

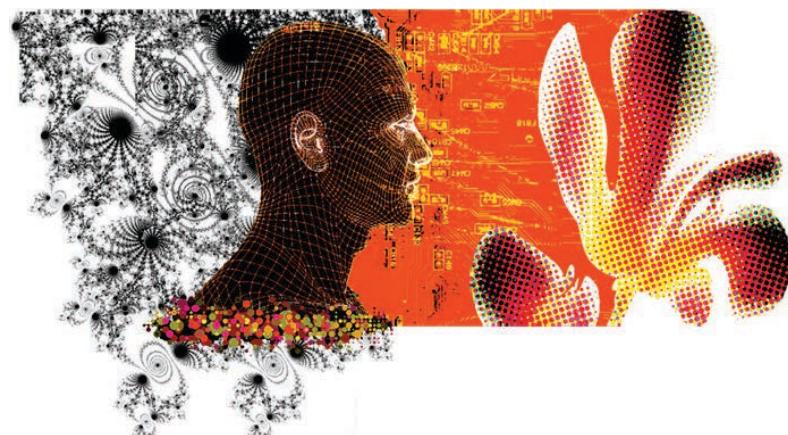
# Data Science Training

November 2017

## Data Visualization

Xavier Bresson

Data Science and AI Research Centre  
NTU, Singapore



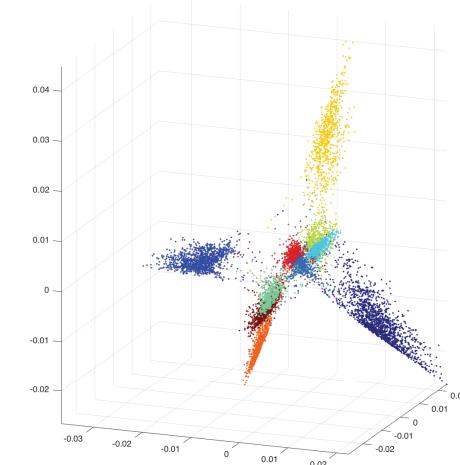
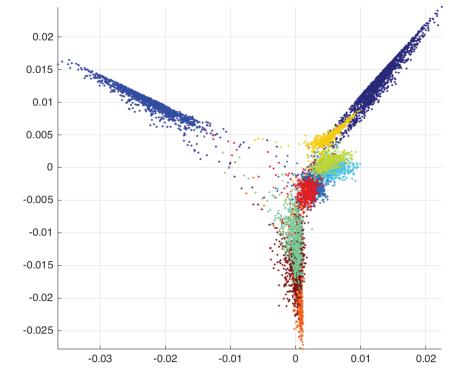
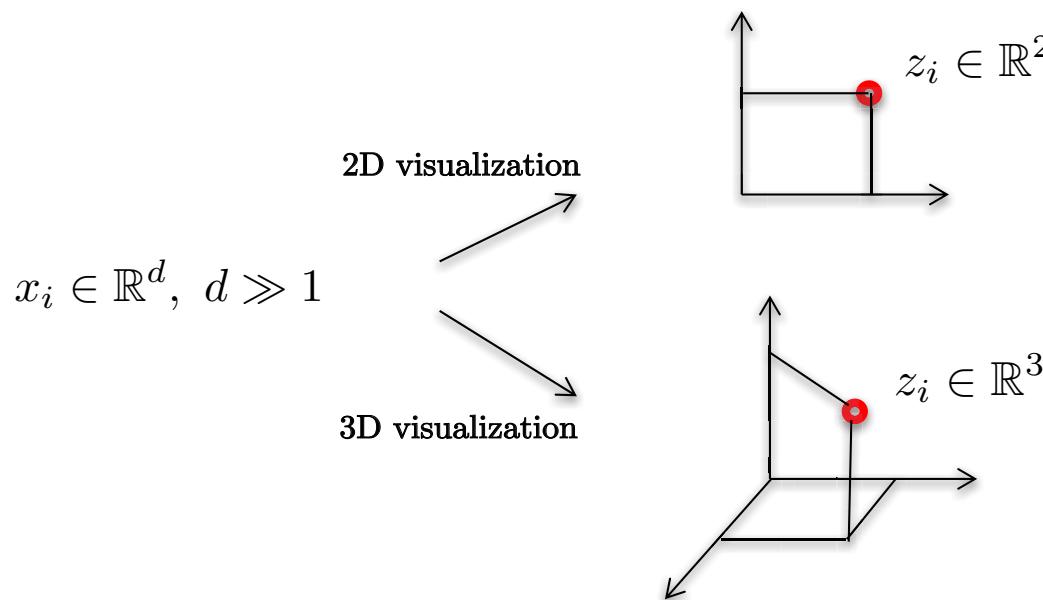
<http://data-science-optum17.tk>

# Outline

- Visualization Problem
- Kernel PCA
- Locally-Linear Embedding (LLE)
- Laplacian Eigenmaps
- T-SNE
- Conclusion

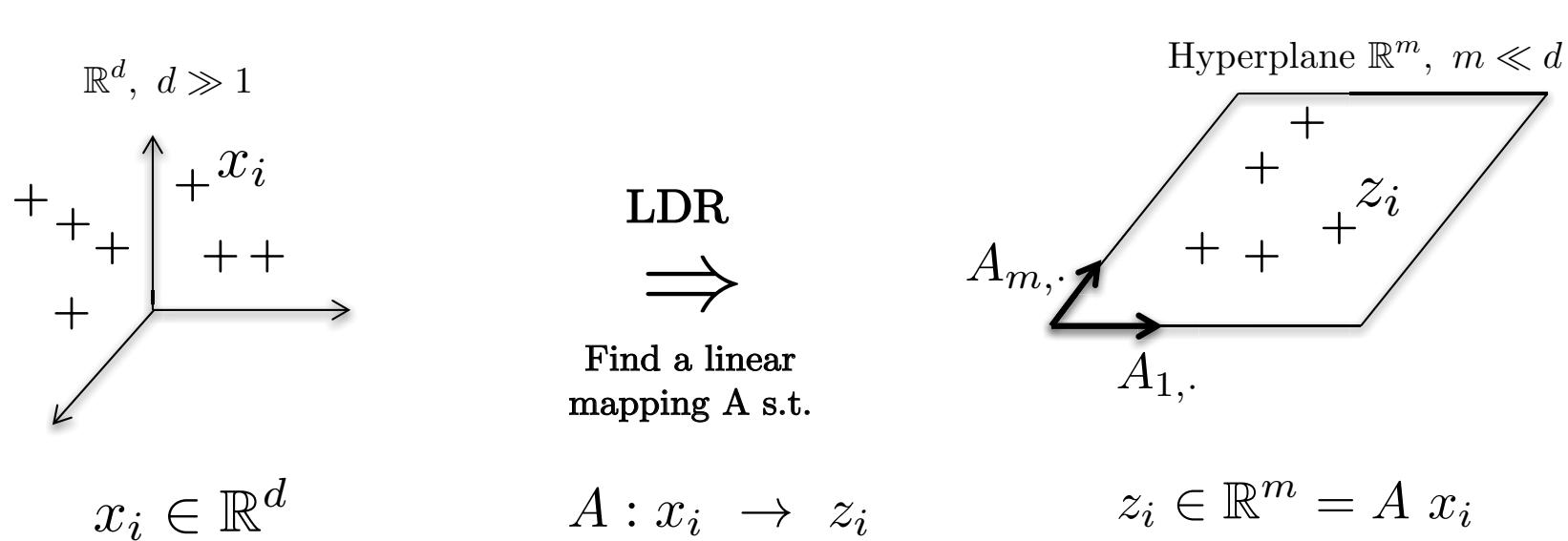
# Visualization Problem

- Data **visualization** is the **same** problem as
  - (1) *Data representation*
  - (2) *Feature extraction*
- **Data representation** looks for the best "dictionary"  $D$  where the data  $x$  can be represented, and the projected data  $z$  on  $D$  are used as coordinates for 2D or 3D visualization:



