

# Clustering Big Data

Anil K. Jain

(with Radha Chitta and Rong Jin)

Department of Computer Science

Michigan State University

November 29, 2012

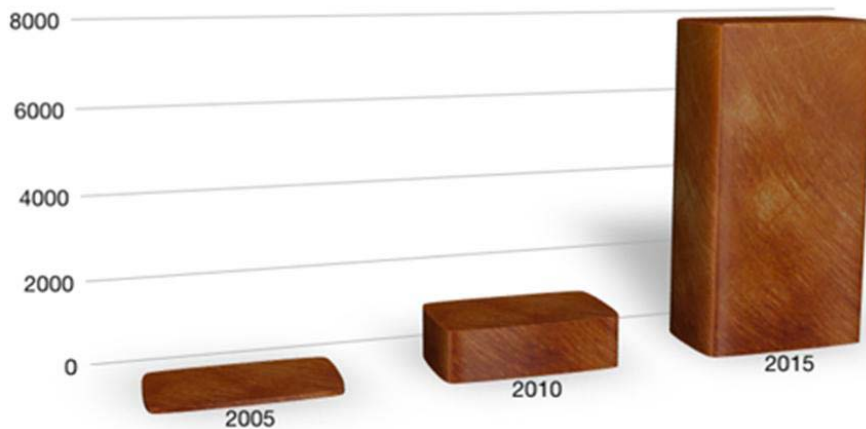
# Outline

- Big Data
- How to extract “information”?
- Data clustering
- Clustering Big Data
- Kernel K-means & approximation
- Summary

# How Big is Big Data?

- **Big** is a fast moving target: kilobytes, megabytes, gigabytes, terabytes ( $10^{12}$ ), petabytes ( $10^{15}$ ), exabytes ( $10^{18}$ ), zettabytes ( $10^{21}$ ),.....
- Over 1.8 zb created in 2011; ~8 zb by 2015

D  
E  
a  
t  
a  
b  
y  
s  
i  
z  
e



Source: IDC's Digital Universe study, sponsored by EMC, June 2011

<http://idcdocserv.com/1142>

<http://www.emc.com/leadership/programs/digital-universe.htm>



As of June 2012

Nature of Big Data: Volume, Velocity and Variety

# Big Data on the Web



~900 million users, 2.5 billion content items, 105 terabytes of data each half hour, 300M photos and 4M videos posted per day

Over 225 million users generating over 800 tweets per second

twitter



<http://techcrunch.com/2012/08/22/how-big-is-facebooks-data-2-5-billion-pieces-of-content-and-500-terabytes-ingested-every-day/>  
<http://royal.pingdom.com/2012/01/17/internet-2011-in-numbers/>  
<http://www.dataversity.net/the-growth-of-unstructured-data-what-are-we-going-to-do-with-all-those-zettabytes/>

# Big Data on the Web



Over 50 billion pages indexed and more than 2 million queries/min



Articles from over 10,000 sources in real time



~4.5 million photos uploaded/day



48 hours of video uploaded/min; more than 1 trillion video views

No. of mobile phones will exceed the world's population by the end of 2012

# What to do with Big Data?

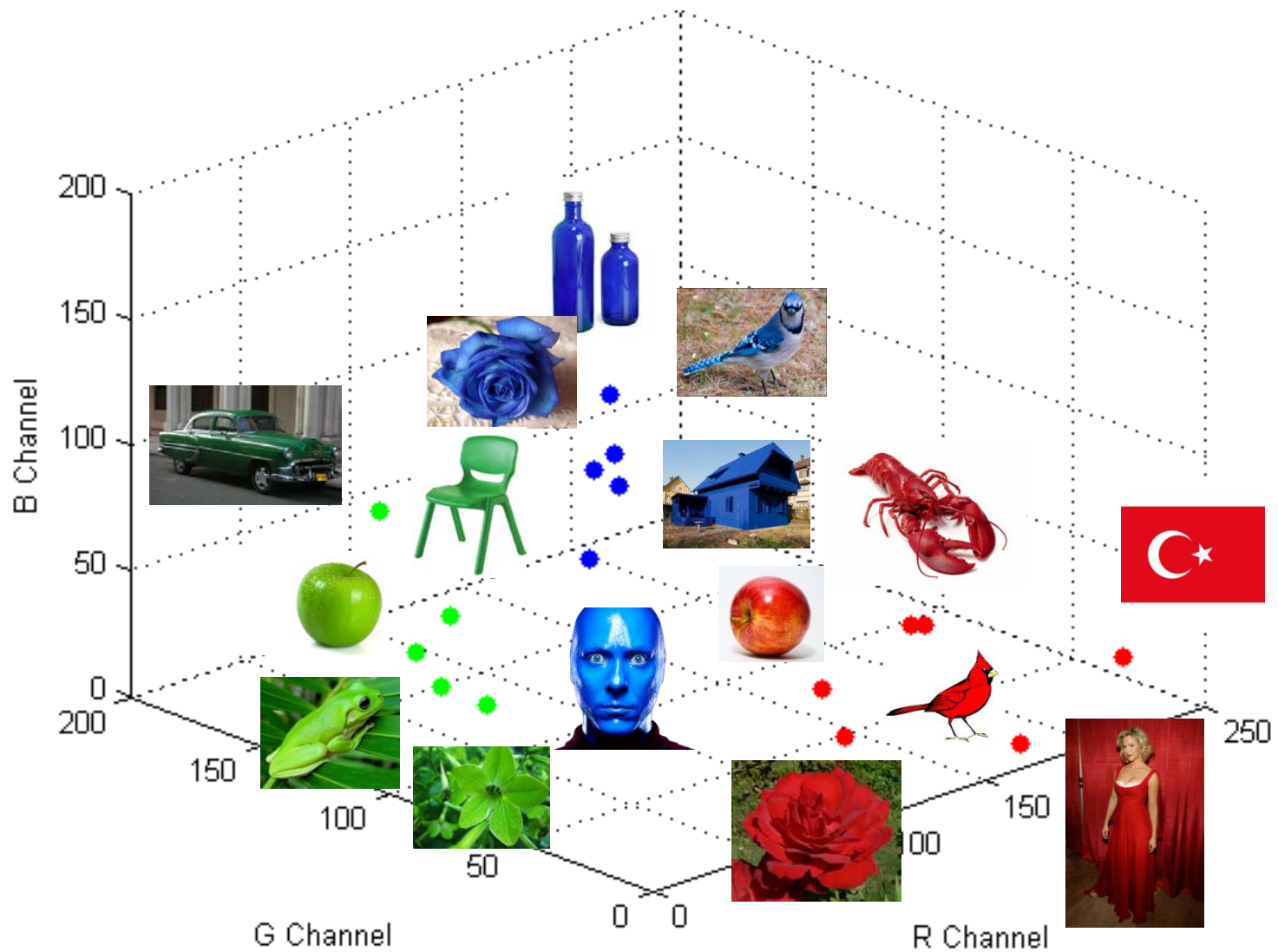
- Extract information to make decisions
- Evidence-based decision: data-driven vs. analysis based on intuition & experience
- Analytics, business intelligence, data mining, machine learning, pattern recognition
- Big Data computing: IBM is promoting Watson (Jeopardy champion) to tackle Big Data in healthcare, finance, drug design,...

# Decision Making

- Data Representation
  - Features and similarity
- Learning
  - Classification (labeled data)
  - Clustering (unlabeled data)

Most big data problems have unlabeled objects

# Pattern Matrix





















$n \times d$  pattern matrix



# Similarity Matrix

Polynomial kernel:  $K(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^T \mathbf{y} + 1)^4$

									
	16	15	14	4	6	6	4	3	1
	15	16	14	4	5	5	6	4	3
	14	14	16	9	9	9	8	7	4
	4	4	9	16	15	15	9	10	6
	6	5	9	15	16	16	7	8	4
	6	5	9	15	16	16	7	8	4
	4	6	8	9	7	7	16	16	14
	3	4	7	10	8	8	16	16	14
	1	3	4	6	4	4	14	14	16

n x n similarity matrix

# Classification



Dogs



Cats

Given a training set of labeled objects, learn a decision rule

# Clustering



Given a collection of (unlabeled) objects, find meaningful groups

# Semi-supervised Clustering

Supervised



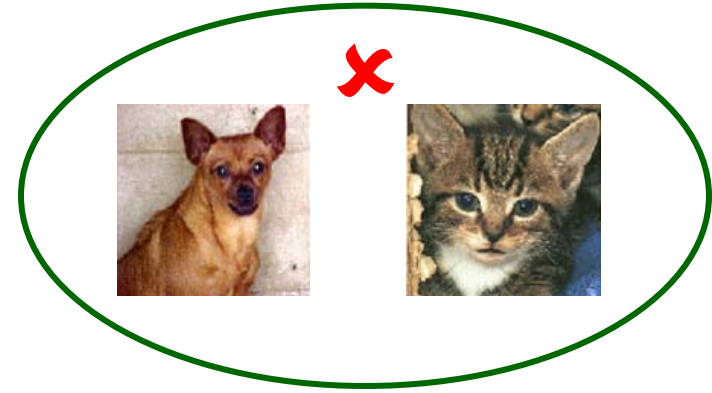
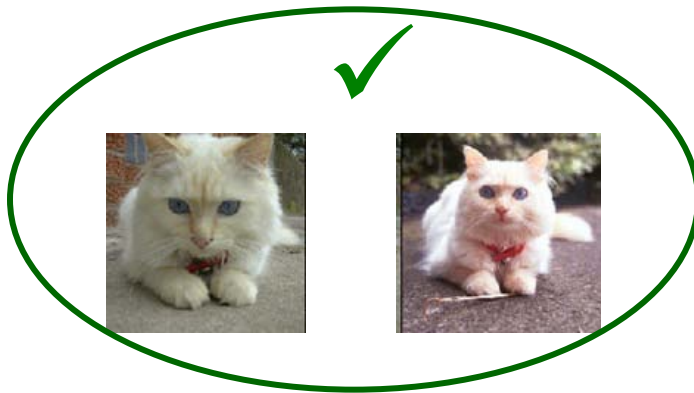
Dogs

