

# Manifold Alignment of High-Dimensional Datasets

PIs: Sridhar Mahadevan, Rui Wang

Blake Foster, Armita Kaboli, Peter Krafft,  
Chang Wang (IBM Research)

Department of Computer Science  
University of Massachusetts, Amherst

# Outline

- Problem: transfer learning across domains
  - How to transfer knowledge across tasks?
- Solution:
  - Find manifold-based projections of original data
  - Align projected data in lower-dimensional space
- Applications:
  - Cross lingual IR, activity recognition,  
reinforcement learning
- Acceleration using GPU

# Example: Cross-Lingual Retrieval



Madam President, on a point of order. You will be aware from the press and television that there have been a number of bomb explosions and killings in Sri Lanka.



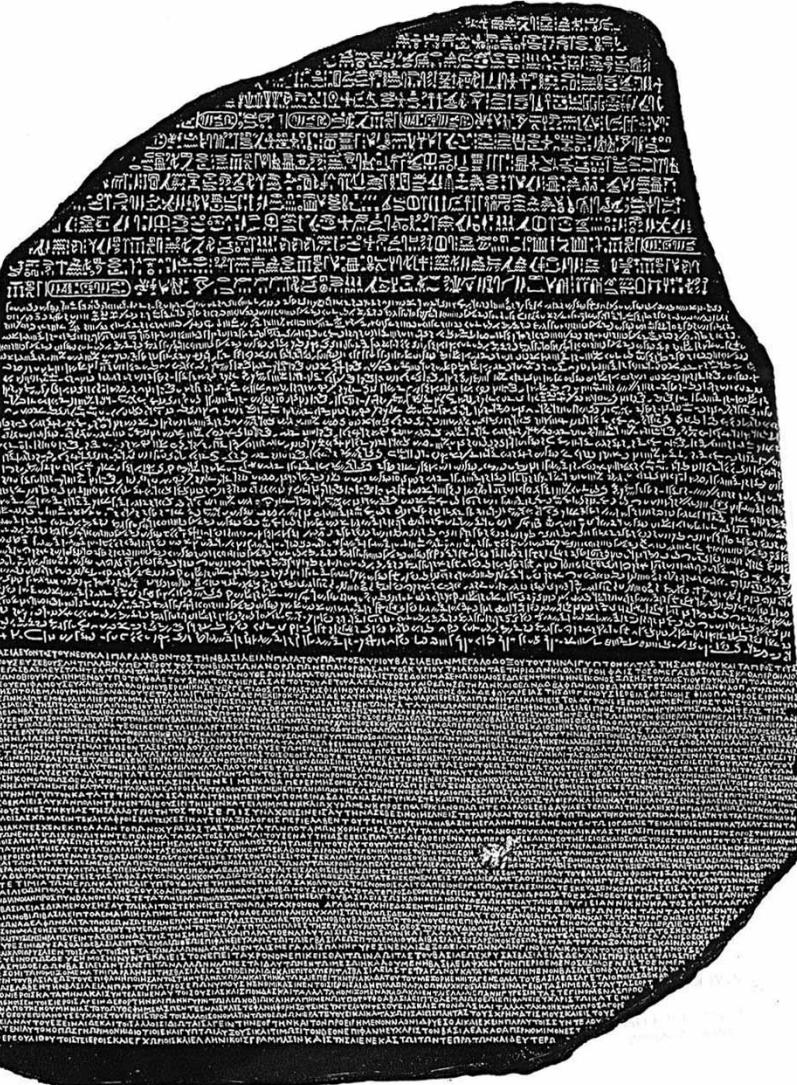
Signora Presidente, intervengo per una mozione d'ordine. Come avrà letto sui giornali o sentito alla televisione, in Sri Lanka si sono verificati numerosi assassinii ed esplosioni di ordigni.



Frau Präsidentin, zur Geschäftsordnung. Wie Sie sicher aus der Presse und dem Fernsehen wissen, gab es in Sri Lanka mehrere Bombenexplosionen mit zahlreichen Toten.



# Rosetta Stone



Carved in 196 BC

Discovered in  
1799 AD

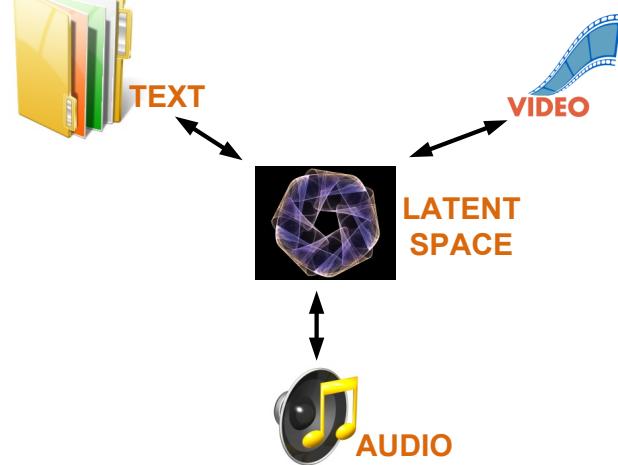
Paved the way to  
decipherment of  
hieroglyphics

Egyptian  
hieroglyphics

Demotic

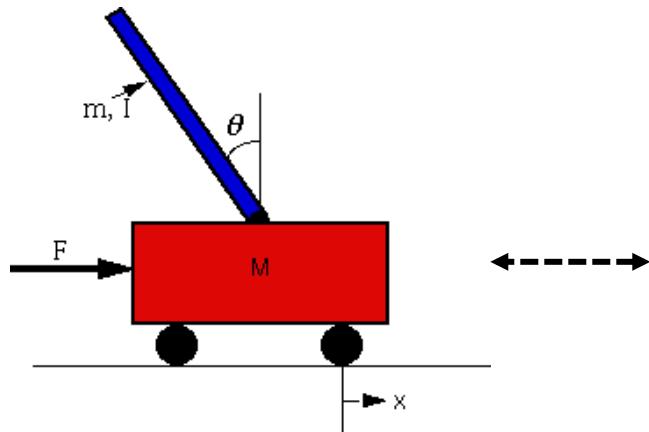
Greek

# Knowledge Transfer



[flowers](#), [grass](#), [tiger](#), [water](#)