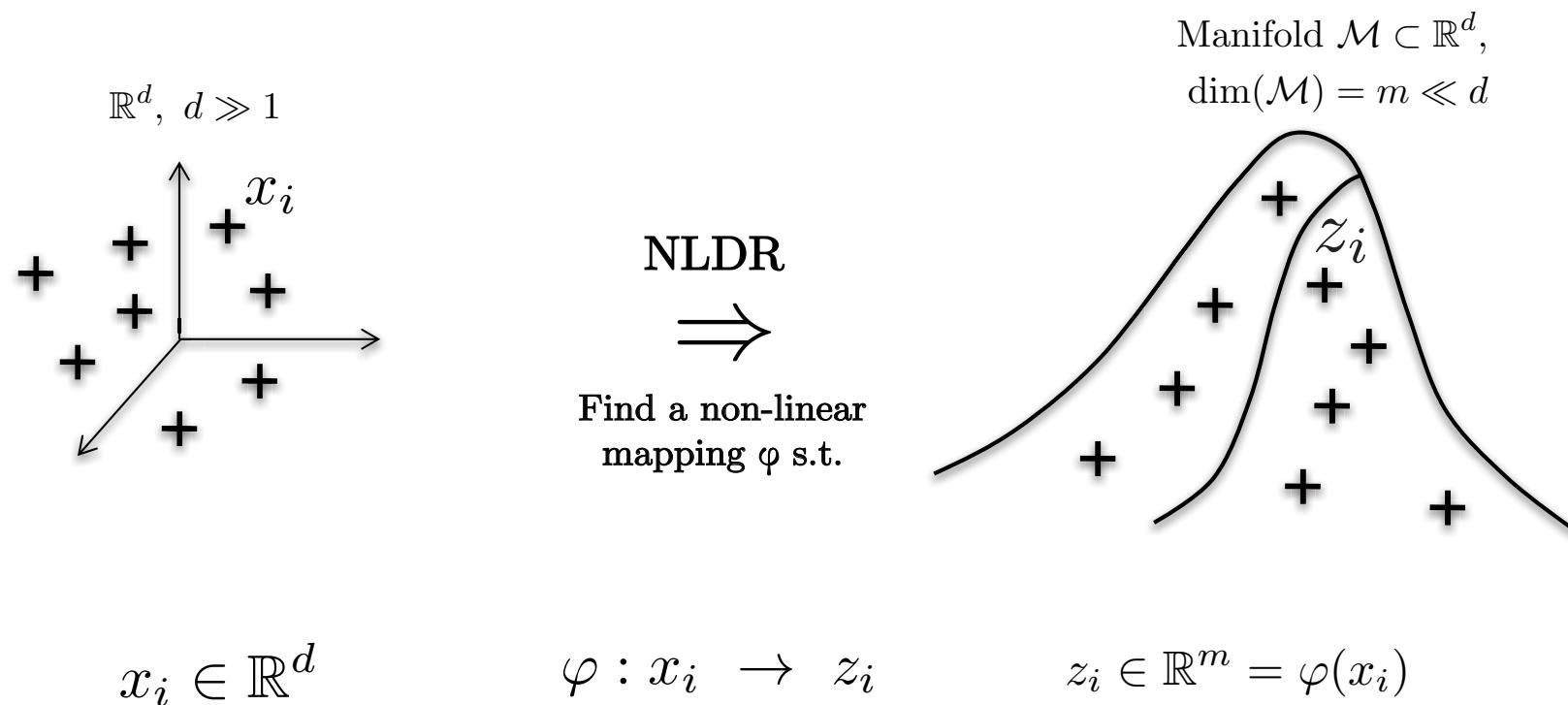
# Visualization Techniques

- Visualization techniques are also **dimensionality reduction techniques** because they aim at mapping data into a much lower-dimensional space, 2D or 3D Euclidean spaces.
- **Linear dimensionality reduction (LDR) techniques.**  
*Assumption:* Data that can be represented on a **low-dimensional hyperplane**.



# Non-Linear Dimensionality Reduction

- *Assumption:* Data that can be represented on **low-dimensional curved spaces**, i.e. **manifolds**:



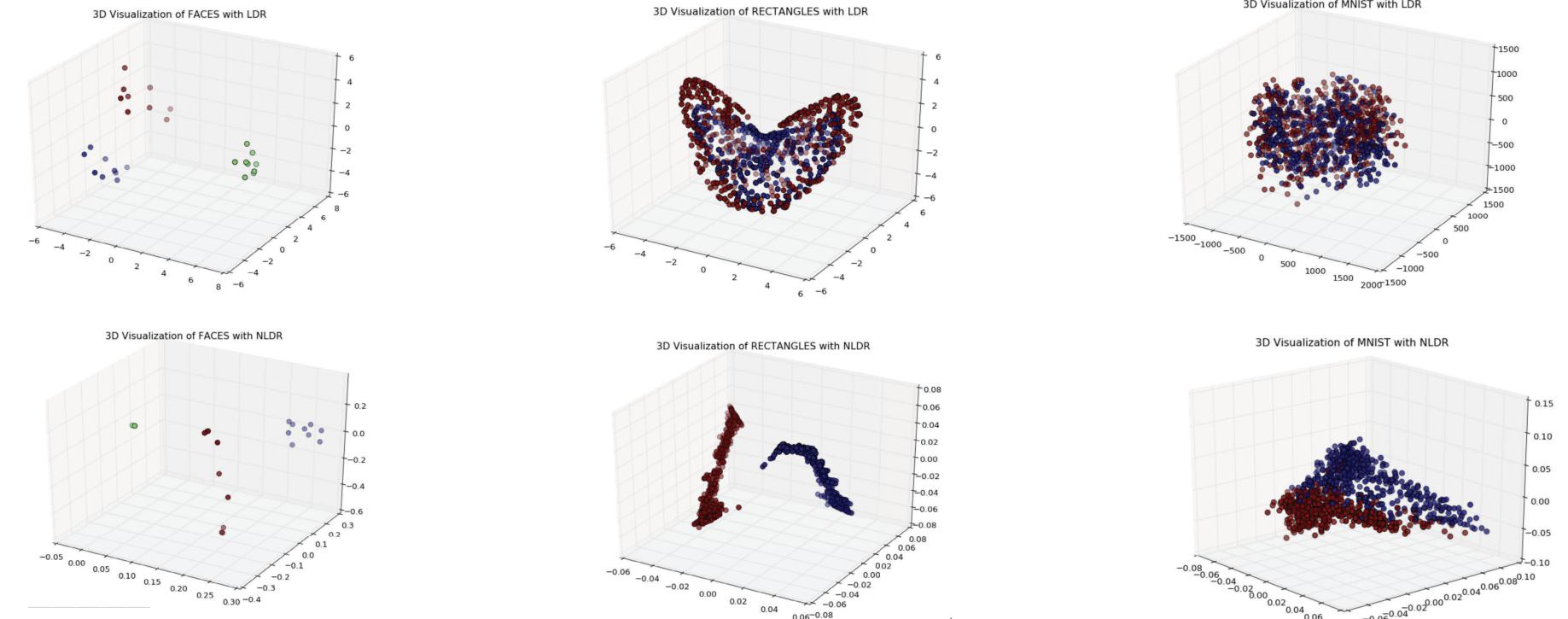
- Non-linear dimensionality reduction (NLDR) techniques are also called **manifold learning** techniques.