Cats

Unsupervised



Semi-supervised



Pairwise constraints improve the clustering performance

# What is a cluster?

*"A group of the same or similar elements gathered or occurring closely together"*



Galaxy clusters



Birdhouse clusters



Cluster munition



Cluster computing

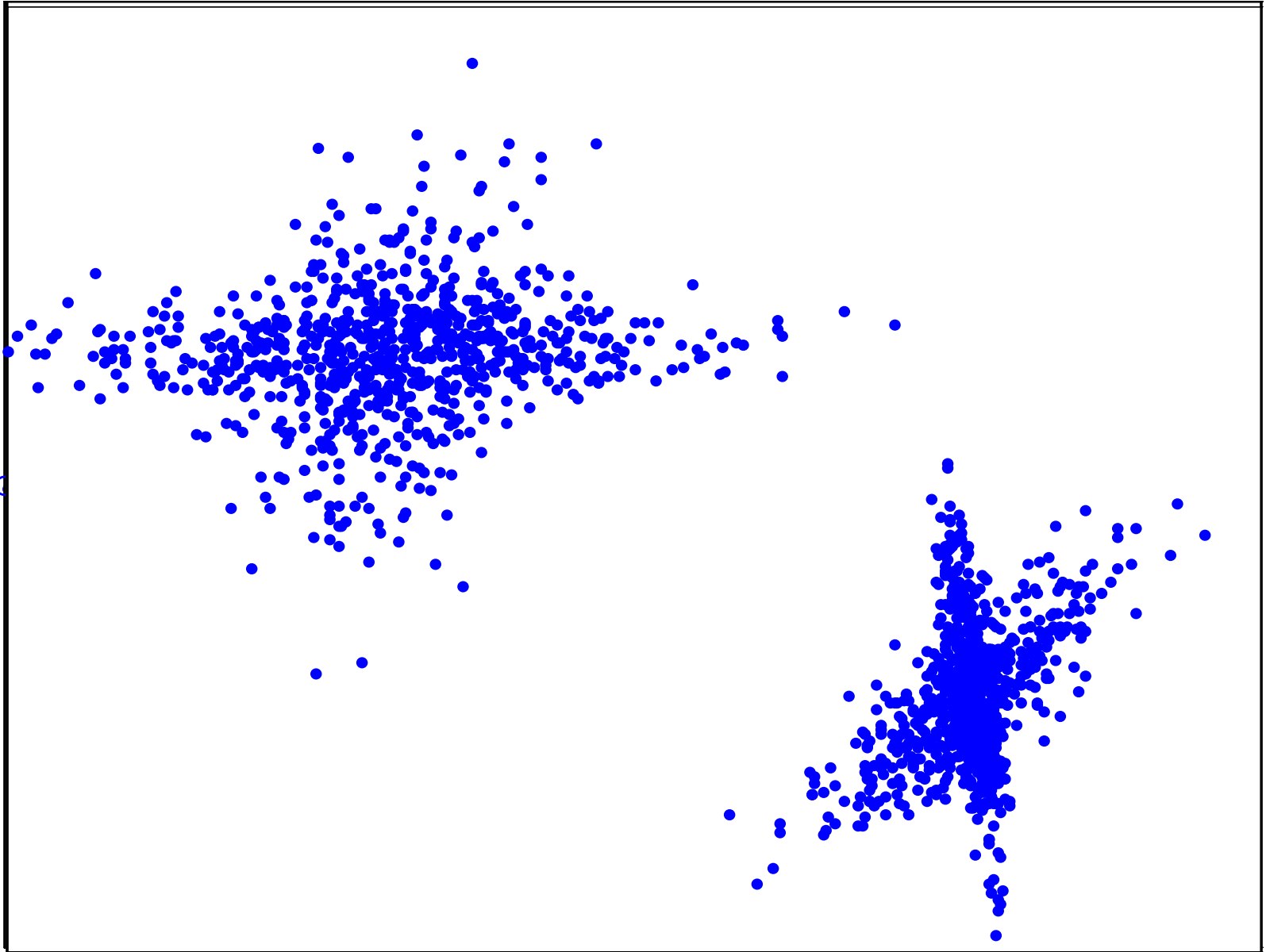


Cluster lights



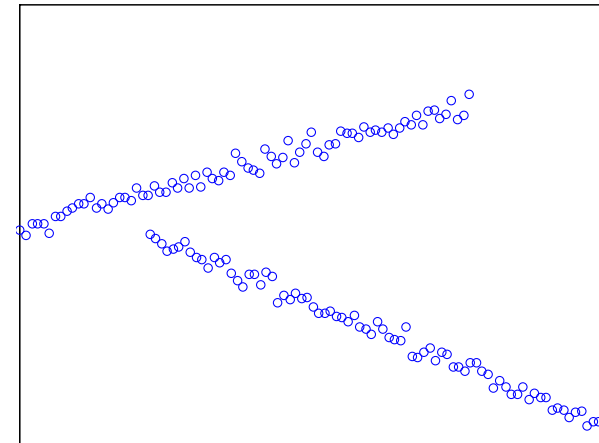
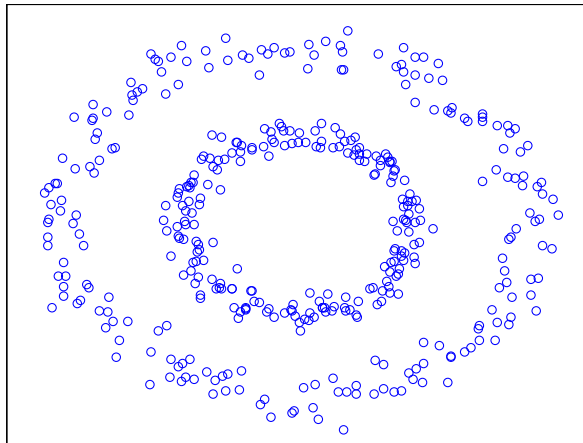
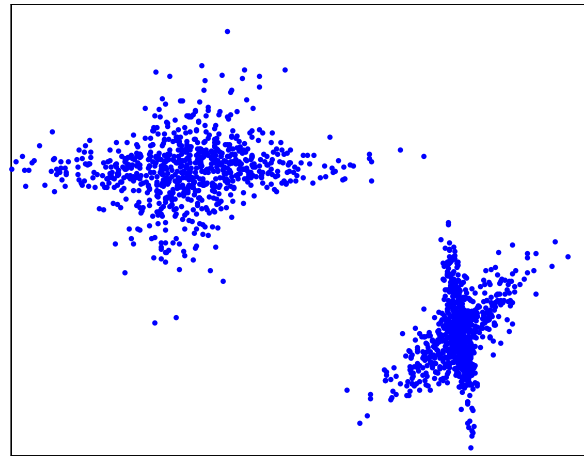
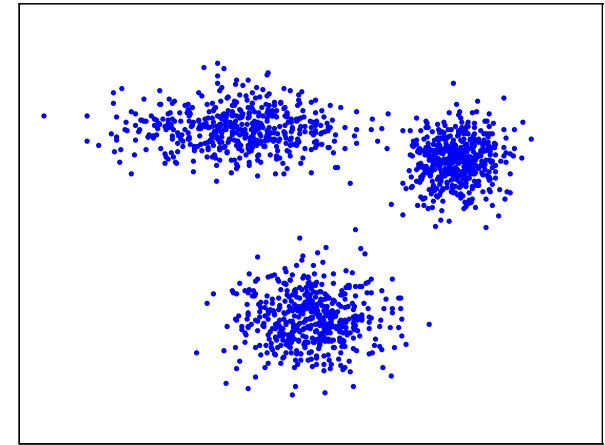
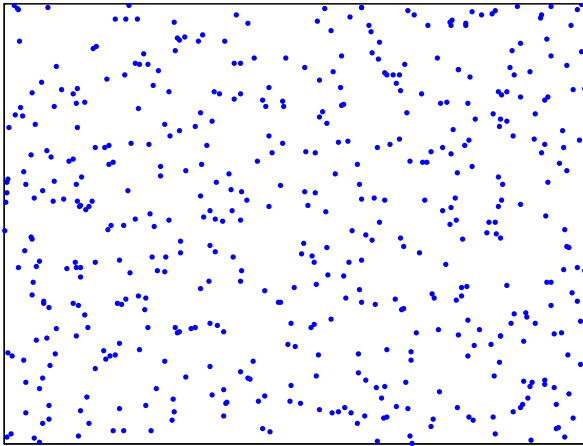
Hongkeng Tulou cluster