# Knowledge Transfer in Reinforcement Learning



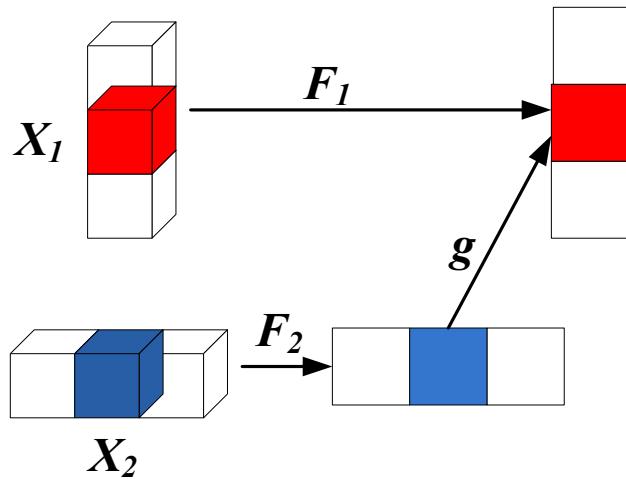
# Why is Transfer Difficult

- Data is usually **high-dimensional** ( $\sim 100,000$  dimensions or more)
- The source and target domains may have **distinct features**
- Correspondences in the original space are difficult to find due to the high dimensionality

# Classical Solutions to Alignment

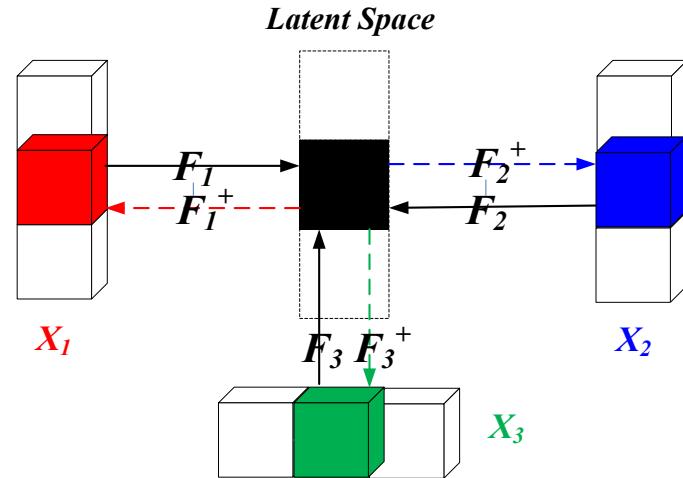
- Canonical Correlational Analysis (Hotelling)
  - Analogous to PCA across two data sets
  - Find a linear projection of each domain that maximizes correlation across domains
- Dynamic Time Warping (Juang)
  - Used in aligning time series data
  - Based on dynamic programming
- Our approach: find alignments by explicitly modeling the underlying data manifold

# Manifold Alignment



## Two-step alignment

*Example: Procrustes alignment  
(Wang and Mahadevan, ICML 2008)*



## One-step alignment

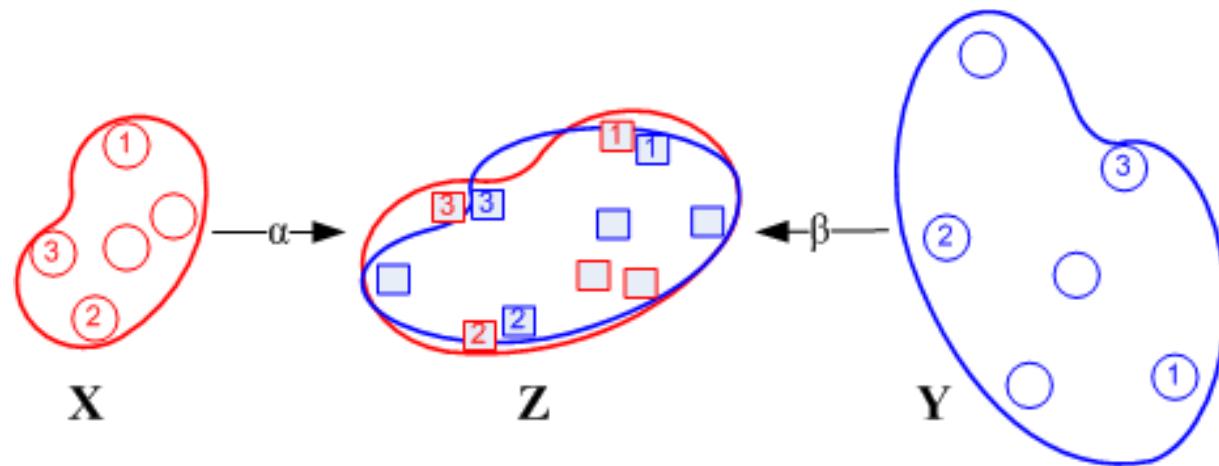
*Example: Manifold Projections  
(Wang and Mahadevan, IJCAI 2009)*

# A Spectrum of Manifold Alignment Approaches

	<i>Given correspondences</i>	<i>Given labels</i>	<i>Unsupervised alignment</i>
<b>Preserve Local geometry</b>			
<b>Preserve Global geometry</b>			
<b>One-step alignment</b>			
<b>Two-step alignment</b>			
<b>Feature-level</b>			
<b>Instance-level</b>			

*Procrustes alignment*    *Manifold Projections (MP)*    *Extensions of MP*

# Feature-level Manifold Projection

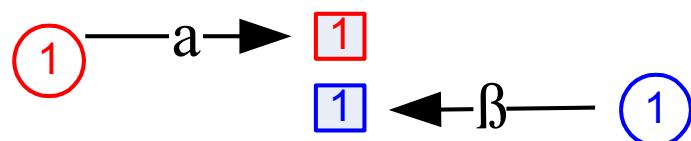


$$X = [x_1, \dots, x_m], x_i \in R^p.$$
$$Y = [y_1, \dots, y_n], y_j \in R^q$$
$$x_i \Leftrightarrow y_i \text{ for } i \in [1, l]$$

