data.*

•Agnostic data geometerization, enables automation of data organization and fusion +analytical intelligence.
Like a good memory organization, we would have the first step to ab initio learning, learning in which we  have a feedback mechanism to reorganize the data according to the inferences we wish to achieve.

- The main analytical challenge is to simultaneously build a graph of columns and a graph of rows so that the matrix entries are as smooth (or predictable )as possible, relative to the tensor product of these geometries.  This smoothness is measured in terms of an appropriate  tensor Besov norm or entropy .

- The next challenge is to enable simple reorganization to achieve regression or machine learning, or fast numerical analysis.

The underlying analytical methods enables filtering out anomalous responses , and provides detailed quantitative assessments of consistency of responses .

The analysis-synthesis tools, that enable the geometric construction, are useful to provide a metric to assess success in organizing the data base.

We extend ideas of Harmonic Analysis and approximation theory to the study of general matrices,( or higher order tensors,) whether the goal is organization of a data base to extract knowledge, or to build a representation relative to which a matrix is efficiently described.
We illustrate the outcome of such organization on the MMPI ( Minnesota Multiphasic Psychological Inventory) questionnaire .
The Tensor Haar Bases enable filtering out anomalous responses , and provide detailed "analysis" (pun intended) .

Stromberg's observations about the efficiency of approximation of functions of bounded mixed variation in the tensor Haar basis is particularly useful in the statistical data analysis context of analysing a data base

*Start by considerimg the problem of unraveling the geometric structure in a matrix. We view the columns or the rows as collections of points in high dimension whose geometry we need to define.*
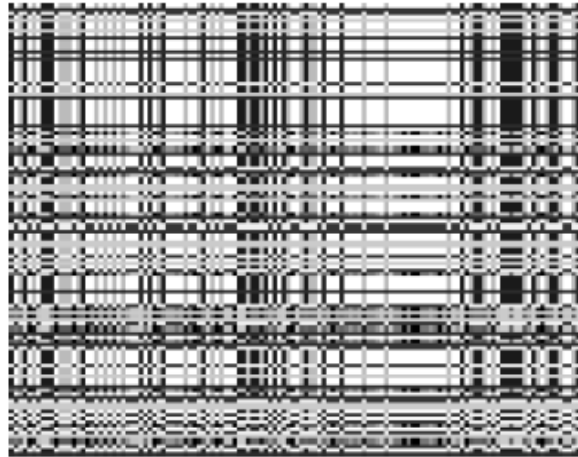
Figure 4.2: A permutation of the matrix $A$.

*The matrix on the left is a permutation in rows and columns of the matrix below it .*

*The challenge is to unravel the various simple submatrices .*

More generally assume that the function represents a probability field which has be garbled by permuting rows and columns. At each pixel we toss a coin with corresponding probability .

The Challenge is to recover the underlying field with some accuracy control.



Probability Field

Shuffled Probability Field

Data Realization

Questionnaire Reconstruction (error = 56.03)

The simplest joint organization is achieved as follows

Assuming an initial hierarchical organization of the columns of the database (see later) into contextual folders ( for example groups of responders which are similar at different "scales" ) use these folders to assign new response coordinates to each row (question), for example an average response of the demographic group.
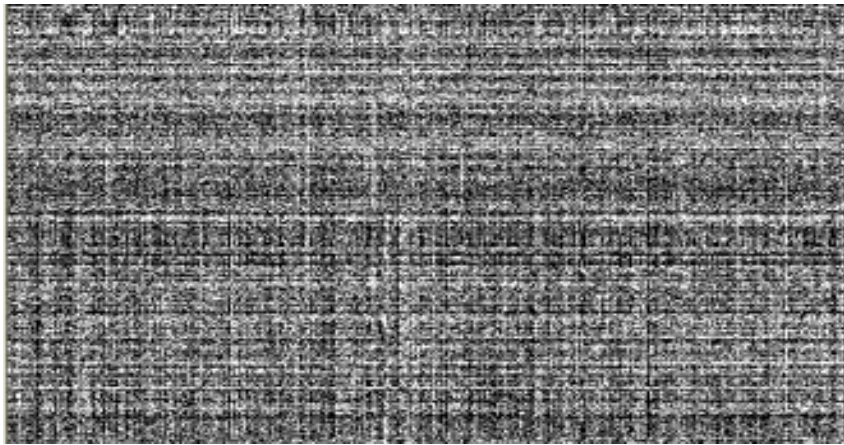
Use the augmented response coordinates to organize responses into a conceptual hierarchy of folders of rows which are similar across the population of columns.