# Demo: LDR vs. NLDR

- Run code01.ipynb

```
# Run LDR = Standard PCA
nb_pca = 3
[PC,PD,EnPD] = compute_pca(X,nb_pca)
Xldr = PC[:,0]
Yldr = PC[:,1]
Zldr = PC[:,2]

# Run NLDR = Laplacian Eigenmaps
W = construct_knn_graph(X,5,'euclidean')
Xnldr,Ynldr,Znldr = nldr_visualization(W)
```



# Outline

- Visualization Problem
- Kernel PCA
- Locally-Linear Embedding (LLE)
- Laplacian Eigenmaps
- T-SNE
- Conclusion

# Kernel PCA [Scholkopf-Smola-Muller'97]

- Standard PCA:

Gram matrix:  $G = XX^T = \begin{bmatrix} \langle x_1, x_1 \rangle & \langle x_1, x_2 \rangle & \dots \\ \langle x_2, x_1 \rangle & \langle x_2, x_2 \rangle & \ddots \\ \vdots & & \ddots \\ & & \langle x_n, x_n \rangle \end{bmatrix} \stackrel{EVD}{=} UDU^T \Rightarrow X_{pca} = UD^{1/2}$

- Kernel PCA: Gram matrix in higher-dim space:

$$G = \begin{bmatrix} \langle \phi(x_1), \phi(x_1) \rangle & \langle \phi(x_1), \phi(x_2) \rangle & \dots \\ \langle \phi(x_2), \phi(x_1) \rangle & \langle \phi(x_2), \phi(x_2) \rangle & \ddots \\ \vdots & & \ddots \\ & & \langle \phi(x_n), \phi(x_n) \rangle \end{bmatrix} \stackrel{EVD}{=} UDU^T \Rightarrow X_{kPCA} = UD^{1/2}$$

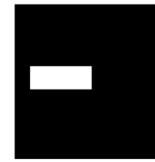
Apply the **kernel trick** to the Gram matrix:

$$K(x, y) = \underbrace{\langle \phi(x), \phi(y) \rangle}_{\text{Never computed}} = e^{-\|x-y\|_2^2/\sigma^2}$$

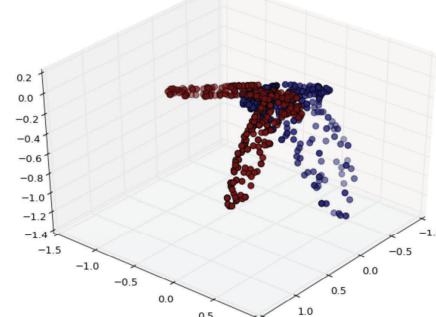
# Demo: Kernel PCA

- Run code02.ipynb

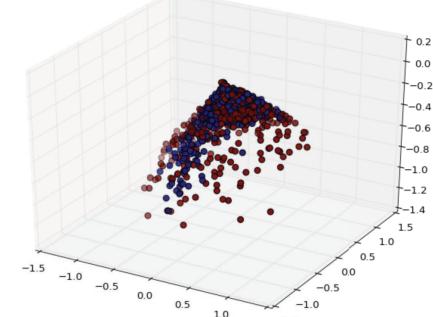
```
# Run NLDR = Kernel PCA
Ker = construct_kernel(X,'polynomial',[2.,0.,3.])
nb_kernel_pca = 3
lamb, U = scipy.sparse.linalg.eigsh(Ker, k=nb_kernel_pca, which='LM')
D = np.diag(lamb)
kernel_PCs = U.dot( pow(D[:,nb_kernel_pca], 0.5) )
Xnldr_kernelPCA = kernel_PCs[:,0]
Ynldr_kernelPCA = kernel_PCs[:,1]
Znldr_kernelPCA = kernel_PCs[:,2]
```



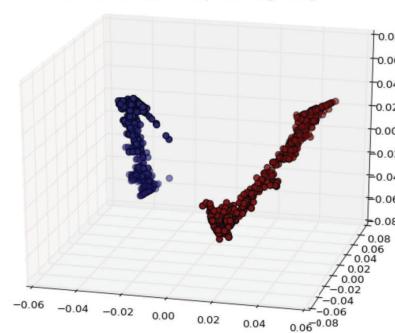
3D Visualization with Kernel PC



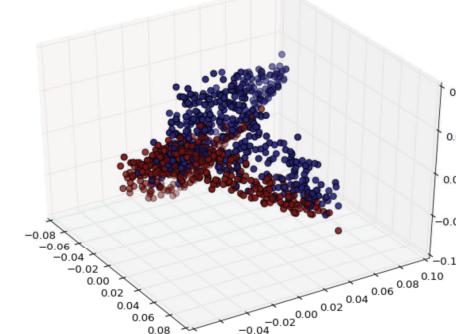
3D Visualization with Kernel PC



3D Visualization of Laplacian Eigenmaps



3D Visualization of Laplacian Eigenmaps

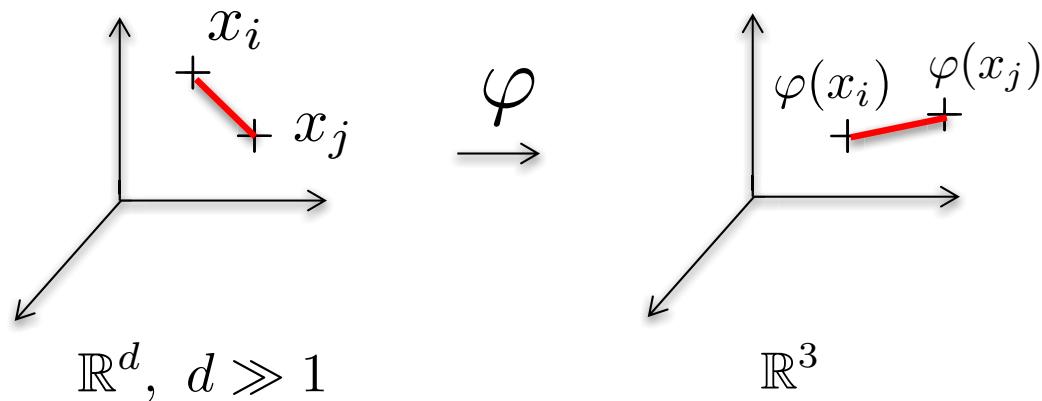


# Outline

- Visualization Problem
- Kernel PCA
- Locally-Linear Embedding (LLE)
- Laplacian Eigenmaps
- T-SNE
- Conclusion

# Locally-Linear Embedding (LLE) [Roweis, Saul'00]

- Motivation: Design a **mapping** from high-dimensional space to low-dimensional space such that the **geometric distances between neighbor data are preserved**.



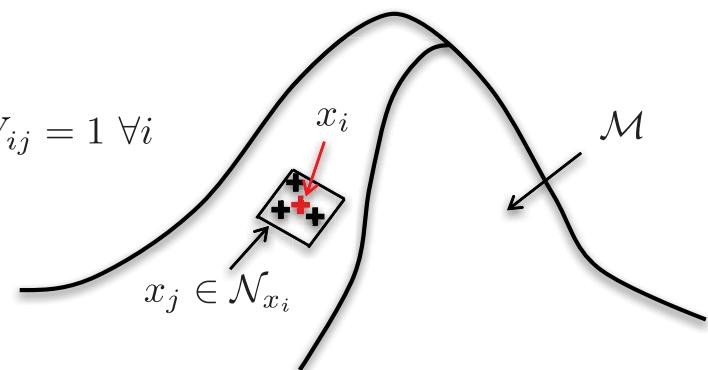
- Description: LLE aims at computing a manifold  $M$  by locally linear fits, that is each data on  $M$  and its neighbors lie on a locally linear patch of  $M$ .