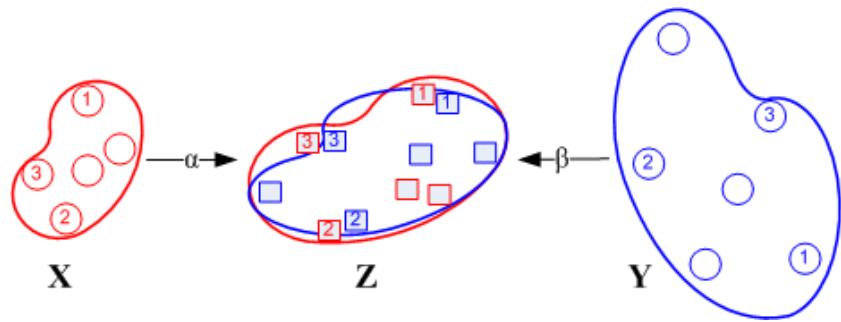
Projection Results:

$$x_i \rightarrow \alpha^T x_i$$

$$y_j \rightarrow \beta^T y_j$$



# Manifold Projection and Alignment

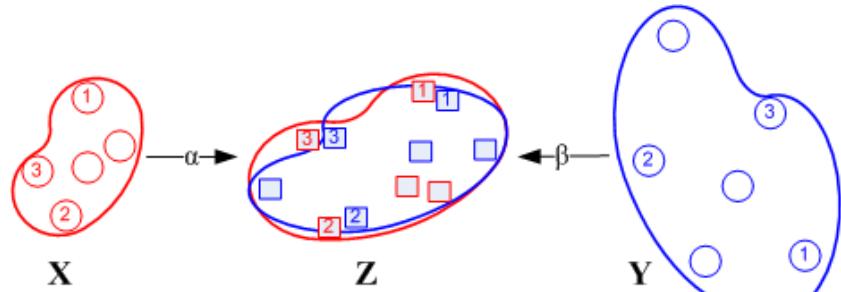


$$X = [x_1, \dots, x_m], x_i \in R^p.$$

$$Y = [y_1, \dots, y_n], y_j \in R^q$$

$$x_i \Leftrightarrow y_i \text{ for } i \in [1, l]$$

# Manifold Projection and Alignment



$$X = [x_1, \dots, x_m], x_i \in R^p.$$

$$Y = [y_1, \dots, y_n], y_j \in R^q$$

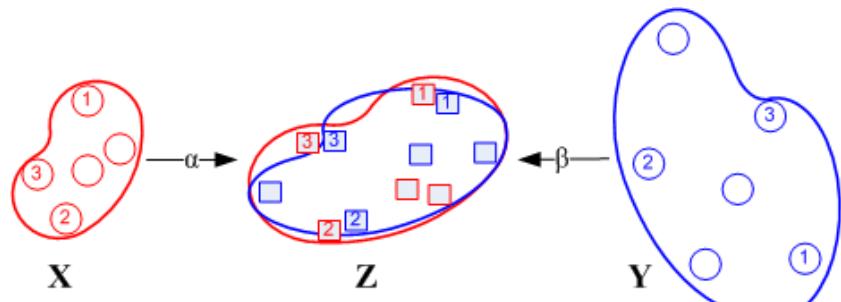
$$x_i \Leftrightarrow y_i \text{ for } i \in [1, l]$$

We want to find mapping functions  $\alpha, \beta$  to minimize the cost function  $C(\alpha, \beta)$ , where

$$C(\alpha, \beta) = \mu \sum_i \sum_j (\alpha^T x_i - \beta^T y_j)^2 W^{i,j} + 0.5 \sum_{i,j} (\alpha^T x_i - \alpha^T x_j)^2 W_x^{i,j} + 0.5 \sum_{i,j} (\beta^T y_i - \beta^T y_j)^2 W_y^{i,j}$$

---

# Manifold Projection and Alignment



$$X = [x_1, \dots, x_m], x_i \in R^p.$$

$$Y = [y_1, \dots, y_n], y_j \in R^q$$

$$x_i \Leftrightarrow y_i \text{ for } i \in [1, l]$$

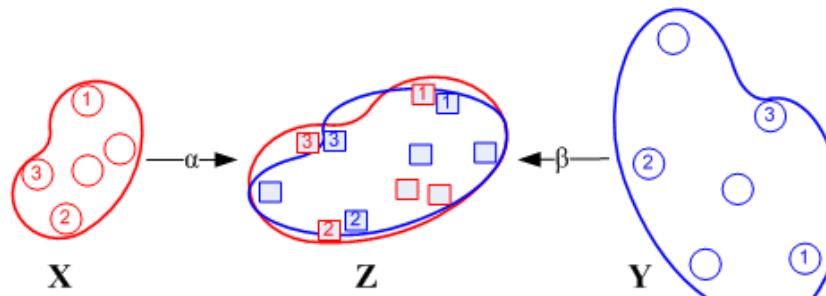
We want to find mapping functions  $\alpha, \beta$  to minimize the cost function  $C(\alpha, \beta)$ , where

$$C(\alpha, \beta) = \mu \sum_i \sum_j (\alpha^T x_i - \beta^T y_j)^2 W^{i,j} + 0.5 \sum_{i,j} (\alpha^T x_i - \alpha^T x_j)^2 W_x^{i,j} + 0.5 \sum_{i,j} (\beta^T y_i - \beta^T y_j)^2 W_y^{i,j}$$

The **first** term encourages the corresponding instances from different domains to be projected to similar locations.

$W^{i,j}=1$ , when  $x_i$  and  $y_j$  are in correspondence; 0, otherwise.

# Manifold Projection and Alignment



$$X = [x_1, \dots, x_m], x_i \in R^p.$$

$$Y = [y_1, \dots, y_n], y_j \in R^q$$

$$x_i \Leftrightarrow y_i \text{ for } i \in [1, l]$$

We want to find mapping functions  $\alpha, \beta$  to minimize the cost function  $C(\alpha, \beta)$ , where

$$C(\alpha, \beta) = \boxed{\mu \sum_i \sum_j (\alpha^T x_i - \beta^T y_j)^2 W^{i,j}} + 0.5 \sum_{i,j} (\alpha^T x_i - \alpha^T x_j)^2 W_x^{i,j} + 0.5 \sum_{i,j} (\beta^T y_i - \beta^T y_j)^2 W_y^{i,j}$$

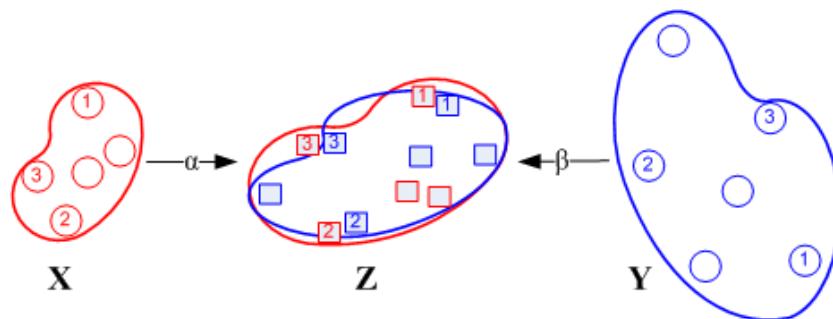
The **first** term encourages the corresponding instances from different domains to be projected to similar locations.