# Clusters in 2D



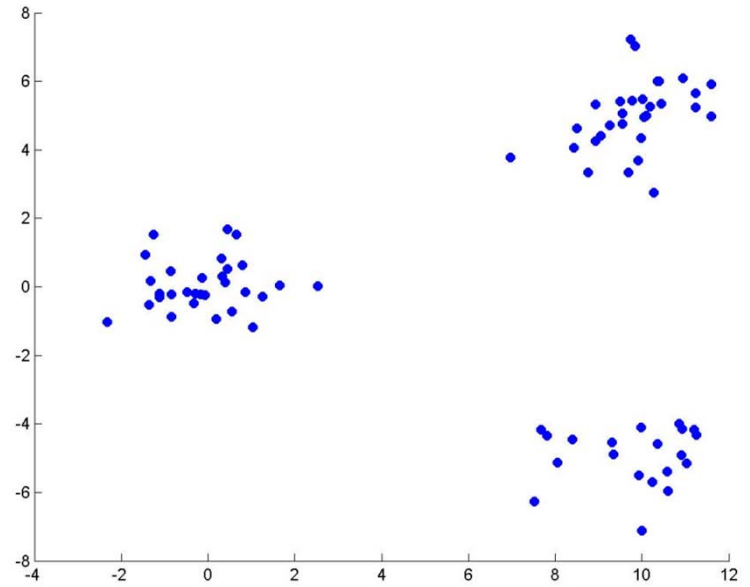
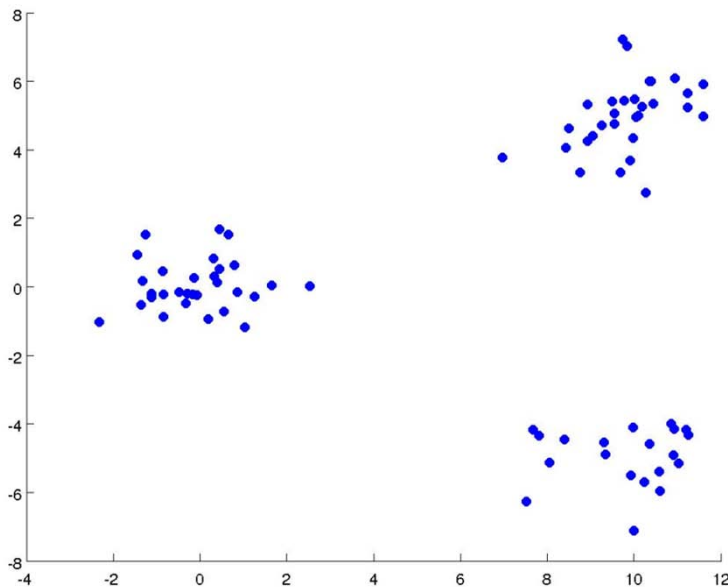
# Challenges in Data Clustering

- Measure of similarity
- No. of clusters
- Cluster validity
- Outliers



# Data Clustering

Organize a collection of  $n$  objects into a **partition** or a **hierarchy** (nested set of partitions)



**“Data clustering” returned ~6,100 hits for 2011 (Google Scholar)**

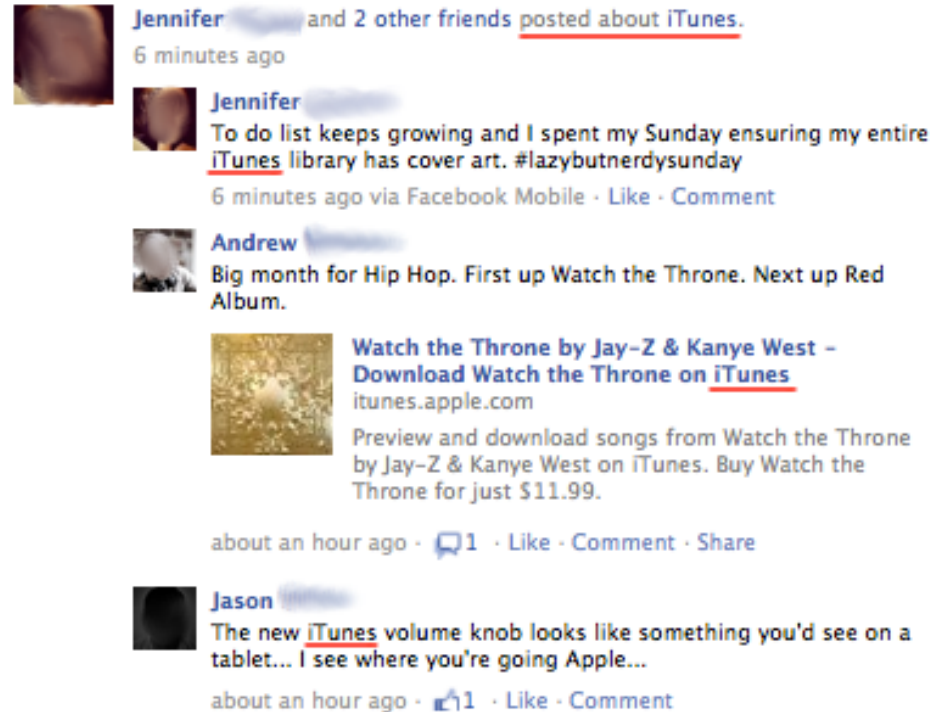


# Clustering is the Key to Big Data Problem

- Not feasible to “label” large collection of objects
- No prior knowledge of the number and nature of groups (clusters) in data
- Clusters may evolve over time
- Clustering provides efficient browsing, search, recommendation and organization of data

# Clustering Users on Facebook

- ~300,000 status updates per minute on tens of thousands of topics
- Cluster users based on topic of status messages




The screenshot shows a vertical feed of Facebook posts. The top post is from Jennifer and 2 other friends, posted about iTunes 6 minutes ago. The second post is from Jennifer, posted 6 minutes ago via Facebook Mobile, mentioning her iTunes library and cover art. The third post is from Andrew, posted about an hour ago, mentioning Hip Hop and the album Watch the Throne. Below this is a promotional link for the album 'Watch the Throne by Jay-Z & Kanye West' on iTunes, with a preview and download option. The bottom post is from Jason, posted about an hour ago, mentioning the new iTunes volume knob.

Jennifer [redacted] and 2 other friends posted about iTunes.  
6 minutes ago

Jennifer [redacted]  
To do list keeps growing and I spent my Sunday ensuring my entire iTunes library has cover art. #lazybutnerdysunday  
6 minutes ago via Facebook Mobile · Like · Comment

Andrew [redacted]  
Big month for Hip Hop. First up Watch the Throne. Next up Red Album.

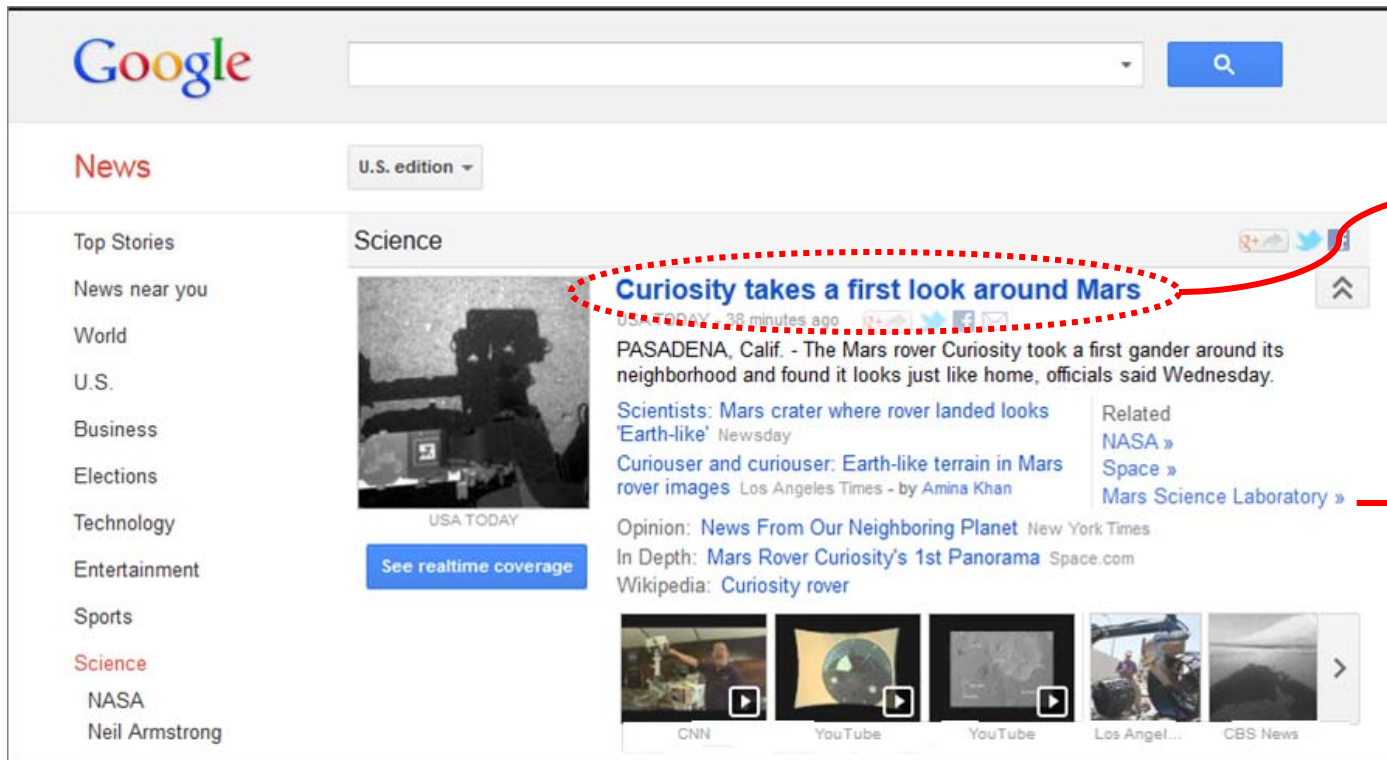
 **Watch the Throne by Jay-Z & Kanye West - Download Watch the Throne on iTunes**  
itunes.apple.com  
Preview and download songs from Watch the Throne by Jay-Z & Kanye West on iTunes. Buy Watch the Throne for just \$11.99.

about an hour ago · 1 · Like · Comment · Share

Jason [redacted]  
The new iTunes volume knob looks like something you'd see on a tablet... I see where you're going Apple...

about an hour ago · 1 · Like · Comment

# Clustering Articles on Google News



Topic cluster

Article Listings

# Clustering Videos on Youtube

YouTube

The Dark Knight Rises - Official Trailer #3 [HD]

WarnerBrosPictures

Subscribe

934 videos

THE FOLLOWING **PREVIEW** HAS BEEN APPROVED FOR **APPROPRIATE AUDIENCES** BY THE MOTION PICTURE ASSOCIATION OF AMERICA, INC.