We then use the conceptual folders to augment the response of the columns and to reorganize them into a more precise contextual hierarchy .
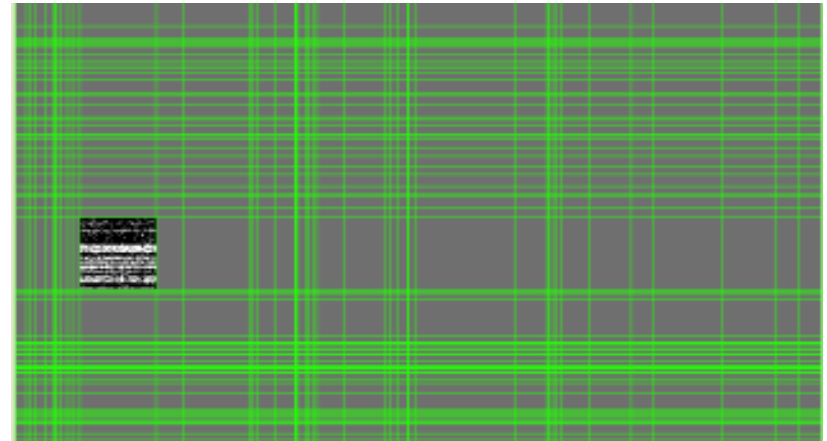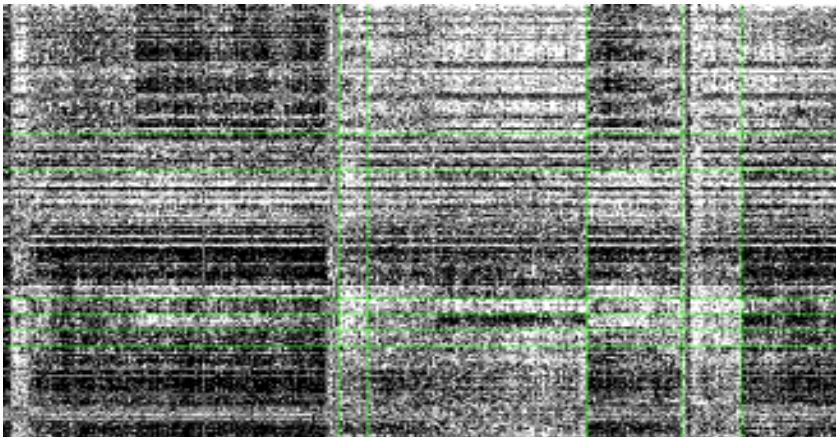
 This process is iterated as long as an "entropy " of the database is being reduced .

The challenge is to organize a data base by organizing both rows and columns simultaneously , if the columns are observations and the rows are features or responses. We organize observations "contextually" and responses "conceptually " each organization informs the other iteratively.



A disorganized questionnaire ,on the left, the columns represent people , the row are binary questions. Mutual multiscale bi learning , organizes the data, bottom left , The questionnaire is split on a two scale grid below. Showing  in the highlighted rectangle , the consistency of responses of a demographic group (context) to a group of questions (concept)

Consider the example of a database of documents , in which the coordinates of each document , are the frequency of occurrence of individual words in a lexicon. Usually the documents are assumed to be related if their vocabulary distributions are "close" to each other.

The problem is that we should be able to interchange words having similar meaning and similarity of meaning should be part of the document comparison .

By duality if we wish to compare two words by conceptual similarity we should look at similarity of frequency of occurrence in documents, here again we should be able to interchange documents if their topical difference is small.

*There are at least three challenges which we claim can be resolved through Harmonic Analysis ;*