$W^{i,j}=1$ , when  $x_i$  and  $y_j$  are in correspondence; 0, otherwise.

$$W = \begin{bmatrix} 1 & & & \\ 1 & \cdots & & \\ \cdots & & 1 & 0 \\ & & 0 & \cdots \\ & & \cdots & 0 \end{bmatrix}$$

- **When 1:1 correspondence is given** ( $x_i \leftrightarrow y_i$  for  $i \leq l$ ):
- **When many:many correspondence is given**, set corresponding entries to 1.
- **When nothing is given**, we can use local geometry information to fill in this matrix. (IJCAI 2009)

# Manifold Projection and Alignment



$$X = [x_1, \dots, x_m], x_i \in R^p.$$

$$Y = [y_1, \dots, y_n], y_j \in R^q$$

$$x_i \Leftrightarrow y_i \text{ for } i \in [1, l]$$

We want to find mapping functions  $\alpha, \beta$  to minimize the cost function  $C(\alpha, \beta)$ , where

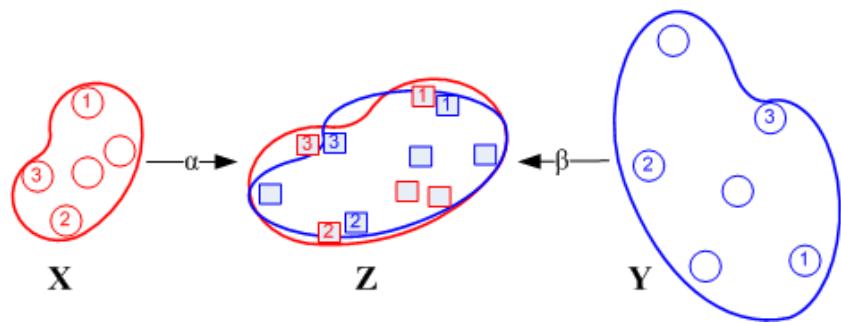
$$C(\alpha, \beta) = \mu \sum_i \sum_j (\alpha^T x_i - \beta^T y_j)^2 W^{i,j} + 0.5 \sum_{i,j} (\alpha^T x_i - \alpha^T x_j)^2 W_x^{i,j} + 0.5 \sum_{i,j} (\beta^T y_i - \beta^T y_j)^2 W_y^{i,j}$$

The second and third terms encourage the neighborhood relationship within  $X$  and  $Y$  to be preserved in the mappings.

$W_x^{i,j}$ : similarity of  $x_i$  and  $x_j$  in the original space.

$W_y^{i,j}$ : similarity of  $y_i$  and  $y_j$  in the original space.

# Comparison with Canonical Correlation Analysis (CCA)



$$X = [x_1, \dots, x_m], x_i \in R^p.$$

$$Y = [y_1, \dots, y_n], y_j \in R^q$$

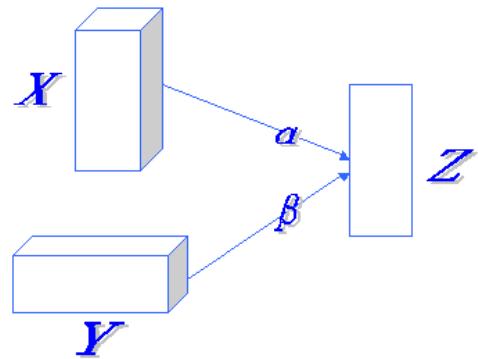
$$x_i \Leftrightarrow y_i \text{ for } i \in [1, l]$$

We want to find mapping functions  $\alpha, \beta$  to minimize the cost function  $C(\alpha, \beta)$ , where

$$C(\alpha, \beta) = \mu \sum_i \sum_j (\alpha^T x_i - \beta^T y_j)^2 \cancel{W^{i,j}} + 0.5 \sum_{i,j} (\alpha^T x_i - \alpha^T x_j)^2 \cancel{W_x^{i,j}} + 0.5 \sum_{i,j} (\beta^T y_i - \beta^T y_j)^2 \cancel{W_y^{i,j}}$$

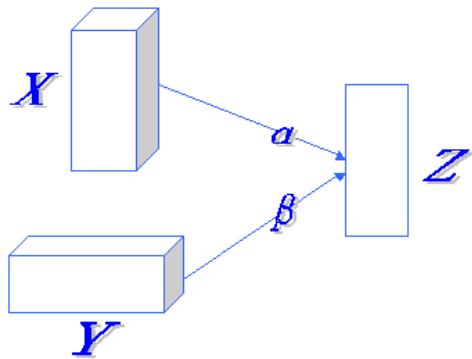
# How to compute projections?

**Optimal Solution:**



# How to compute projections?

## ***Optimal Solution:***



- (1) Construct  $Z, L, D$  using  $X, Y$  and  $W$  (the correspondences).

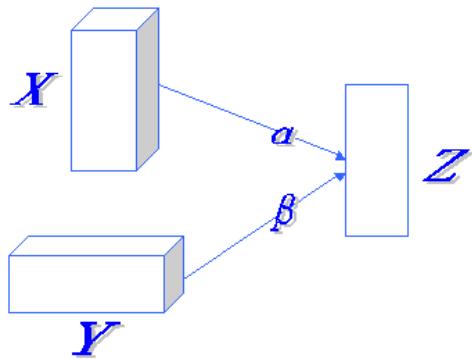
$D_x$  is a diagonal matrix:  $D_x^{ii} = \sum_j W_x^{ij}$ .  
 $L_x = D_x - W_x$ .  
 $D_y$  is a diagonal matrix:  $D_y^{ii} = \sum_j W_y^{ij}$ .  
 $L_y = D_y - W_y$ .  
 $\Omega_1$  is an  $m \times m$  diagonal matrix, and  $\Omega_1^{ii} = \sum_j W^{i,j}$ .  
 $\Omega_2$  is an  $m \times n$  matrix, and  $\Omega_2^{i,j} = W^{i,j}$ .  
 $\Omega_3$  is an  $n \times m$  matrix, and  $\Omega_3^{j,i} = W^{j,i}$ .  
 $\Omega_4$  is an  $n \times n$  diagonal matrix, and  $\Omega_4^{ii} = \sum_j W^{j,i}$ .