# Algorithm

- **Step 1:** For each data  $x_i$ , compute the  $k$  nearest neighbors.
- **Step 2:** Compute linear patches: find the weights  $W_{ij}$  which best linearly reconstruct  $x_i$  from its neighbors:

$$\min_W \sum_{i=1}^n \left\| x_i - \sum_j W_{ij} x_j \right\|_2^2 \text{ s.t. } \sum_j W_{ij} = 1 \quad \forall i$$

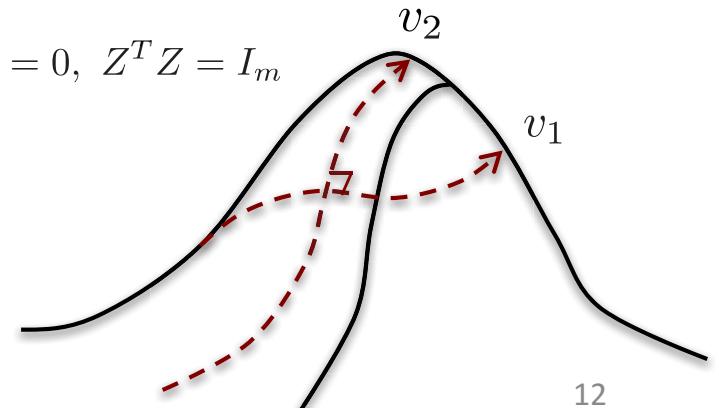
$\Rightarrow$  Solution:  $Ax=b$



- **Step 3:** Compute the low-dimensional embedding data  $z_i$ , which is best reconstructed by the weights  $W_{ij}$ :

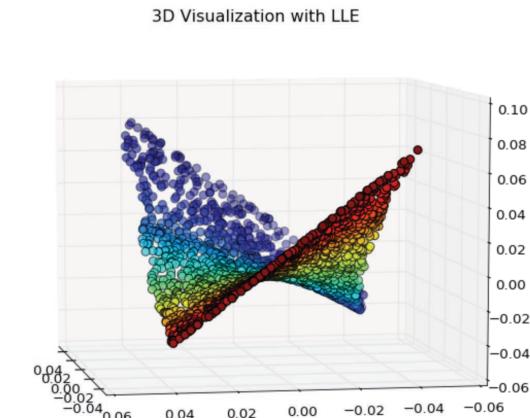
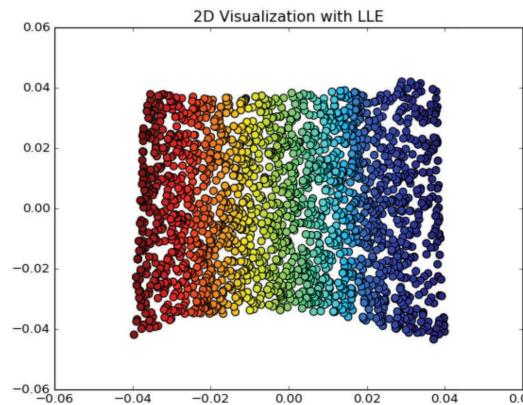
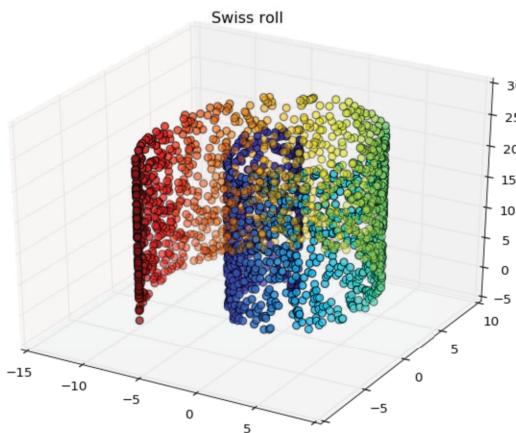
$$\min_{Z=[z_1, \dots, z_m]} \sum_{i=1}^n \left\| z_i - \sum_j W_{ij} z_j \right\|_2^2 \text{ s.t. } \sum_i z_i = 0, \quad Z^T Z = I_m$$

$\Rightarrow$  Solution: EVD



# Demo: LLE

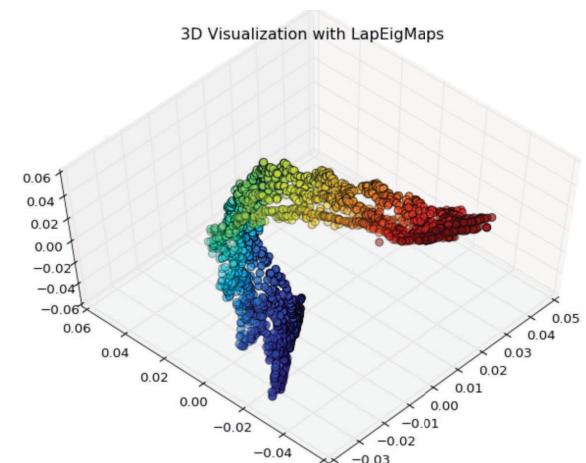
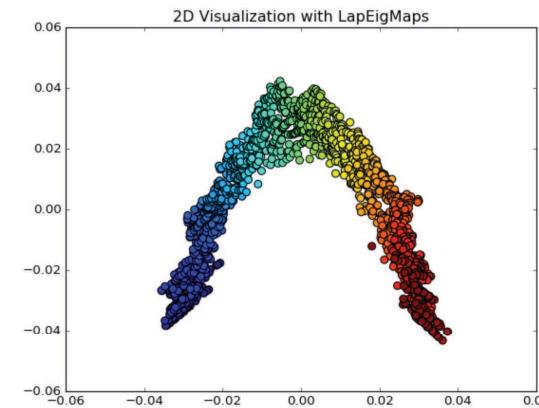
- Run `code03.ipynb`



```
# Run NLDR = LLE
X = Xref
X = X - np.mean(X, axis=0) # zero-centered data

# Step 1: Compute k-NN
kNN = 20
WkNN = construct_knn_graph(X, kNN, 'euclidean').todense()

# Step 2: Compute locally linear patches
W = np.zeros([n,n])
for i in range(n):
    # Find neighbors of data i
    idx_kNN = np.where(WkNN[i,:]>0.0)[1]
    K = len(idx_kNN)
    if K>kNN:
        K = kNN
        idx_kNN = idx_kNN[:K]
    XkNN = X[idx_kNN,:]
```

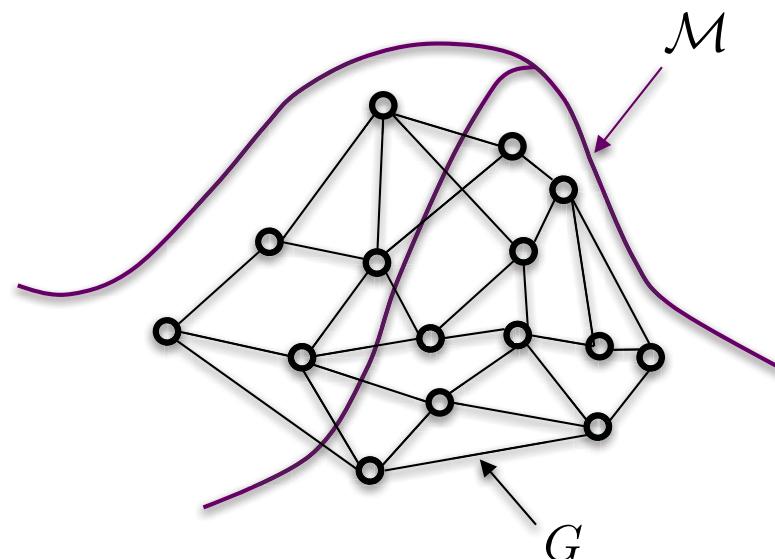


# Outline

- Visualization Problem
- Kernel PCA
- Locally-Linear Embedding (LLE)
- Laplacian Eigenmaps
- T-SNE
- Conclusion

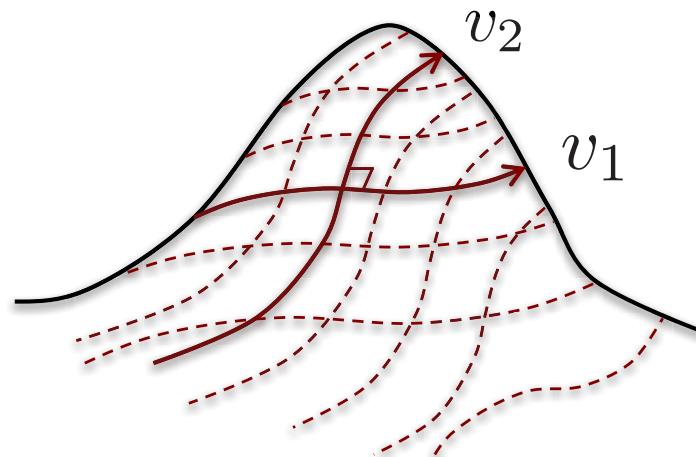
# Laplacian Eigenmaps [Belkin, Niyogi'03]

- Very popular visualization technique.
- **Motivation:** Same as LLE but stronger mathematical analysis and understanding.
- **Manifold assumption:** Data are sampled from a manifold  $M$  represented by a  $k$  - nearest neighbors graph.