THE FILM ADVERTISED HAS BEEN RATED

**PG-13** PARENTS STRONGLY CAUTIONED

SOME MATERIAL MAY BE INAPPROPRIATE FOR CHILDREN UNDER 13

INTENSE SEQUENCES OF VIOLENCE AND ACTION, SOME SENSUALITY AND LANGUAGE

www.filmratings.com www.mpa.org

24,883,922

100,177 likes, 2,304 dislikes

Published on Apr 30, 2012 by WarnerBrosPictures

<http://www.thedarkknightriserises.com/>

<http://www.facebook.com/thedarkknightriserises>

"The Dark Knight Rises" In theaters July 20.

Warner Bros. Pictures® and Legendary Pictures® "The Dark Knight Rises" is the epic conclusion to filmmaker Christopher Nolan's Batman trilogy.

The Dark Knight Rises - Official Trailer #2 [HD] by WarnerBrosPictures 2,655,480 views 2:13

The Dark Knight Rises - Official Trailer #4 [HD] by WarnerBrosPictures 1,402,729 views 2:21

MAGIC MIKE - OFFICIAL TRAILER [HD] by WarnerBrosPictures 3,662,666 views 2:33

Dark Shadows - Vampire History by WarnerBrosPictures 242,352 views 3:01

THE ROCK (1996) FULL MOVIE HD by MrNazier100 241,843 views 2:16-21

I AM LEGEND ALTERNATE ENDING by renjensel 853,155 views 7:12

ULTIMATE TRILOGY The Dark Knight Rises Ultimate Trilogy Trailer

- Keywords
- Popularity
- Viewer engagement
- User browsing history

# Clustering for Efficient Image retrieval

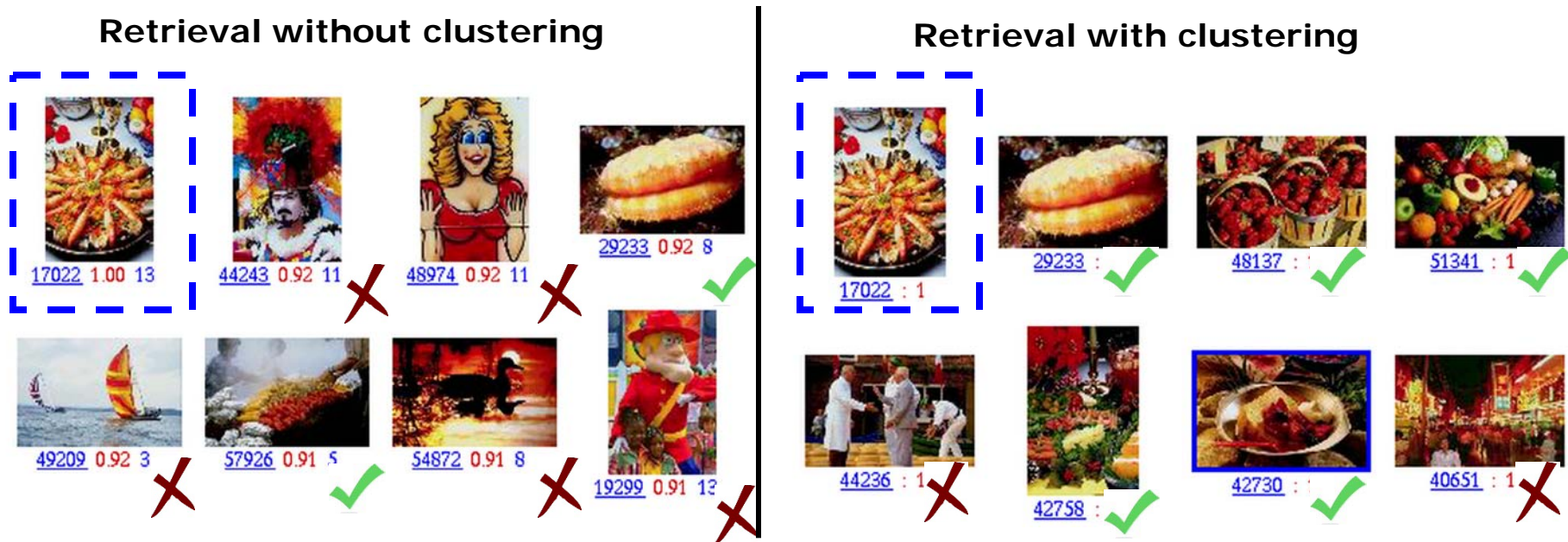


Fig. 1. Upper-left image is the query. Numbers under the images on left side: image ID and cluster ID; on the right side: Image ID, matching score, number of regions.

Retrieval accuracy for the “food” category (average precision):

Without clustering: **47%**

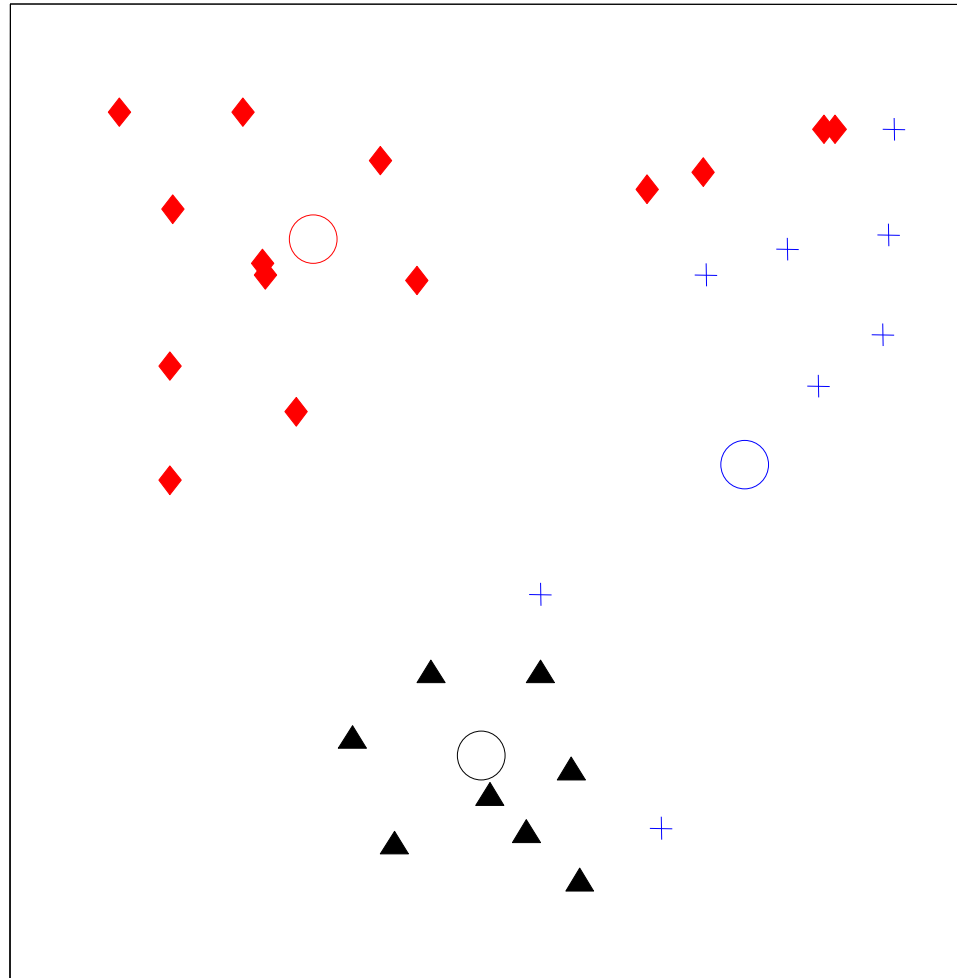
With clustering: **61%**

# Clustering Algorithms

Hundreds of clustering algorithms are available; many are “admissible”, **but no algorithm is “optimal”**

- K-means
- Gaussian mixture models
- Kernel K-means
- Spectral Clustering
- Nearest neighbor
- Latent Dirichlet Allocation