$$Z = \begin{pmatrix} X & 0 \\ 0 & Y \end{pmatrix}.$$
$$D = \begin{pmatrix} D_x & 0 \\ 0 & D_y \end{pmatrix}.$$
$$L = \begin{pmatrix} L_x + \mu\Omega_1 & -\mu\Omega_2 \\ -\mu\Omega_3 & L_y + \mu\Omega_4 \end{pmatrix}.$$

Create a joint domain.  
( use correspondences  
if available)

# How to compute projections

## ***Optimal Solution:***



- (1) Construct  $Z, L, D$  using  $X, Y$  and  $W$  (the correspondences).

$D_x$  is a diagonal matrix:  $D_x^{ii} = \sum_j W_x^{ij}$ .  
 $L_x = D_x - W_x$ .  
 $D_y$  is a diagonal matrix:  $D_y^{ii} = \sum_j W_y^{ij}$ .  
 $L_y = D_y - W_y$ .  
 $\Omega_1$  is an  $m \times m$  diagonal matrix, and  $\Omega_1^{ii} = \sum_j W^{i,j}$ .  
 $\Omega_2$  is an  $m \times n$  matrix, and  $\Omega_2^{i,j} = W^{i,j}$ .  
 $\Omega_3$  is an  $n \times m$  matrix, and  $\Omega_3^{j,i} = W^{j,i}$ .  
 $\Omega_4$  is an  $n \times n$  diagonal matrix, and  $\Omega_4^{ii} = \sum_j W^{j,i}$ .

$$Z = \begin{pmatrix} X & 0 \\ 0 & Y \end{pmatrix}.$$
$$D = \begin{pmatrix} D_x & 0 \\ 0 & D_y \end{pmatrix}.$$
$$L = \begin{pmatrix} L_x + \mu\Omega_1 & -\mu\Omega_2 \\ -\mu\Omega_3 & L_y + \mu\Omega_4 \end{pmatrix}.$$

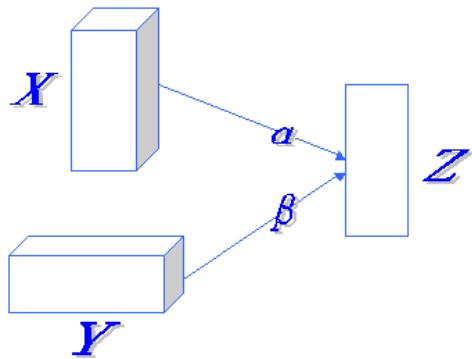
Create a joint domain.  
( use correspondences  
if available)

Project the joint domain to  
a lower dimensional space.

- (2) *Theorem 1 :  $\alpha, \beta$  to minimize  $C(\alpha, \beta)$  are given by the eigenvectors corresponding to the smallest eigenvalues of  $ZLZ^T \gamma = \lambda ZDZ^T \gamma$ .*

# How to compute projections?

## **Optimal Solution:**



- (1) Construct  $Z, L, D$  using  $X, Y$  and  $W$  (the correspondences).

$D_x$  is a diagonal matrix:  $D_x^{ii} = \sum_j W_x^{ij}$ .  
 $L_x = D_x - W_x$ .  
 $D_y$  is a diagonal matrix:  $D_y^{ii} = \sum_j W_y^{ij}$ .  
 $L_y = D_y - W_y$ .  
 $\Omega_1$  is an  $m \times m$  diagonal matrix, and  $\Omega_1^{ii} = \sum_j W^{i,j}$ .  
 $\Omega_2$  is an  $m \times n$  matrix, and  $\Omega_2^{i,j} = W^{i,j}$ .  
 $\Omega_3$  is an  $n \times m$  matrix, and  $\Omega_3^{j,i} = W^{j,i}$ .  
 $\Omega_4$  is an  $n \times n$  diagonal matrix, and  $\Omega_4^{ii} = \sum_j W^{j,i}$ .

$$Z = \begin{pmatrix} X & 0 \\ 0 & Y \end{pmatrix}.$$
$$D = \begin{pmatrix} D_x & 0 \\ 0 & D_y \end{pmatrix}.$$
$$L = \begin{pmatrix} L_x + \mu\Omega_1 & -\mu\Omega_2 \\ -\mu\Omega_3 & L_y + \mu\Omega_4 \end{pmatrix}.$$

Create a joint domain.  
( use correspondences  
if available)

Project the joint domain to  
a lower dimensional space.

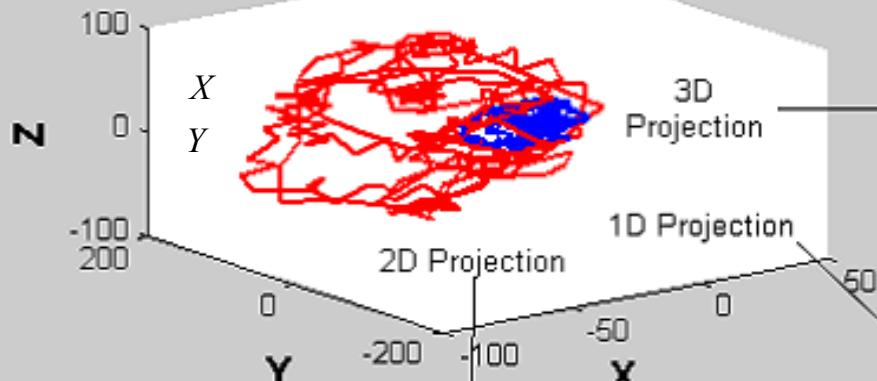
- (2) *Theorem 1 :  $\alpha, \beta$  to minimize  $C(\alpha, \beta)$  are given by the eigenvectors corresponding to the smallest eigenvalues of*