**1.Define good document content flexible-distances , and simultaneously good conceptual vocabulary distances.**

**2. Develop a method which is purely data driven and data agnostic ,**

**3.The complexity of calculations should scale linearly with data size.**
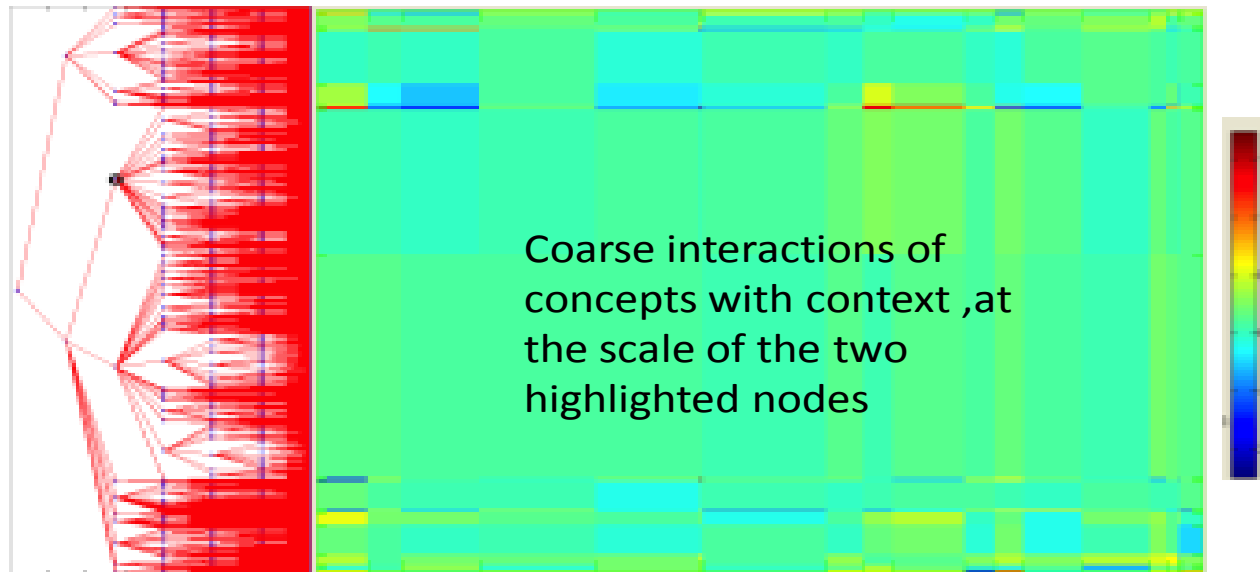
We start by discussing metrics

**Mutual Organization / Tree Structures for context- concept duality,**
**Although we use linguistic analogies these trees were built on time series of observations of 500 objects , the concepts are scenarios of times with similar responses among the population while the contexts are group of objects with similar temporal responses.**