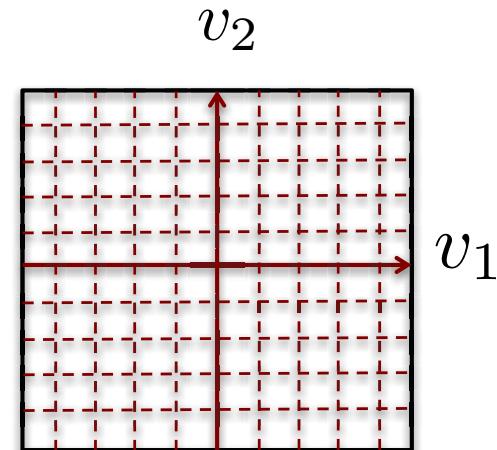


# Differential Geometry

- Eigenfunctions  $v_k$  of continuous Laplace-Beltrami  $\Delta_M$  serve as **embedding coordinates** of  $M$ :



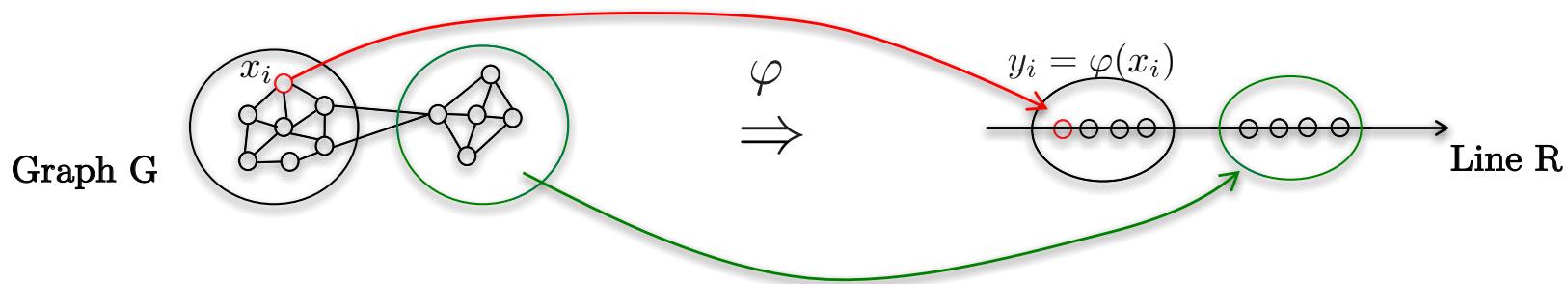
$\varphi$   
⇒  
Unwrap



- Note: **Discretization** of  $\Delta_M$  = graph Laplacian  $L$

# Formalization

- 1D Visualization: Map a graph  $G=(V,E,W)$  to a line such that **neighbor data** on  $G$  stay as close as possible on the line.



- Note: We look for the mapping  $\varphi$  but we **never** compute it explicitly.  
We look for the coordinates  $y_i$  of  $x_i$  on the low-dimensional manifold  $M$  such that  $y_i = \varphi(x_i)$ , that is:

$$\min_y \sum_{ij} W_{ij} (y_i - y_j)^2 \quad (1)$$

- Interpretation: As  $W_{ij}=1$  if  $x_i$  close to  $x_j$ , then  $\min_y \sum W_{ij}(y_i - y_j)^2$  implies that  $y_i$  to be close to  $y_j$ .

# Generalization to 2D and 3D

- **K-D Visualization:** Generalizing (1) to K dimensions is straightforward:

$$\min_Y \sum_k Y_{\cdot,k}^T L Y_{\cdot,k} = \text{tr}(Y^T \cancel{LY}) \quad \text{s.t.} \quad Y^T Y = I_K$$

Graph Laplacian

- **Spectral Solution:** Top K eigenvectors of graph Laplacian L.

$$L \stackrel{EV D}{=} U \Lambda U^T \rightarrow Y = U_K$$

- **Advantages:**

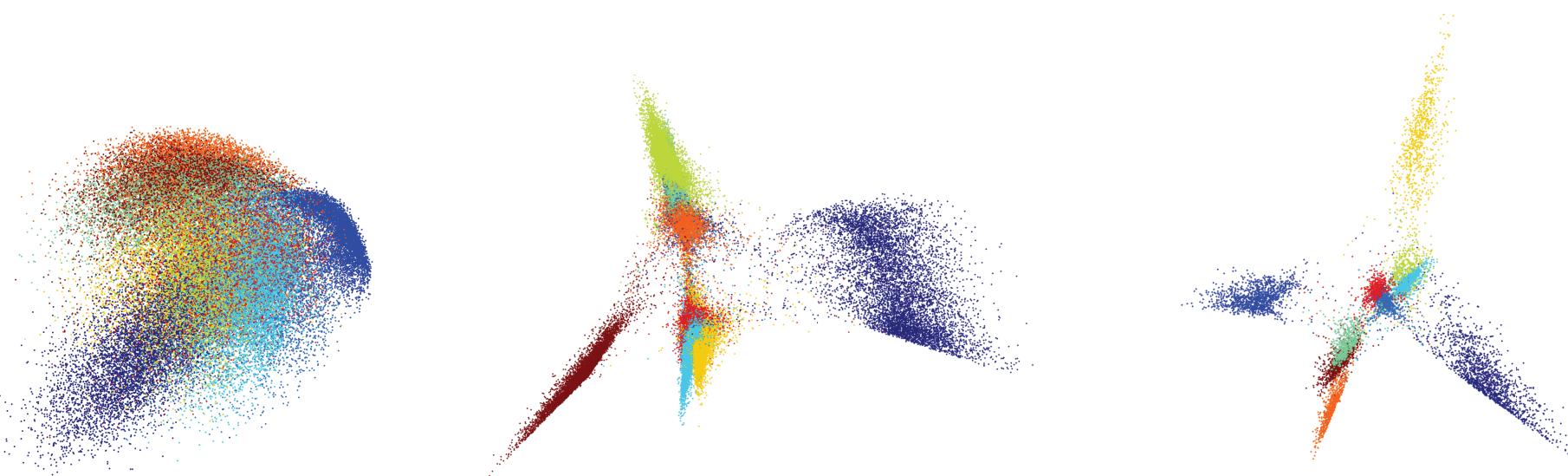
- (1) *Global solutions (independent of initialization)*
- (2) *Fast algorithms (for small n)*

# Demo: Laplacian Eigenmaps

- Run `code04.ipynb`

```
mat = scipy.io.loadmat('datasets/MNIST.mat')
W = mat['W']
n = W.shape[0]
Cgt = mat['C'].squeeze()
print(n)

# Run NLDR = Laplacian Eigenmaps
Xnldr_lapeigmap,Ynldr_lapeigmap,Znldr_lapeigmap = nldr_visualization(W)
```



MNIST  
PCA

MNIST  
Lap Eigenmaps

USPS  
Lap Eigenmaps

# Outline

- Visualization Problem
- Kernel PCA
- Locally-Linear Embedding (LLE)
- Laplacian Eigenmaps
- T-SNE
- Conclusion

# T-SNE [van der Maaten, Hinton'08]



*Q: What visualization technique is the most used?*

- t-Distribution Stochastic Neighbor Embedding (t-SNE) is the **most popular visualization** technique.
- **Model description:** t-SNE learns the mapping/embedding  $\varphi$  function, such that  $y_i = \varphi(x_i)$ , by **minimizing the Kullback-Leibler distance** between the distribution of high-dim data and the distribution of the computed low-dim data:

$$\left\{ \begin{array}{l} p_{ij} = \frac{e^{-\|x_i - x_j\|_2^2 / \sigma_i^2}}{\sum_k e^{-\|x_i - x_k\|_2^2 / \sigma_i^2}} \\ \quad \sigma_i = k^{th} \text{ nearest neighbor distance from } x_i \\ q_{ij}(y) = \frac{(1 + \|y_i - y_j\|_2^2)^{-1}}{\sum_k (1 + \|y_i - y_k\|_2^2)^{-1}} \\ \quad \text{Embedding coordinates of high-dim data} \end{array} \right.$$

# Optimizing Kullback-Leibler

- Problem:

$$\min_y \text{KL}(P, Q(y)) = \sum_{ij} p_{ij} \log \frac{p_{ij}}{q_{ij}(y)}$$

$\Downarrow$  Gradient descent technique

$$y_i^{m+1} = y_i^m - 4\tau \sum_j (p_{ij} - q_{ij})(1 + \|x_i - x_j\|_2)^{-1} (y_i^m - y_j)$$

- Advantages:

- (1) *Local distance preservation (as LapEig, LLE):* minimizing KL forces  $q_{ij}$  to be close to  $p_{ij}$ , the distribution of high-dim data.
- (2) *t-SNE does not assume the existence of a manifold:* More **flexibility** to visualize more complex hidden structures.