$$ZLZ^T \gamma = \lambda ZDZ^T \gamma.$$

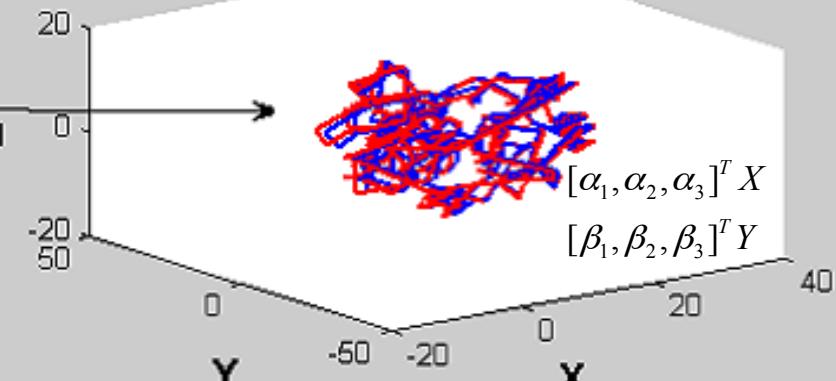
- (3)  $\begin{bmatrix} \alpha \\ \beta \end{bmatrix} = [\gamma_1, \dots, \gamma_d]$ , where  $\gamma_i$  is the  $i^{th}$  minimum eigenvector.

# Protein Alignment

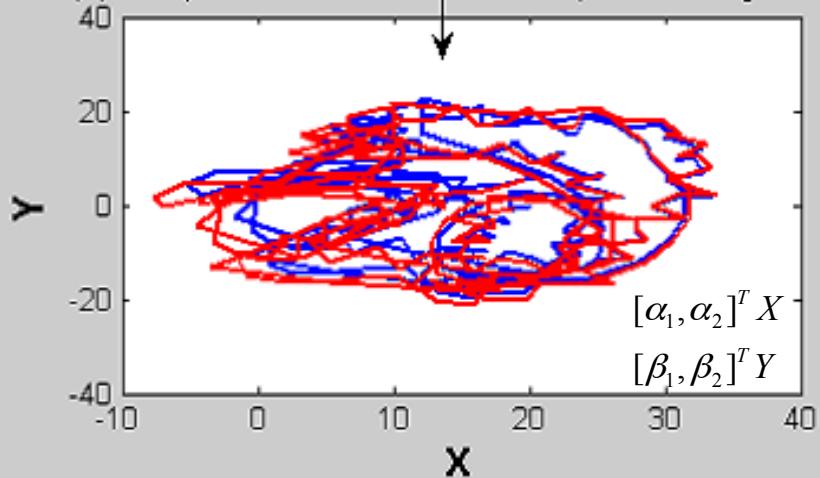
(A) Comparison of Manifold A and B (Before Alignment)



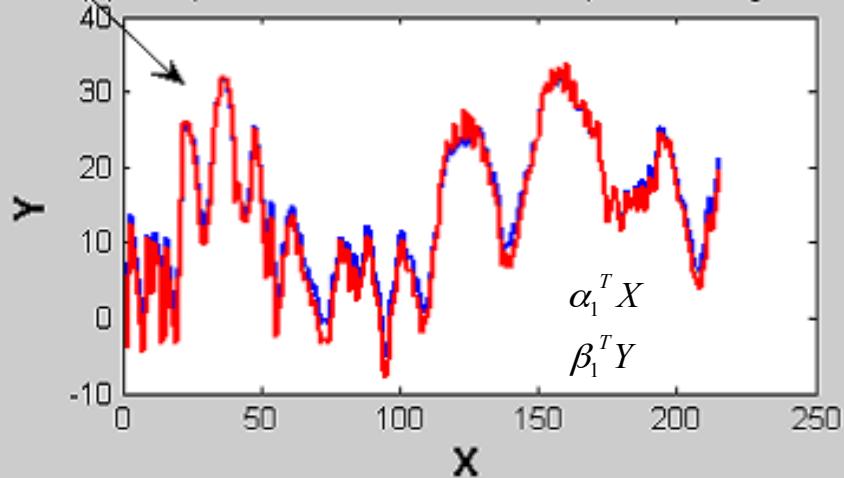
(B) Comparison of Manifold A and B (After 3D Alignment)



(C) Comparison of Manifold A and B (After 2D Alignment)



(D) Comparison of Manifold A and B (After 1D Alignment)



# General Framework for *Manifold Projection*

- Instance-level

$$C(\mathcal{Y}_1, \dots, \mathcal{Y}_c) = 0.5\mu_1 \sum_{a=1}^c \sum_{b=1}^c \sum_{i=1}^{m_a} \sum_{j=1}^{m_b} \|\mathcal{Y}_a^i - \mathcal{Y}_b^j\|^2 W_{a,b}^{i,j} + 0.5\mu_2 \sum_{k=1}^c \sum_{i=1}^{m_k} \sum_{j=1}^{m_k} \|\mathcal{Y}_k^i - \mathcal{Y}_k^j\|^2 W_k^{i,j}$$

1.  $c=1$ : Laplacian eigenmaps (Belkin, Niyogi, 2003)
2.  $c=2, \mu_2=1, W_{a,b}$  is an identity matrix: Semi-supervised alignment (Ham, et. al. 2005)
3.  $c=2, \mu_2=0, W_{a,b}$  is an identity matrix: Non-linear Canonical Correlation Analysis
4. .....