# K-means Algorithm

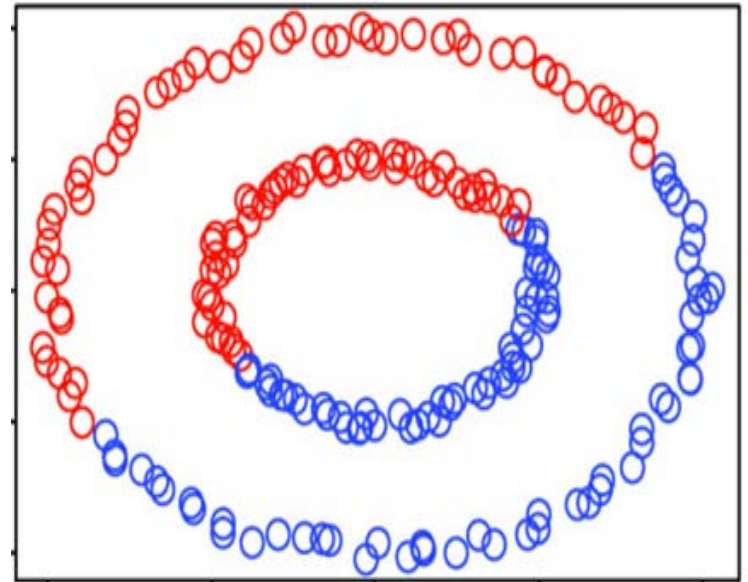
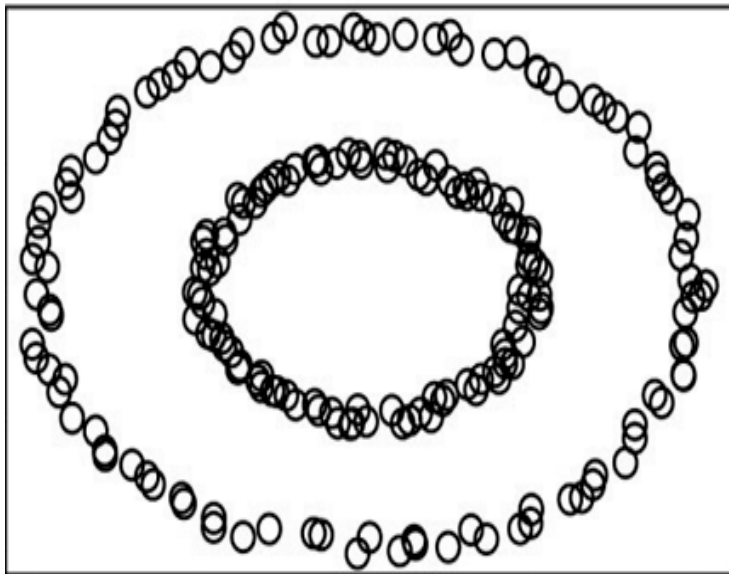


Repeat until points no longer change or until the centroids stop moving

# K-means: Limitations

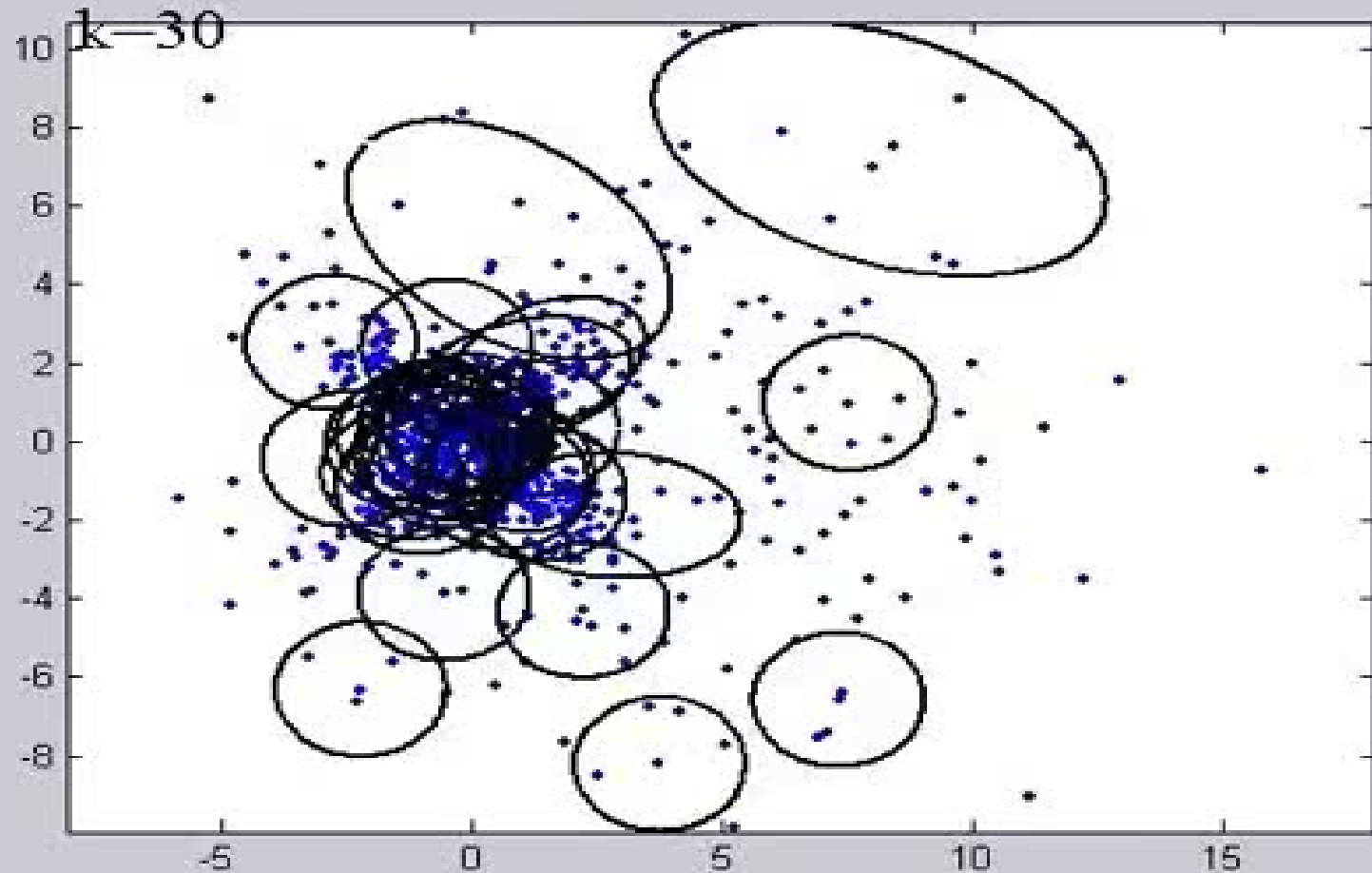
Prefers “compact” and “isolated” clusters

$$\min \sum_{i=1}^n \sum_{k=1}^K u_{ik} \|x_i - c_k\|^2$$





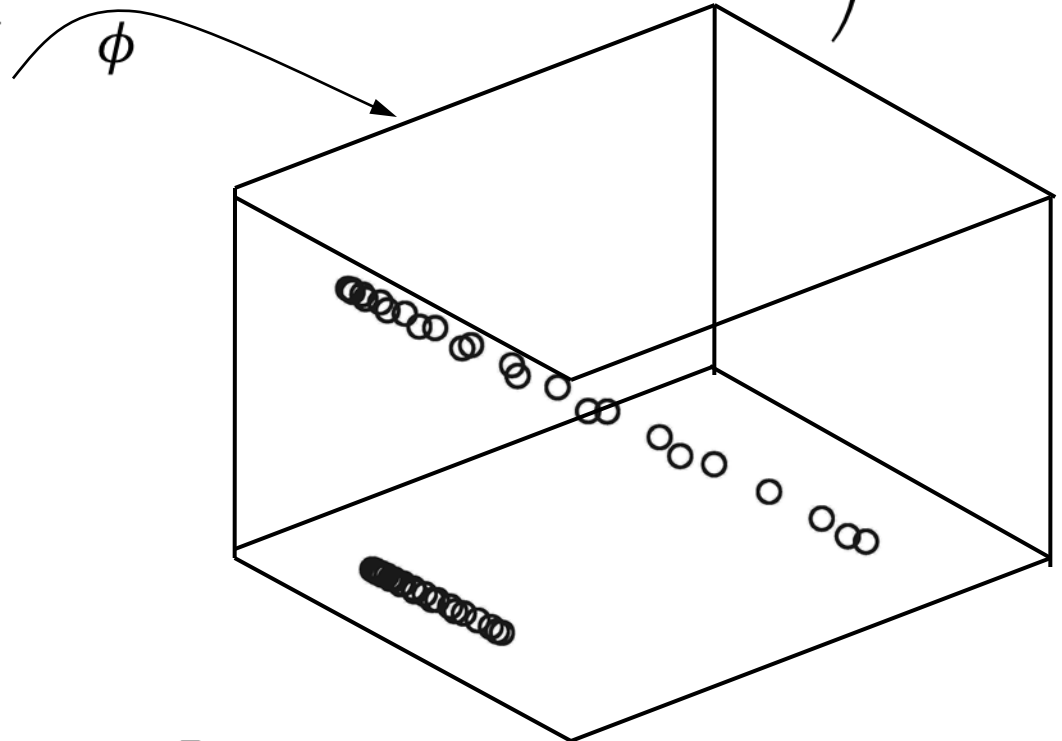
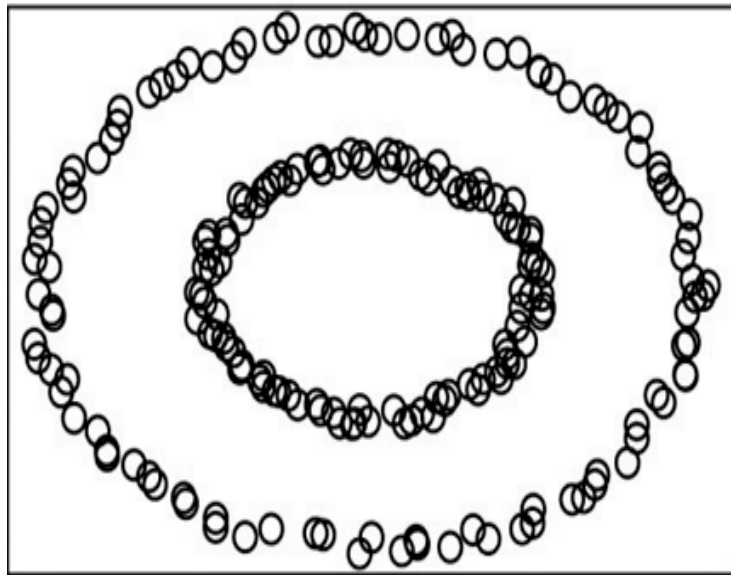
# Gaussian Mixture Model