**Concept tree**

Coarse interactions of concepts with context ,at the scale of the two highlighted nodes
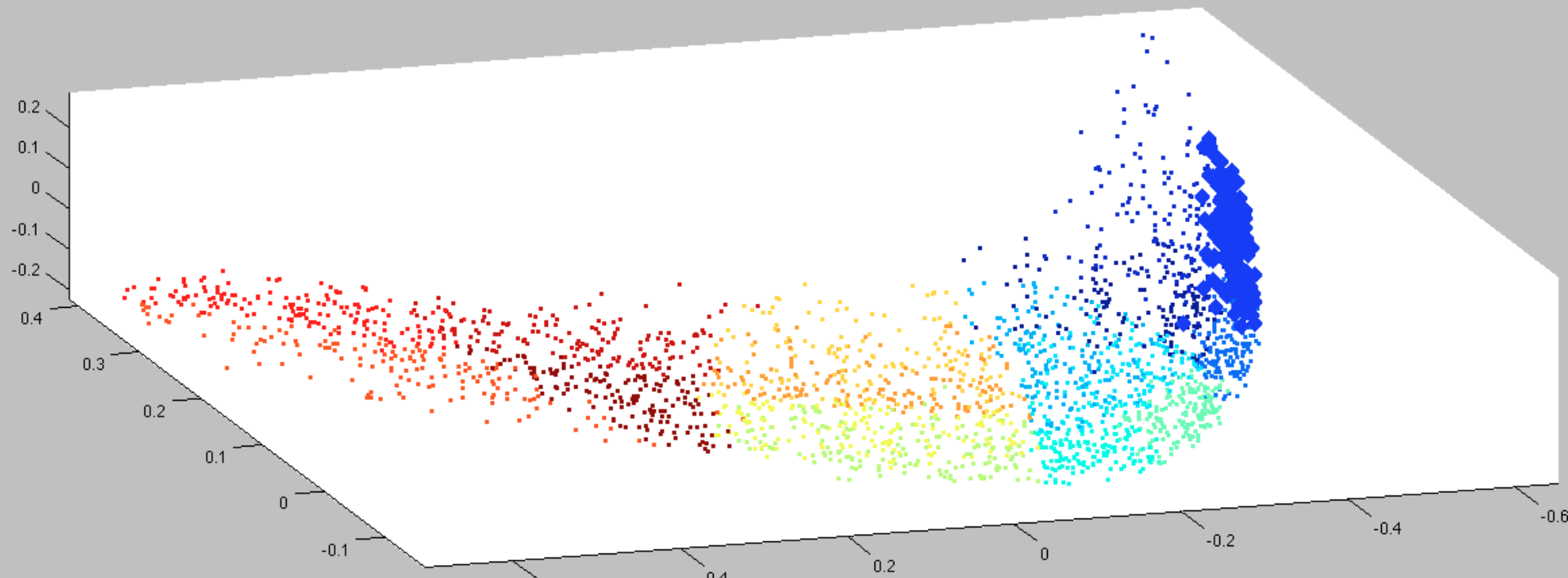
Activity Level Map
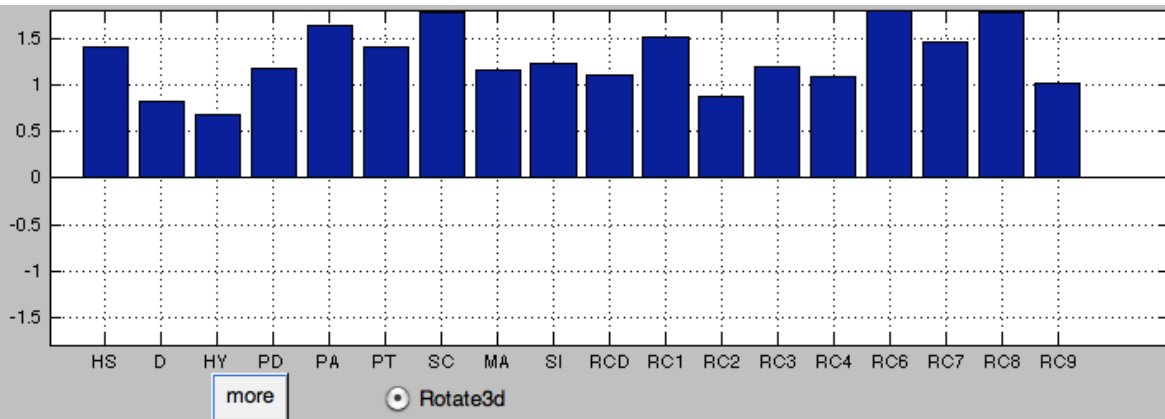
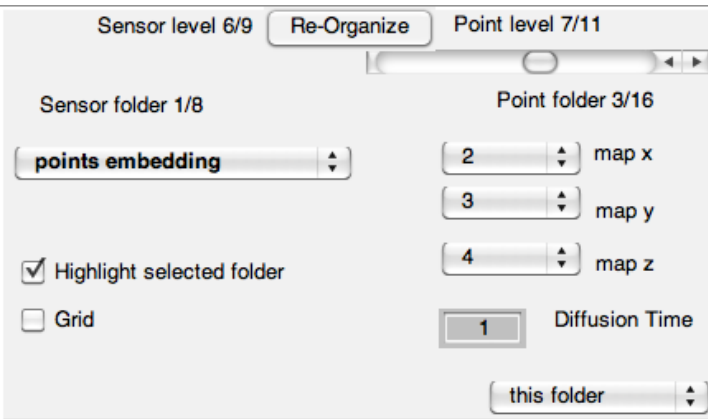**Affinity based Groupings:**