- Feature-level

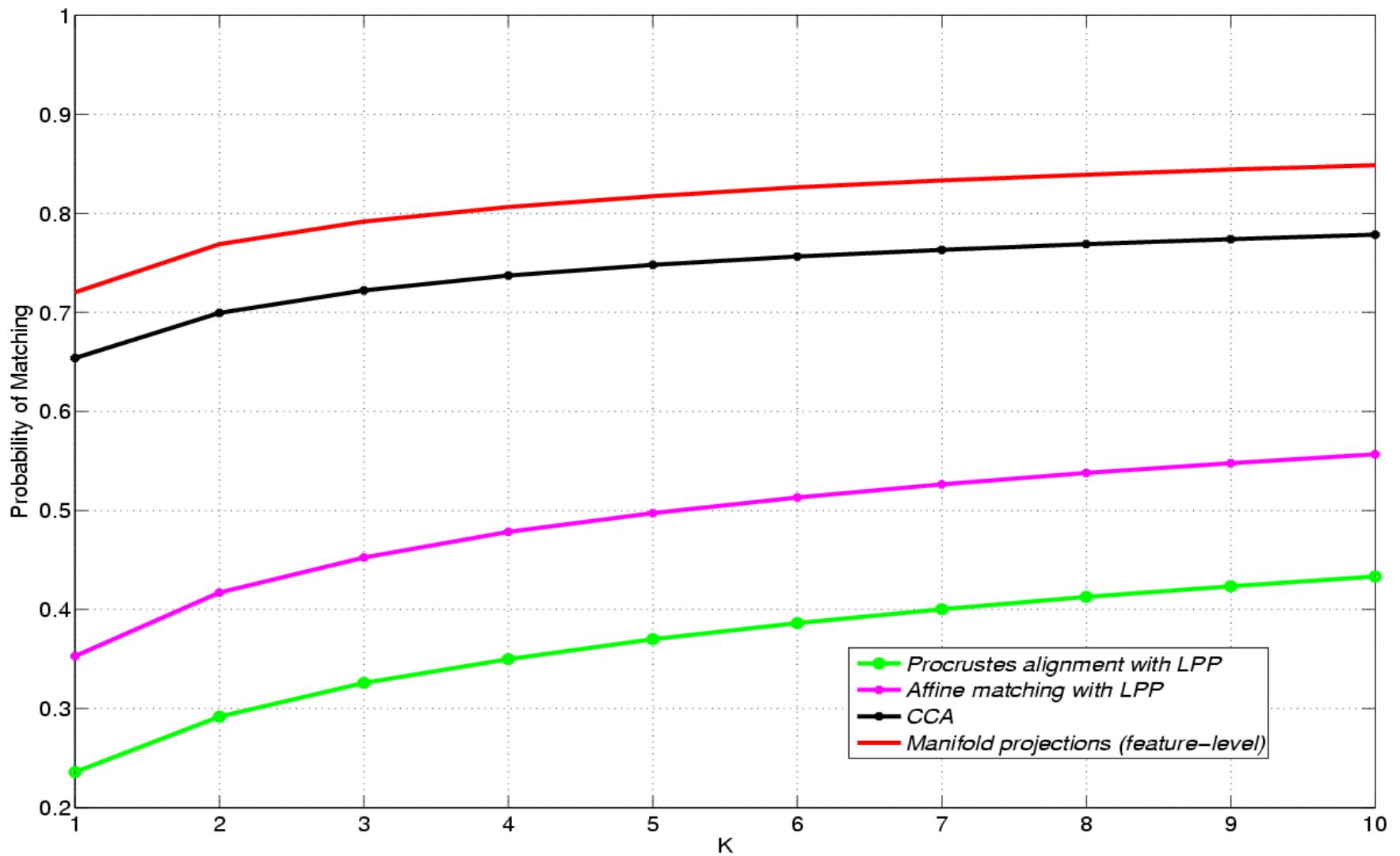
$$C(\mathcal{F}_1, \dots, \mathcal{F}_c) = 0.5\mu_1 \sum_{a=1}^c \sum_{b=1}^c \sum_{i=1}^{m_a} \sum_{j=1}^{m_b} \|\mathcal{F}_a^T x_a^i - \mathcal{F}_b^T x_b^j\|^2 W_{a,b}^{i,j} + 0.5\mu_2 \sum_{k=1}^c \sum_{i=1}^{m_k} \sum_{j=1}^{m_k} \|\mathcal{F}_k^T x_k^i - \mathcal{F}_k^T x_k^j\|^2 W_k^{i,j}$$

2.  $c=2, \mu_2=0, W_{a,b}$  is an identity matrix: Canonical Correlation Analysis (Hotelling, 1936)
3. *Unsupervised manifold alignment* (Wang, Mahadevan, IJCAI 09)
4.  $c=2, \mu_2=1, W_{a,b}$  is a symmetric matrix
5.  $c>2$ : Multiple manifold alignment
6. .....

# European Parliament Parallel Corpus

- **Data:**
  - **70,458** English-Italian-German document triples.
  - Those documents have **> 170,000,000** words.
- **Features:**
  - $X_1$ : English documents are represented by the most commonly used **2,500 English words**.
  - $X_2$ : Italian documents are represented by the most commonly used **2,500 Italian words**.
  - $X_3$ : German documents are represented by the most commonly used **2,500 German words**.

# EU Parallel Corpus: Results



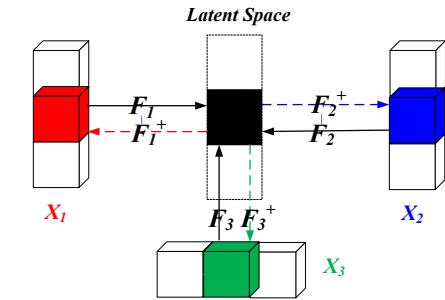
# EU Parallel Corpus

- Interpretation of mapping functions:

$F1=[\alpha_1 \alpha_2 \dots \alpha_{200}]$ : 2500\*200 matrix

$F2=[\beta_1 \beta_2 \dots \beta_{200}]$ : 2500\*200 matrix

$F3=[\gamma_1 \gamma_2 \dots \gamma_{200}]$ : 2500\*200 matrix



Top Terms	
$\alpha_1$	fisheries fishing agency fishermen negotiated applause protocol nos sustainability ports
$\alpha_2$	consumers internet consumer strategies b5 bulgaria behaviour b4 discharge november
$\alpha_3$	strategies swedish courage denmark telecommunications nato credibility wine regional brings
$\alpha_4$	interinstitutional parliaments repeated guarantees century rail finland british choose conciliation
$\alpha_5$	unemployment thursday heads portuguese economies declarations balkans widespread islands india

Top Terms	
$\beta_1$	pesca agenzia a5 applausi protocollo ripartizione pescatori bilaterali sostenibilita tonnellate
$\beta_2$	consumatori reca internet consumatore strategie discarico bulgaria novembre allargamento chiusa
$\beta_3$	strategie svedese interrogazioni occidentali danimarca regionale kyoto coraggio credibilita segretario
$\beta_4$	parlamenti sentenza interistituzionale aprile ferroviario britannica tecnici essenzialmente unanimita indipendenza
$\beta_5$	disoccupazione giovedi scientifica portoghese balcani aeree firmato turco maggio piccoli