# Kernel K-means

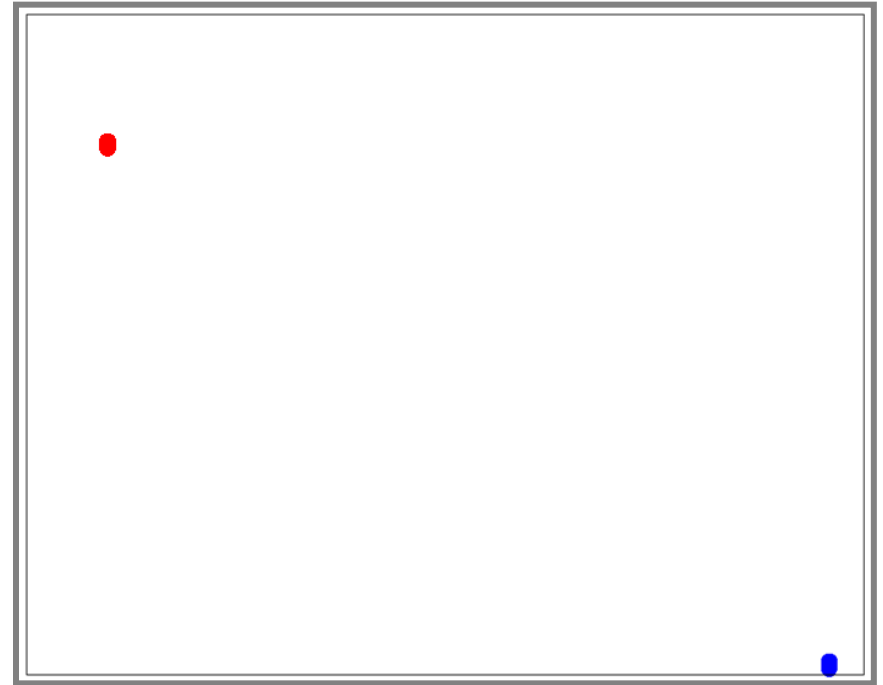
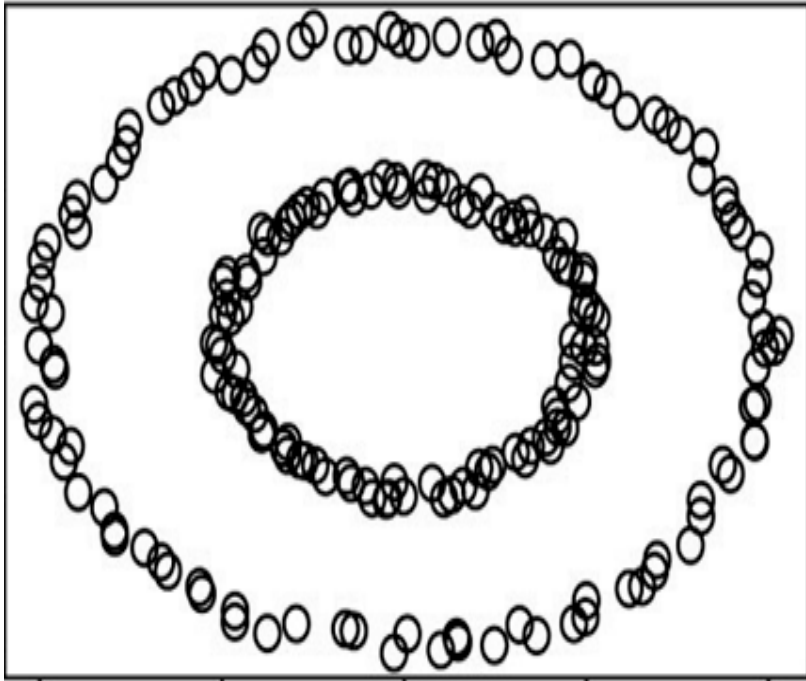
**Non-linear** mapping to find clusters of arbitrary shapes

$$\min \text{Trace} \left( \sum_{i=1}^n \sum_{k=1}^K u_{ik} (\phi(x_i) - c_k^\phi) (\phi(x_i) - c_k^\phi)^T \right)$$



$\phi(x, y) = (x^2, \sqrt{2}xy, y^2)^T$   $K(a, b) = \phi(a)^T \phi(b)$  Polynomial kernel representation

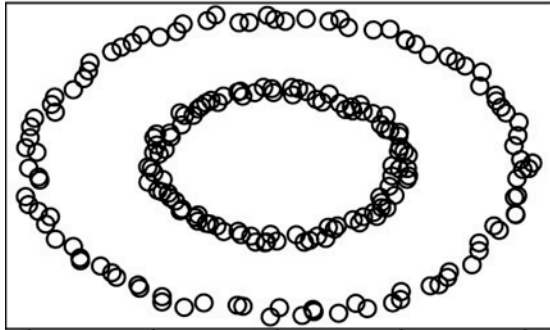
# Spectral Clustering



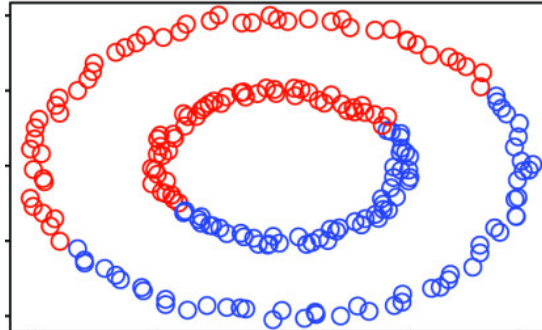
Represent data using the top  $K$  eigenvectors of the kernel matrix; **equivalent to Kernel K-means**

# K-means vs. Kernel K-means

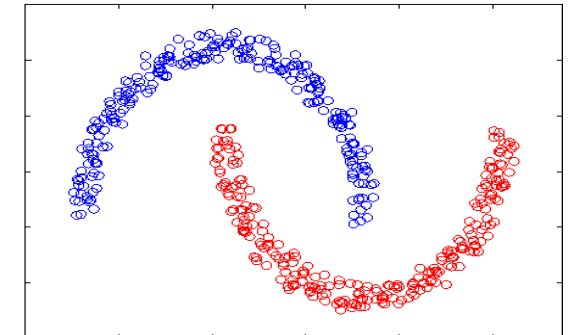
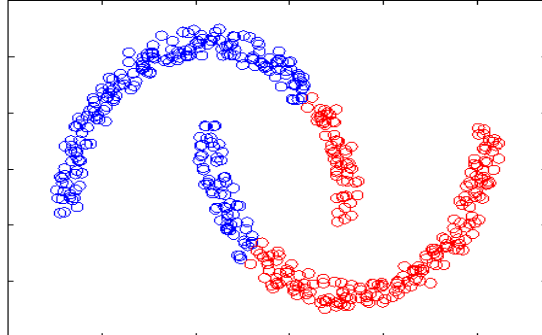
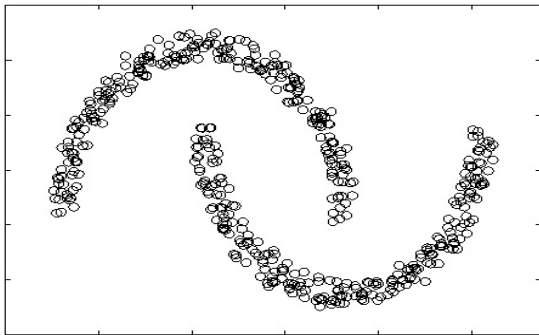
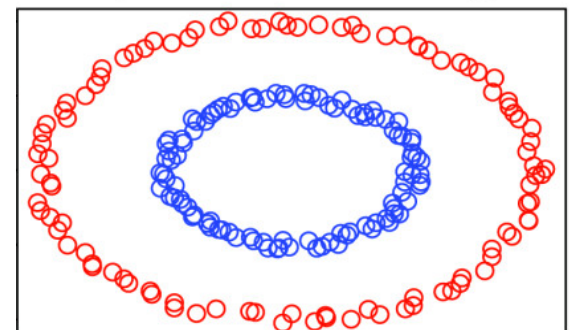
Data