**Concepts - clusters of "words" at different levels of abstraction as they relate to various documents clusters =contexts**
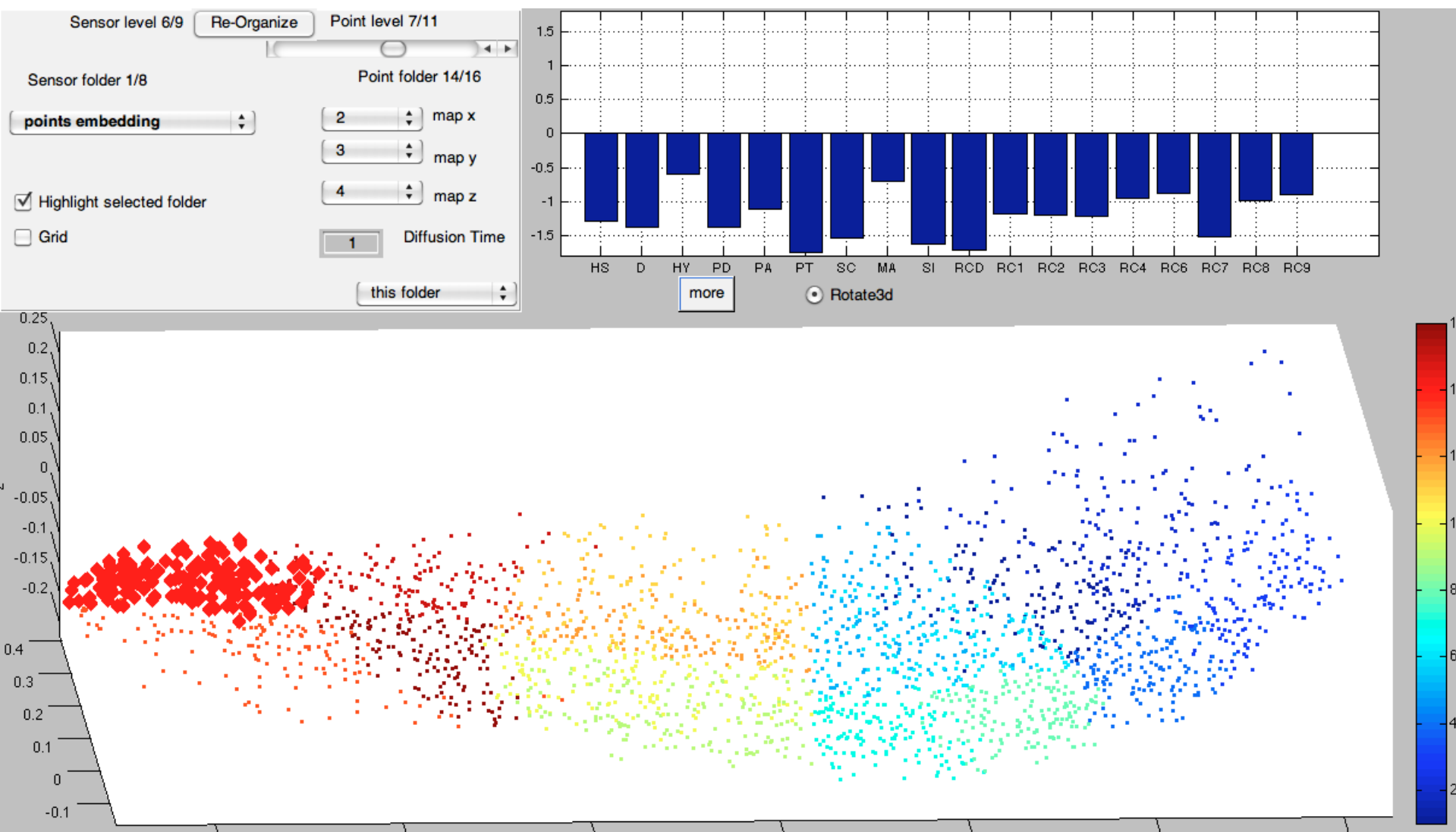**Contexts- clusters of "documents" with similar vocabulary profile**
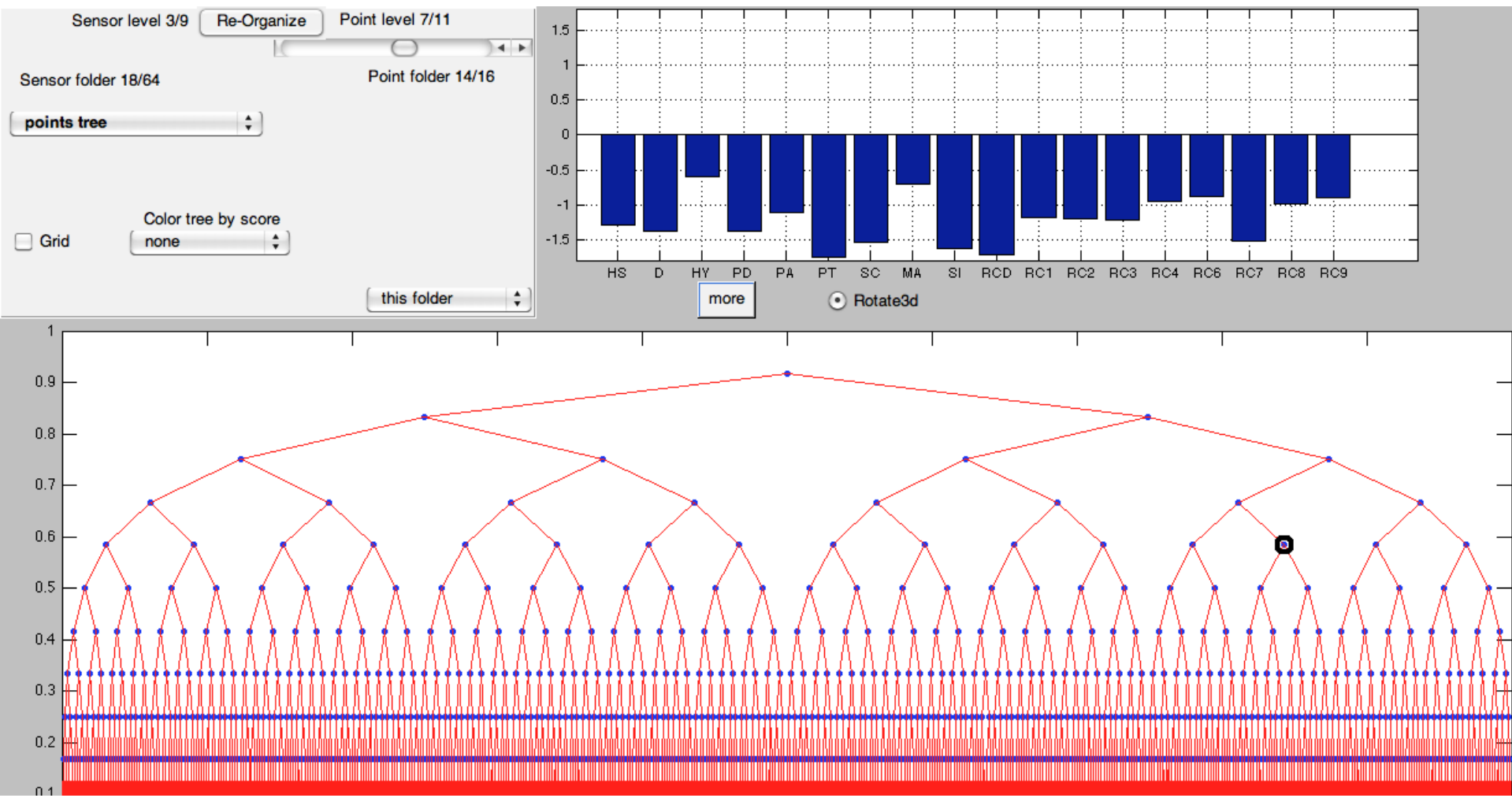
**Context tree**