Top Terms	
$\gamma_1$	fischerei fischereipolitik agentur protokolls fischer protokoll ablehnen a5 tatschlichen arten
$\gamma_2$	verbrauchern verbraucher strategien internet bulgarien entlastung anfragen todesstrafe b4 zukunftigen
$\gamma_3$	schwedischen strategien regionalpolitik frist anfragen westlichen mut nato regionalen minuten
$\gamma_4$	parlamente interinstitutionelle b4 parlamenten urteil spezielle folgt anmerkungen nächster beschluß
$\gamma_5$	portugiesischen personal arbeitslosigkeit offenen wissenschaftlichen polizei verordnungen donnerstag inseln äußerungen

# The Indus Script (2500 BC)



- **Data:**

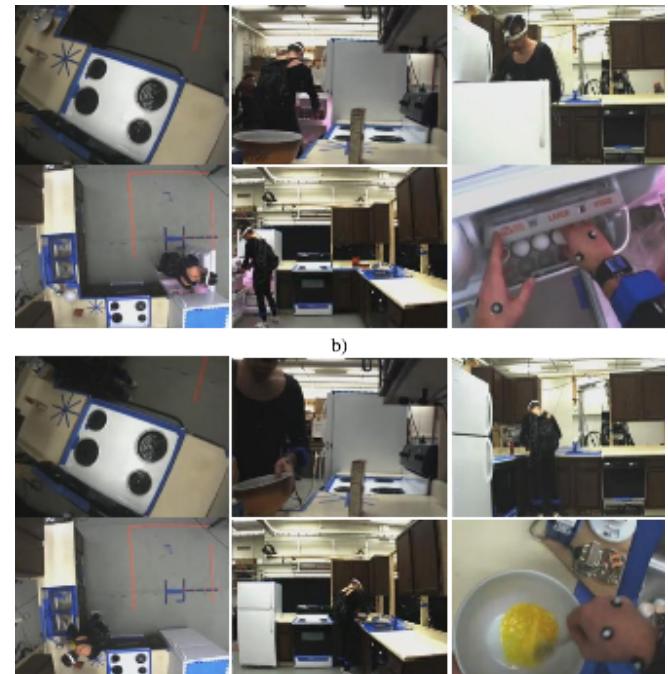
- **1548** “documents” – different seals
- **7000** “words” – symbols on the seals (**377** unique)
- Many competing inconsistent claims of decipherement
- Is this a language (Rao et al., Science, PNAS)

- **Challenges:**

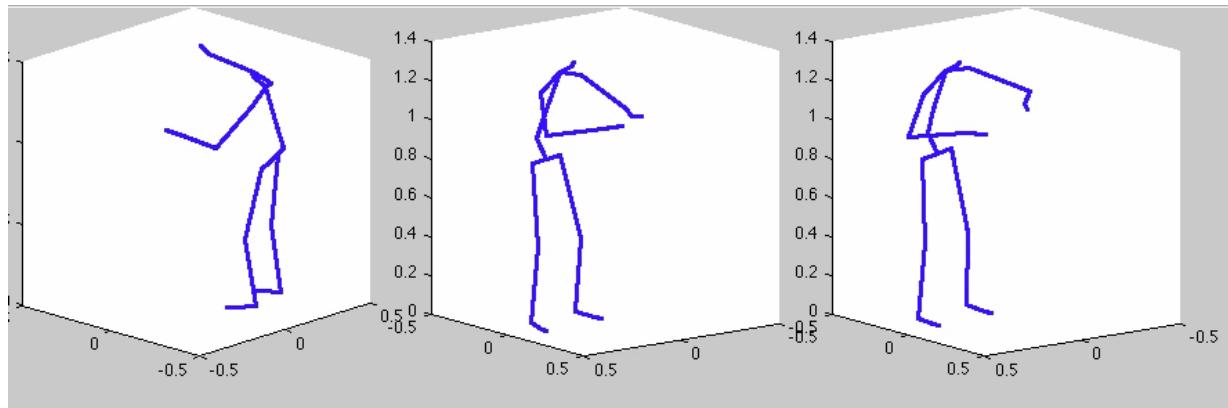
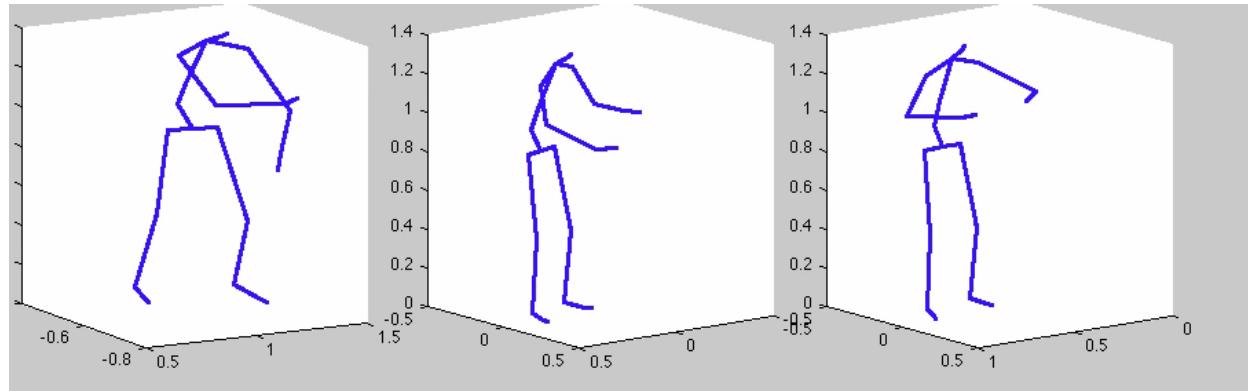
- We have no parallel corpus, no Rosetta Stone
- The document-symbol matrix is sparse
- Many of the documents are very short

# Manifold Alignment over time

- CMU Multimodal activity dataset
  - [kitchen.cs.cmu.edu](http://kitchen.cs.cmu.edu)
- Measure human activity while cooking
  - 26 subjects
  - 5 different recipes
- Many sensors:
  - Cameras
  - RFID
  - Audio
  - Joint angles



# Temporal Alignment of Activities



# Scaling to **Really Large** Problems

- Data mining a million books!
  - NSF Large IIS Project: James Allan, PI, U.Mass
- Transfer learning in Reinforcement Learning
  - Backgammon ( $10^{20}$  states!), Go ( $10^{170}$  states!!)
- CMU Multi-modal activity dataset
  - Many sensor modalities, high frame rate