K-means



Kernel K-means



Kernel clustering is able to find “complex” clusters

How to choose the right kernel? RBF kernel is the default

# Kernel K-means is Expensive

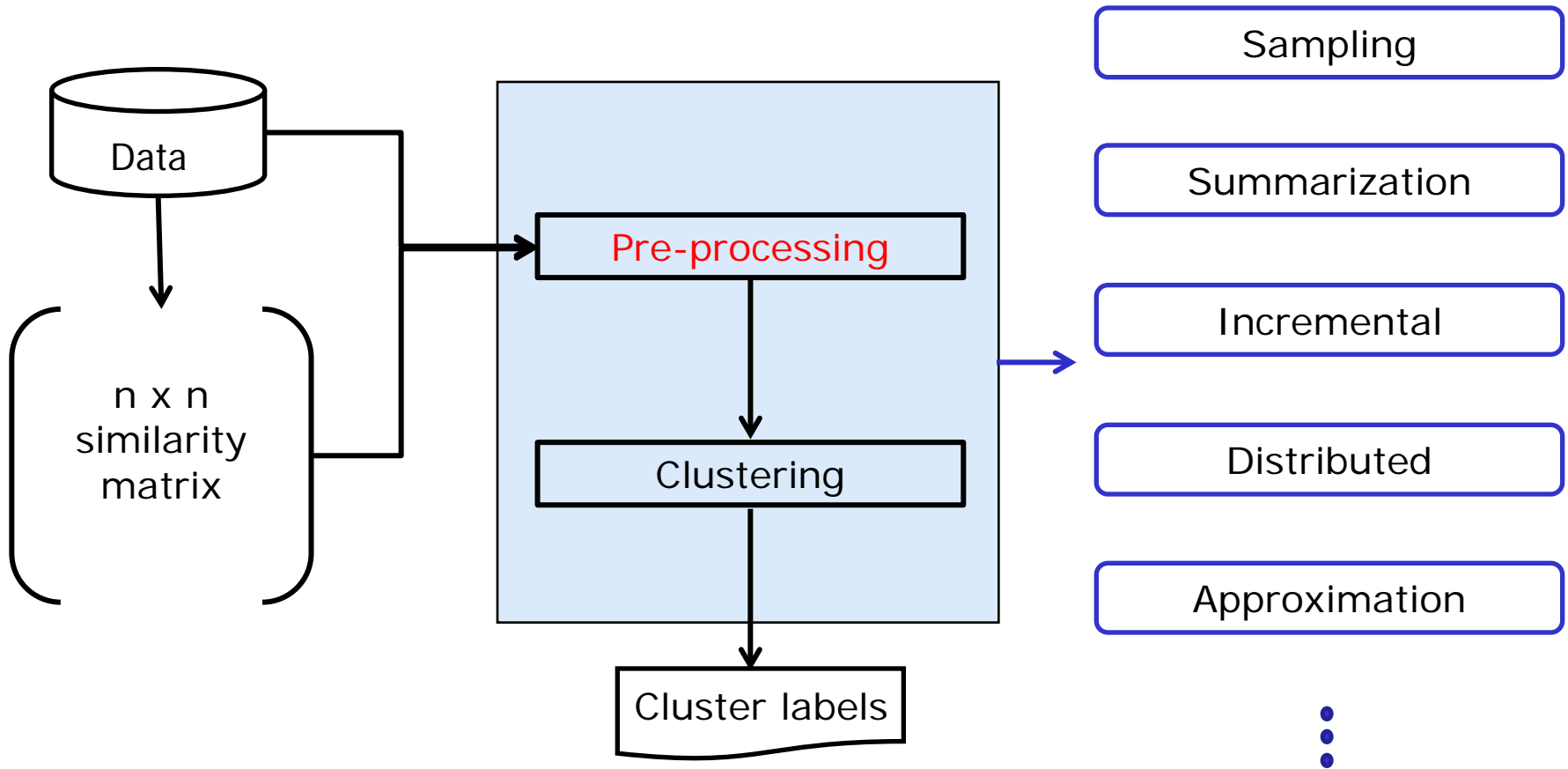
No. of Objects (n)	No. of operations	
	K-means	Kernel K-means
	$O(nKd)$	$O(n^2K)$
1M	$10^{13}$ (6412*)	$10^{16}$
10M	$10^{14}$	$10^{18}$
100M	$10^{15}$	$10^{20}$
1B	$10^{16}$	$10^{22}$

$d = 10,000$ ;  $K=10$

\* Runtime in seconds on Intel Xeon 2.8 GHz processor using 40 GB memory

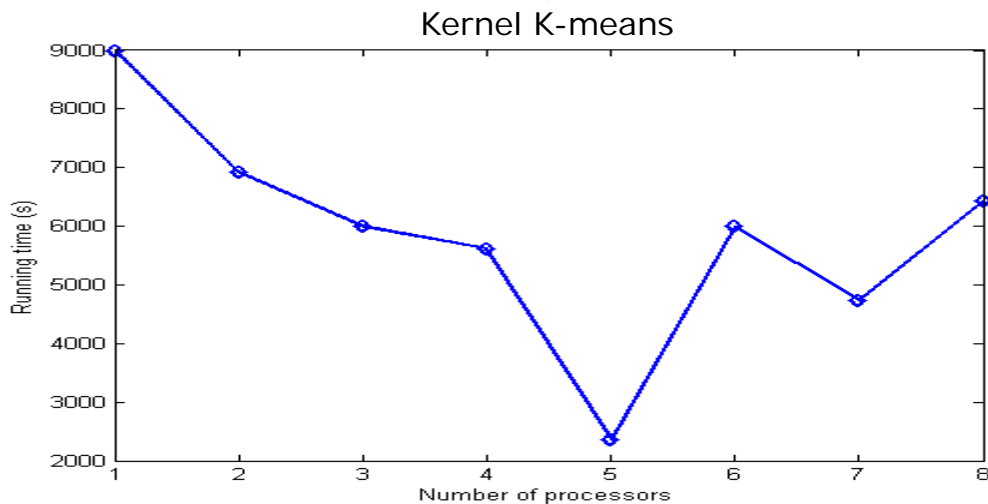
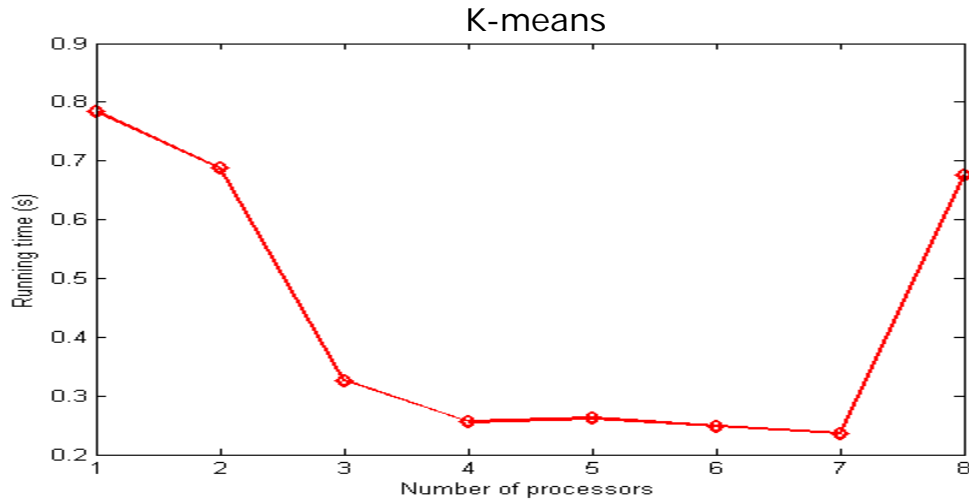
A petascale supercomputer (IBM Sequoia, June 2012) with ~1 exabyte memory is needed to run kernel K-means on 1 billion points!

# Clustering Big Data



# Distributed Clustering

Clustering 100,000 2-D points with 2 clusters on 2.3 GHz quad-core Intel Xeon processors, with 8GB memory in intel07 cluster

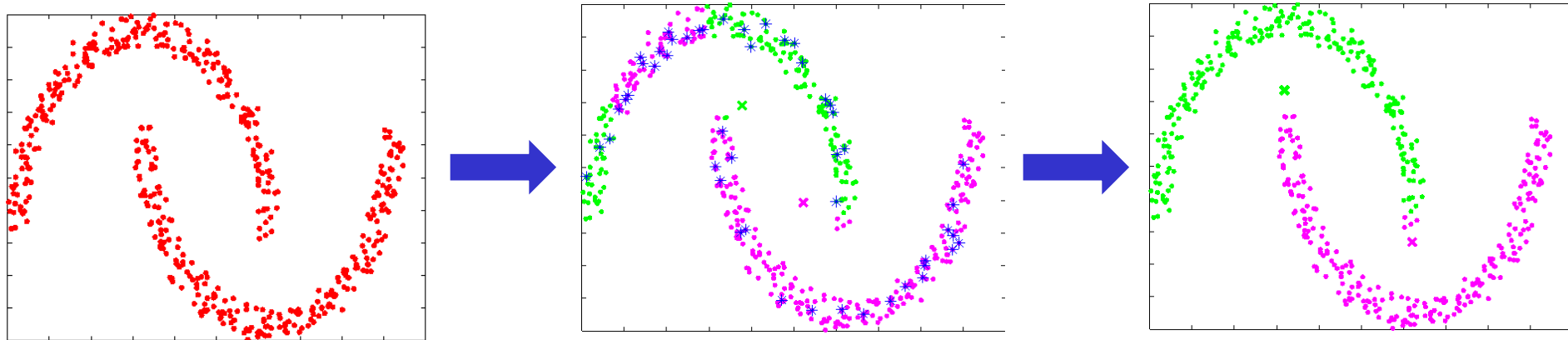


Number of processors	Speedup	
	K-means	Kernel K-means
2	1.1	1.3
3	2.4	1.5
4	3.1	1.6
5	3.0	3.8
6	3.1	1.9
7	3.3	1.5
8	1.2	1.5