Demographic organization by earth mover distance among profiles of the population.
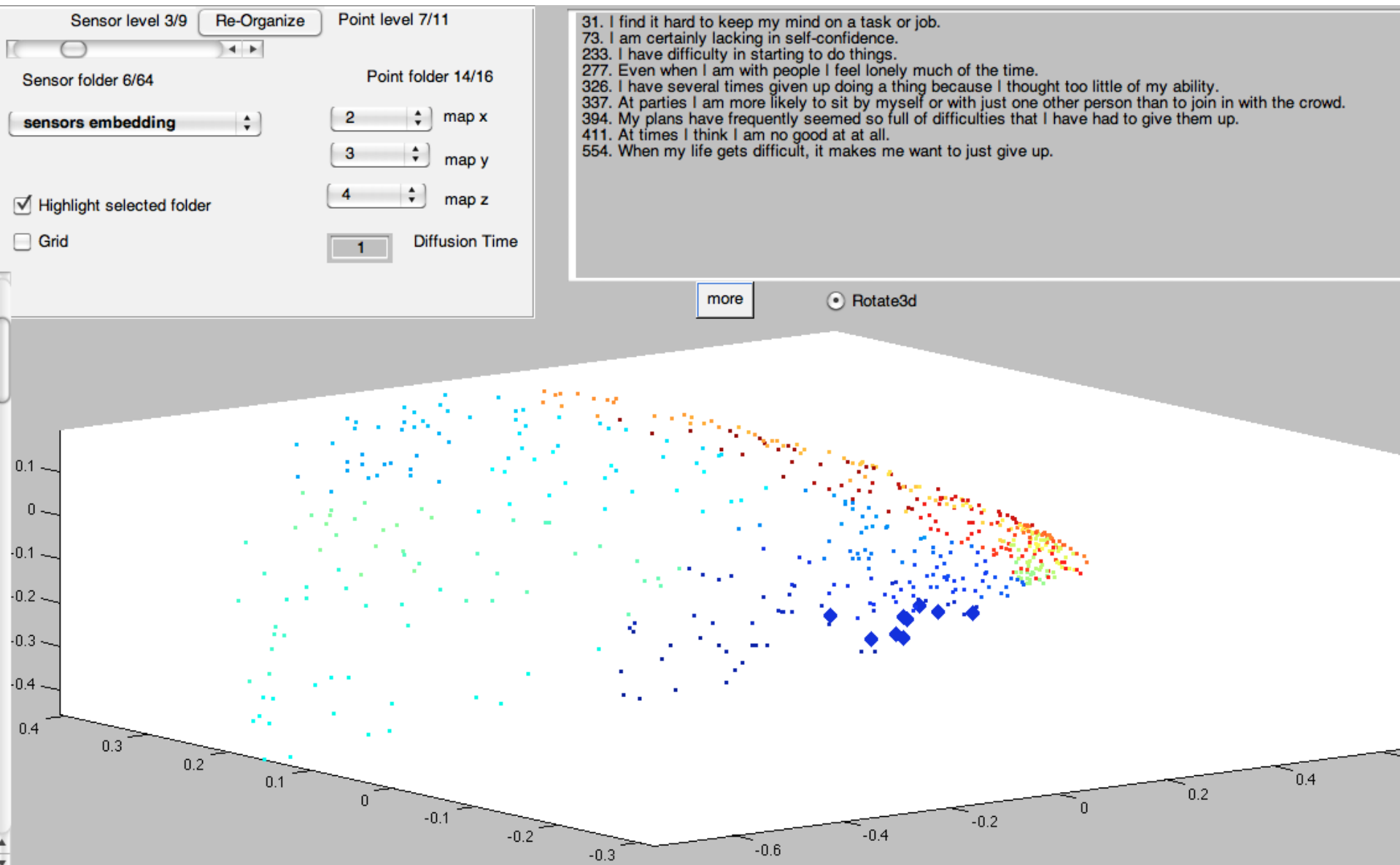The blue highlighted group is on one extremity ,having problems.