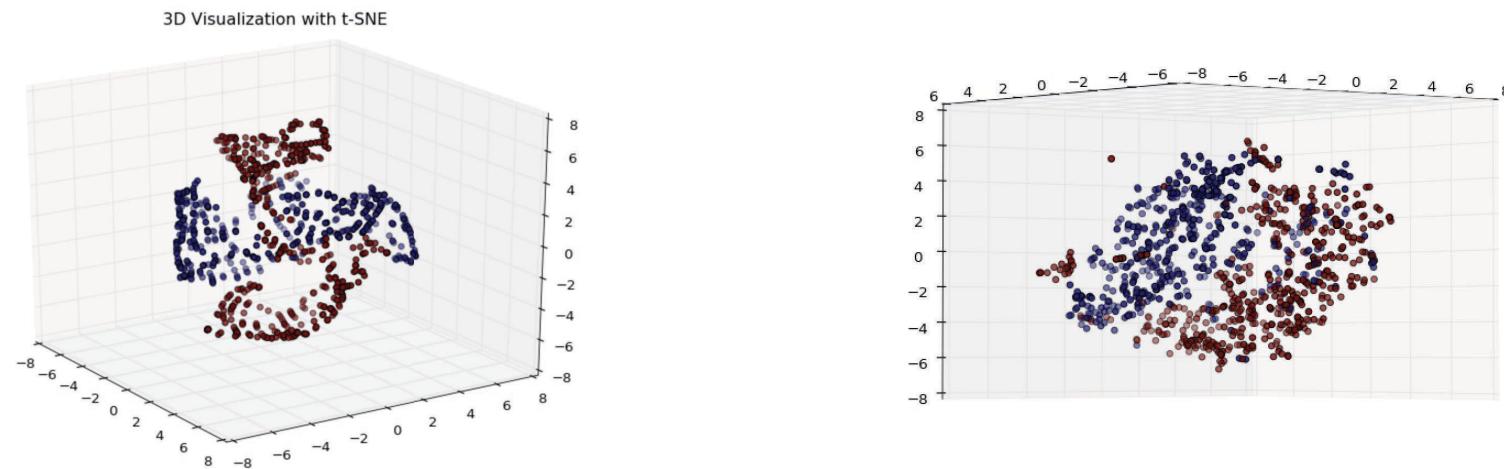
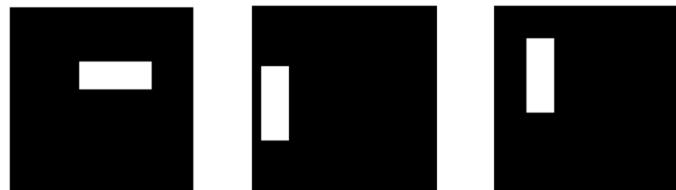
- Limitations:

- (1) *Non-convex energy*  $\Rightarrow$  Existence of **bad** local solutions, problem of initialization (PCA is used as initialization).
- (2) *Slow optimization* (gradient descent).

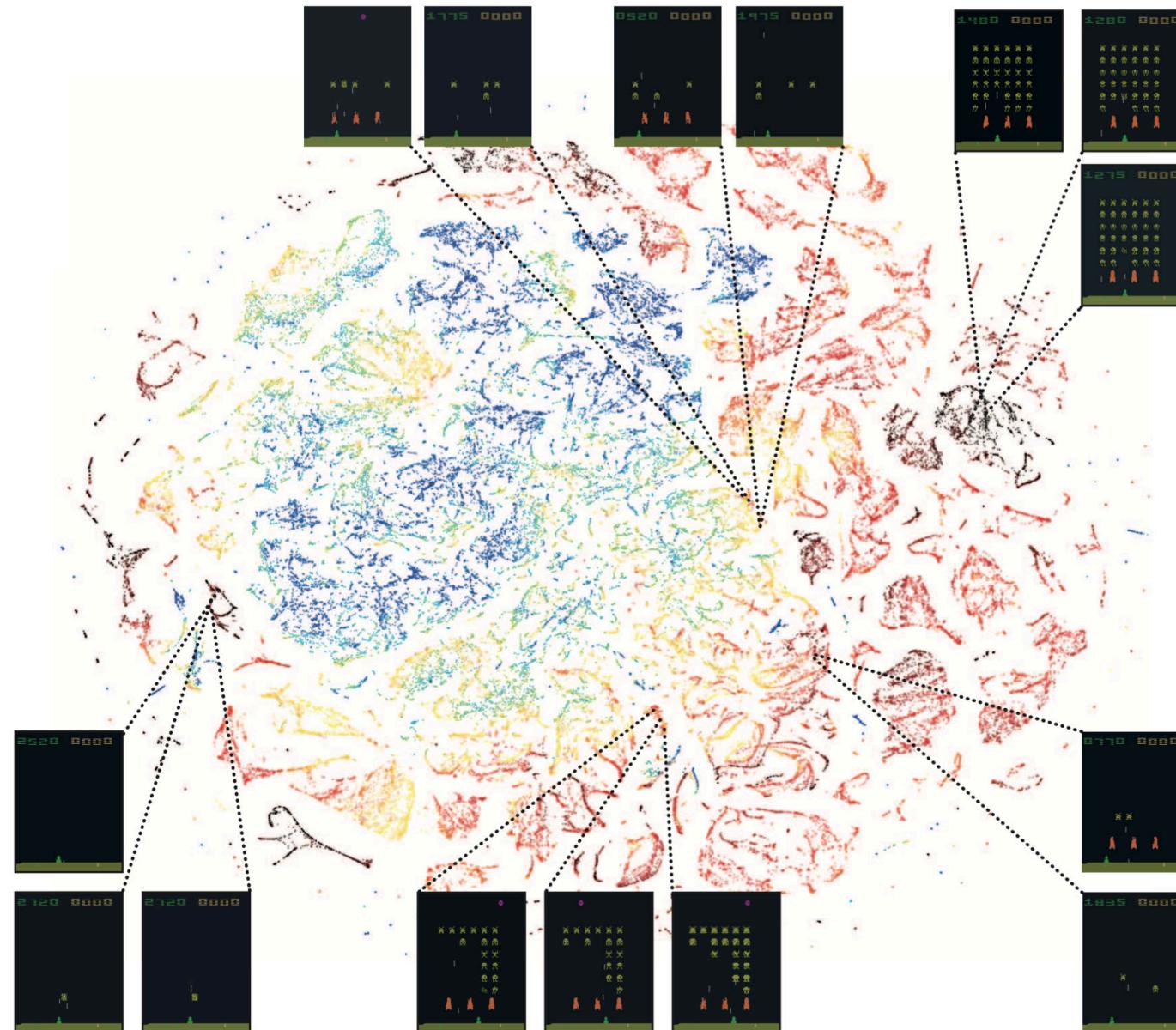
# Demo: T-SNE

- Run code05.ipynb

```
# Run t-SNE
model = TSNE(n_components=3, n_iter=2000, learning_rate=1000)
np.set_printoptions(suppress=True)
Xtsne = model.fit_transform(X) # Xtsne = n x d
Xtsne = Xtsne[:,0]
Ytsne = Xtsne[:,1]
Ztsne = Xtsne[:,2]
```



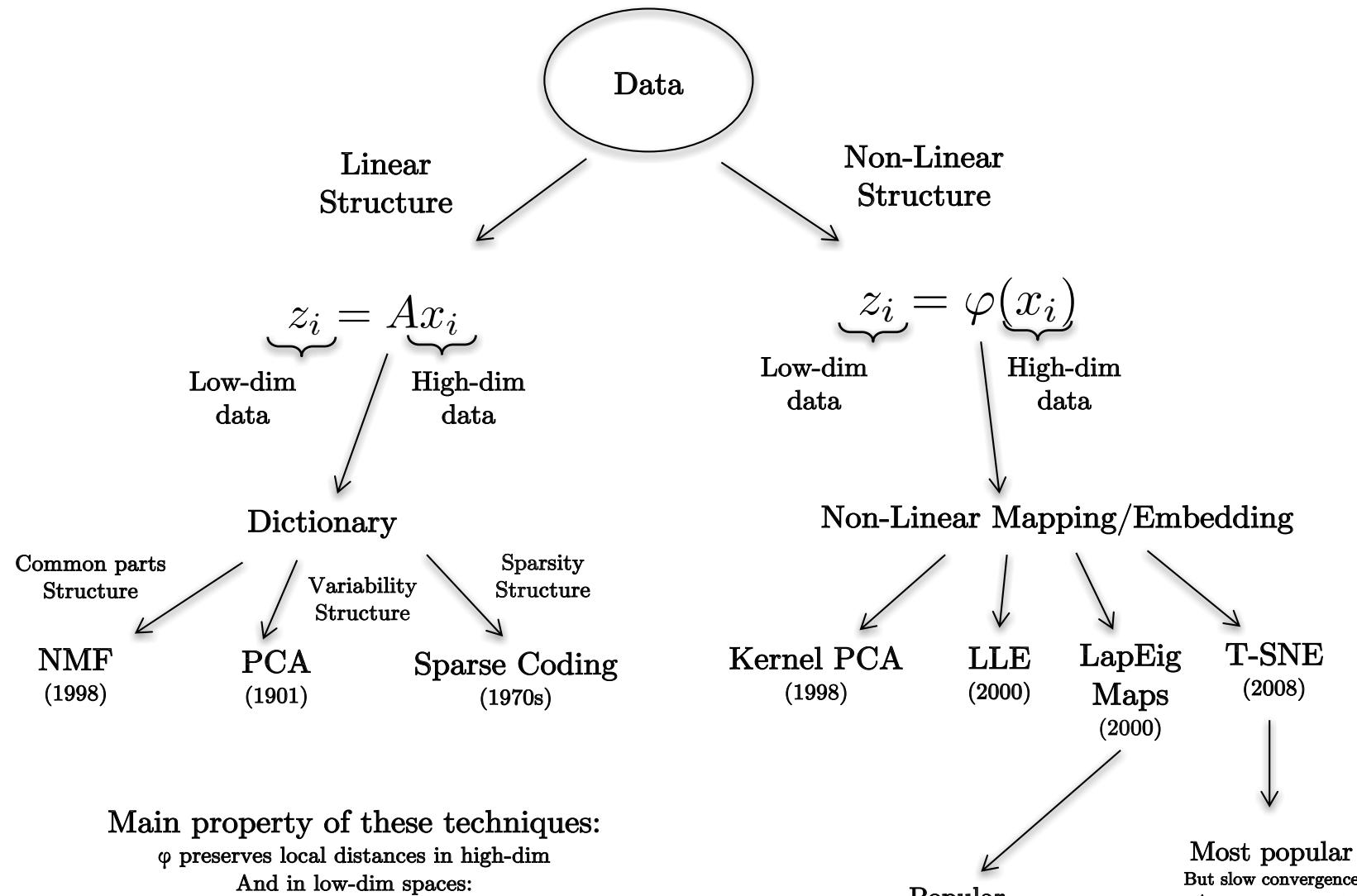
# Visualizing Video Games [Mnih-et.al'14]



# Outline

- Visualization Problem
- Kernel PCA
- Locally-Linear Embedding (LLE)
- Laplacian Eigenmaps
- T-SNE
- Conclusion

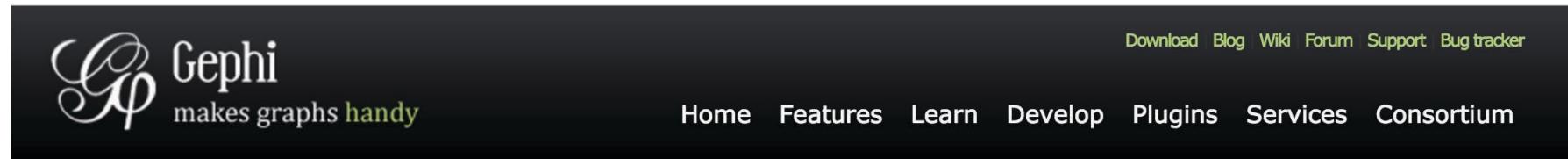
# Summary



**Popular**  
 Math sound  
 Unique solution  
 Manifold assumption  
 Too strong

**Most popular**  
 But slow convergence  
 As non-convex opt  
 Local minimizers

# Gephi



The Gephi website features a dark header bar. On the left is the Gephi logo and tagline "makes graphs handy". On the right are links for "Download", "Blog", "Wiki", "Forum", "Support", and "Bug tracker". Below the header are navigation links for "Home", "Features", "Learn", "Develop", "Plugins", "Services", and "Consortium".

## The Open Graph Viz Platform

**Gephi is the leading visualization and exploration software for all kinds of graphs and networks. Gephi is open-source and free.**

**Runs on Windows, Mac OS X and Linux.**

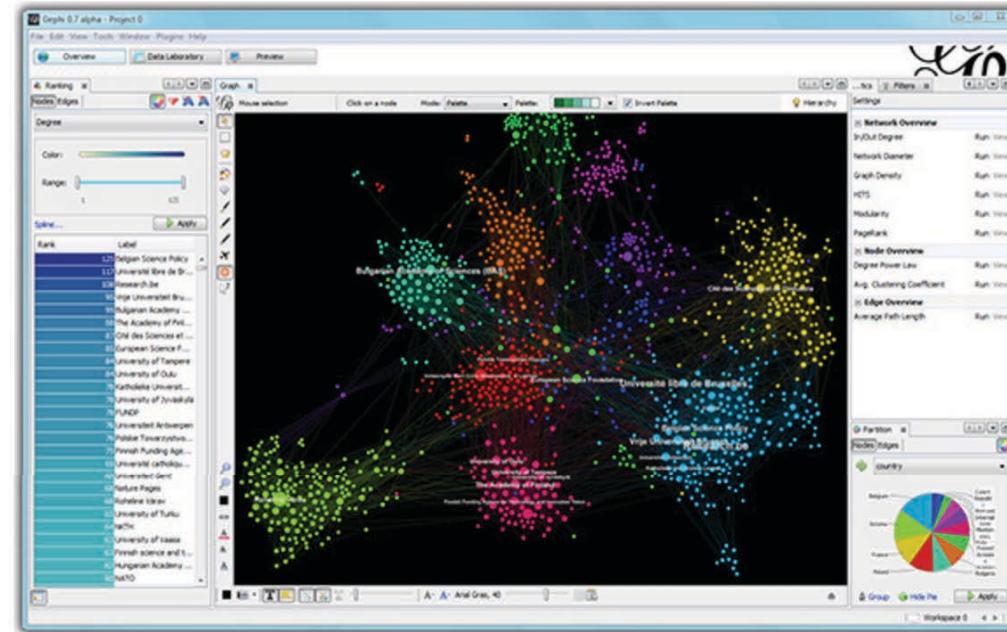
[Learn More on Gephi Platform »](#)

 Download FREE  
Gephi 0.9.1

[Release Notes](#) | [System Requirements](#)

► [Features](#)  
► [Quick start](#)

► [Screenshots](#)  
► [Videos](#)



- Awesome visualization tool!
- Run `code06.ipynb` to convert graphs to Gephi format.