The red group is on the other end , being quite healthy .

The demographic tree , where the previous red group is marked.

# Conceptual organization of the questions into a geometry .

# Another group of questions



Sensor level 3/9 | Re-Organize | Point level 7/11

Sensor folder 18/64          Point folder 14/16

sensors embedding

☑ Highlight selected folder

☐ Grid

2    map x
3    map y
5    map z

1    Diffusion Time

47. I am almost never bothered by pains over my heart or in my chest.
57. I hardly ever feel pain in the back of my neck.
83. I have very few quarrels with members of my famlly.
91. I have little or no trouble with my muscles twitching or jumping.
255. I do not often notice my ears ringing or buzzing. I
295. I have never been paralyzed or had any unusual weakness of any of my muscles.
372. I am not easily angered.
427.  have never seen a vision.
564. I almost never lose self-control.

more        ⦿ Rotate3d

The same questions as above on the metaquestion tree , and the response profile
of various demographic groups , on the left problem groups , on the right healthy
people.

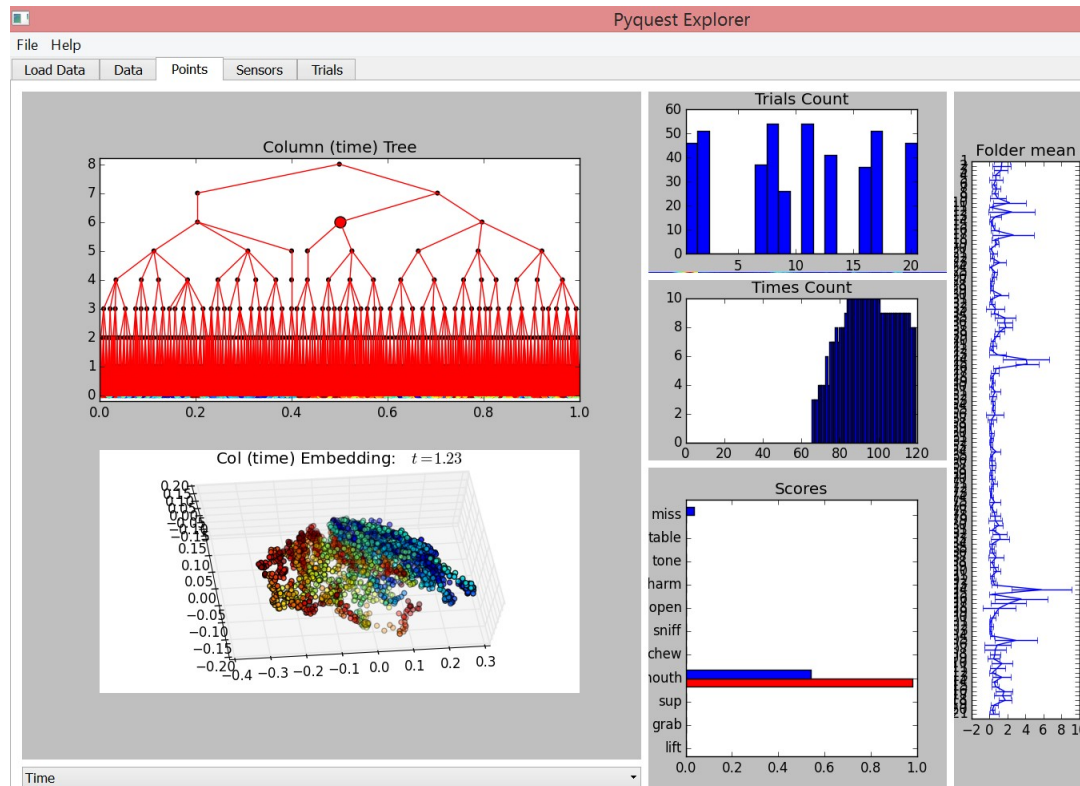The three dimensional plot above organizes time through various repetitive experiments to isolate the time where the sound trigger is activated for the mouse to reach for food .