# Mining Star Wars Universe

100+ international articles including popular **Dailymail** (1.66M Twitter followers), **Gizmodo** (1.48M), **Engadget** (1.57M), **MUY Interesante** (7.1M).

**GIZMODO**  
Computer Analysis Reveals the Stunning Complexity of the Star Wars Expanded Universe

George Dvorsky  
2/10/16 10:43am - Filed to: COMPUTER SCIENCE

A screenshot of a Gizmodo article titled "Computer Analysis Reveals the Stunning Complexity of the Star Wars Expanded Universe". The article features a large image of two characters from Star Wars: a small human-like figure and a large, metallic, cylindrical droid. Below the image is a network graph where nodes represent characters and connections represent interactions. The graph is dense and colorful, with many red, green, and blue nodes.

**SCIENCE ALERT**  
Data scientists map every important character in the Star Wars universe

Even Jar Jar.  
PETER DOCKRILL - 11 FEB 2016

A screenshot of a Science Alert article titled "Data scientists map every important character in the Star Wars universe". The article includes a network graph showing the relationships between characters. A specific node, "Even Jar Jar", is highlighted. The graph is a complex web of colored lines connecting numerous dots.

**engadget**  
Software maps the 'Star Wars' universe

Over 20,000 characters and 36,000 years distilled by clever algorithms.

Jon Fingas, @jonfingas  
02.10.16 in AV

468 Shares

A screenshot of an Engadget article titled "Software maps the 'Star Wars' universe". The article features a large, detailed network graph of Star Wars characters. The graph is highly interconnected, with many nodes and a variety of colors representing different species or groups.

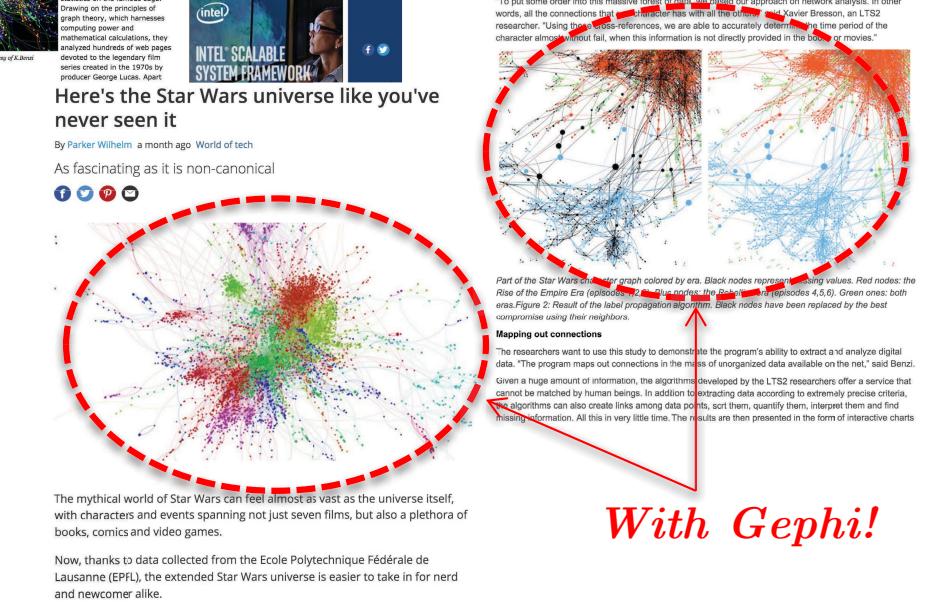
**PHYS.ORG**  
Nanotechnology Physics Earth Astronomy & Space Technology Chemistry

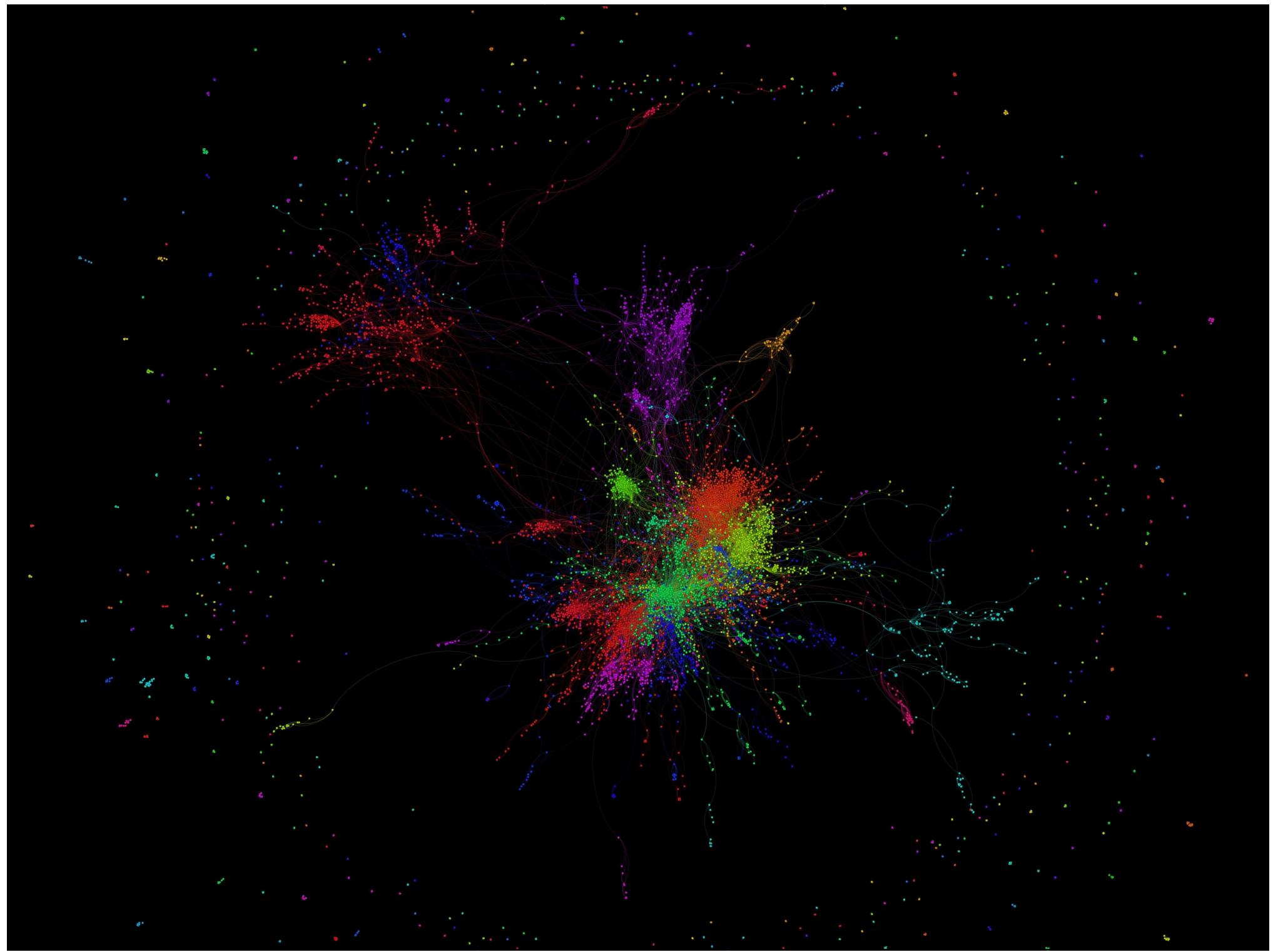
Home > Technology > Computer Sciences > February 10, 2016

**Math reveals unseen worlds of Star Wars**

February 10, 2016 by Sarah Perrin

A screenshot of a Phys.org article titled "Math reveals unseen worlds of Star Wars". The article includes a network graph showing the connections between 7583 Star Wars characters. Below the graph is a pie chart illustrating the distribution of characters across various species: Rodian (blue), Wookiee (orange), Bothan (purple), Yuuzhan Vong (yellow), Trandoshan (green), Zabrak (red), Sullustan (light blue), and Duros (pink).







Questions?