## References

[1] E. Stein, Topics in Harmonic Analysis related to the Littlewood-Paley theory, Princeton University Press, 1970.

[2] R. Coifman and G. Weiss, Analyse Harmonique Noncommutative sur Certains Espaces Homogenes, Springer-Verlag, 1971.}

[3] R. Coifman ,G. Weiss, Extensions of Hardy spaces and their use in analysis. *Bul. Of the A.M.S.,* **83**, **#4**, 1977, 569-645.

[4] Belkin, M., & Niyogi, P. (2001). Laplacian eigenmaps and spectral techniques for embedding and clustering. Advances in Neural Information Processing Systems 14 (NIPS 2001) (p. 585).

[5]Belkin, M., & Niyogi, P. (2003a). Laplacian eigenmaps for dimensionality reduction and data representation. Neural Computation, 6, 1373{1396.

[6]Coifman, R. R., Lafon, S., Lee, A., Maggioni, M.,Nadler, B., Warner, F., & Zucker, S. (2005a) . Geometric diffusions as a tool for harmonic analysis and structure defnition of data. part i: Diffusion maps.Proc. of Nat. Acad. Sci., 7426{7431.

[7] Coifman R.R.,S Lafon, Diffusion maps, *Applied and Computational Harmonic Analysis,* 21: 5-30, 2006.

[8] Coifman R.R., B.Nadler, S Lafon, I G Kevrekidis, Diffusion maps, spectral clustering and reaction coordinates of dynamical systems, *Applied and Computational Harmonic Analysis,* 21:113-127, 2006.

[9] Ronald R Coifman[1], Mauro Maggioni[1], Steven W Zucker[1] and Ioannis G Kevrekidis "Geometric diffusions for the analysis of data from sensor networks" Current Opinion in Neurobiology 2005, 15:576–584

[10] Ham J, Lee DD, Mika S: Scholkopf: "A kernel view of the dimensionality reduction of manifolds". In Proceedings of the XXI Conference on Machine Learning, Banff, Canada, 2004

11.  R. Talmon and R. R. Coifman, "**Empirical intrinsic geometry for nonlinear modeling and time series filtering,**" Proc. Nat. Acad. Sci. (PNAS), vol. 110, no. 31, pp. 12535-12540, Jul. 2013.

12.  **Nonlinear intrinsic variables and state reconstruction in multiscale simulations**
Carmeline J. Dsilva, Ronen Talmon, Neta Rabin, Ronald R. Coifman, and Ioannis G. Kevrekidis . The Journal of Chemical Physics **139**, 184109 (2013); doi: 10.1063/1.4828457

13Talmon, R. & Coifman, R.R. **Intrinsic modeling of stochastic dynamical systems using empirical geometry.** Applied and Computational Harmonic Analysis 39, 138-160 (2015).

14. **Gavish, M. & Coifman, R.R. Sampling, denoising and compression of matrices by coherent matrix organization. Applied and Computational Harmonic Analysis 33, 354-369 (2012).**

15. **Coifman, R.R. & Gavish, M. Harmonic analysis of digital data bases. in Wavelets and Multiscale analysis 161-197 (Springer, 2011).**