Network communication cost increases with the no. of processors

# Approximate kernel K-means

Tradeoff between clustering accuracy and running time



Randomly sample  $m$  points  $\{y_1, y_2, \dots, y_m\}$  in the input space

compute the kernel similarity matrices  $K_A$  ( $m \times m$ ) and  $K_B$  ( $n \times m$ )

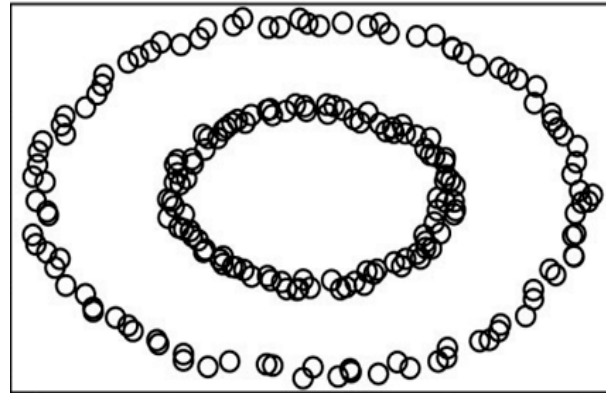
$$\min_{\{u_{ik}\}} \max_{\{a_{jk}\}} \sum_{k=1}^m \sum_{i=1}^n u_{ik} \left\| \phi(x_i) - \sum_{j=1}^m a_{jk} \phi(y_j) \right\|^2$$

(Linear runtime and memory complexity)

(equivalent to running K-means on  $K_B K_A^{-1} K_B^T$ )



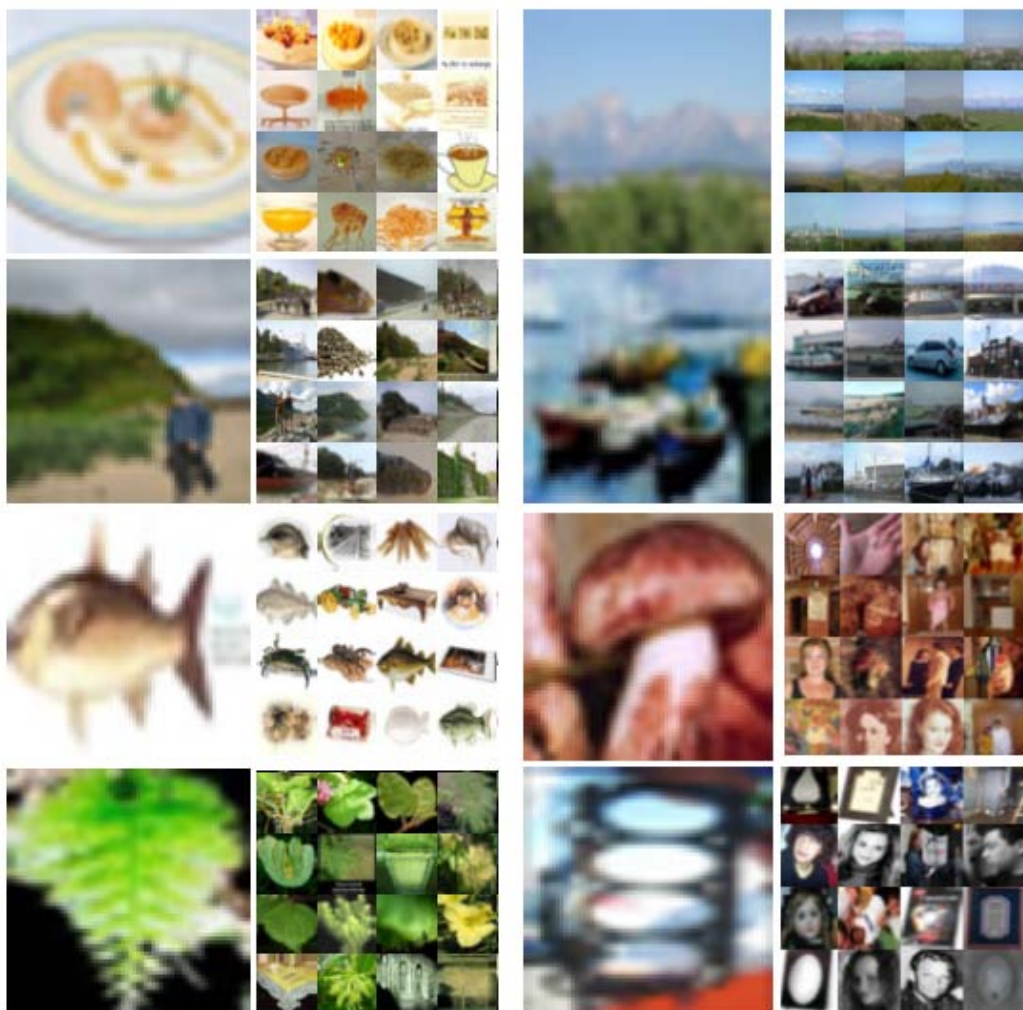
# Approximate Kernel K-Means



No. of objects (n)	Running time (seconds)			Clustering accuracy (%)		
	Kernel K-means	Approximate kernel K-means (m=100)	K-means	Kernel K-means	Approximate kernel K-means (m=100)	K-means
10K	3.09	0.20	0.03	100	93.8	50.1
100K	320.10	1.18	0.17	100	93.7	49.9
1M	-	15.06	0.72	-	95.1	50.0
10M	-	234.49	12.14	-	91.6	50.0

2.8 GHz processor, 40 GB

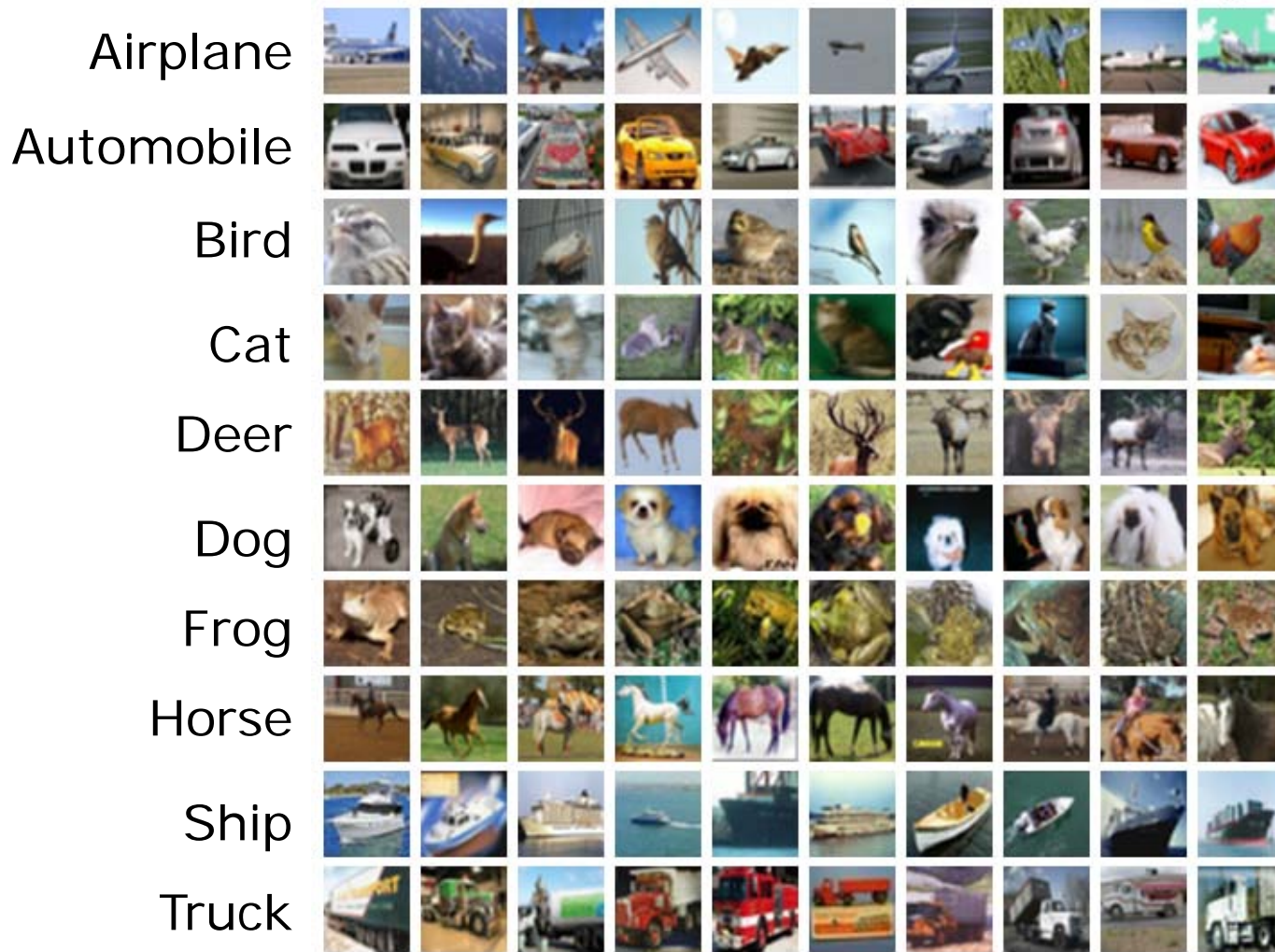
# Tiny Image Data set



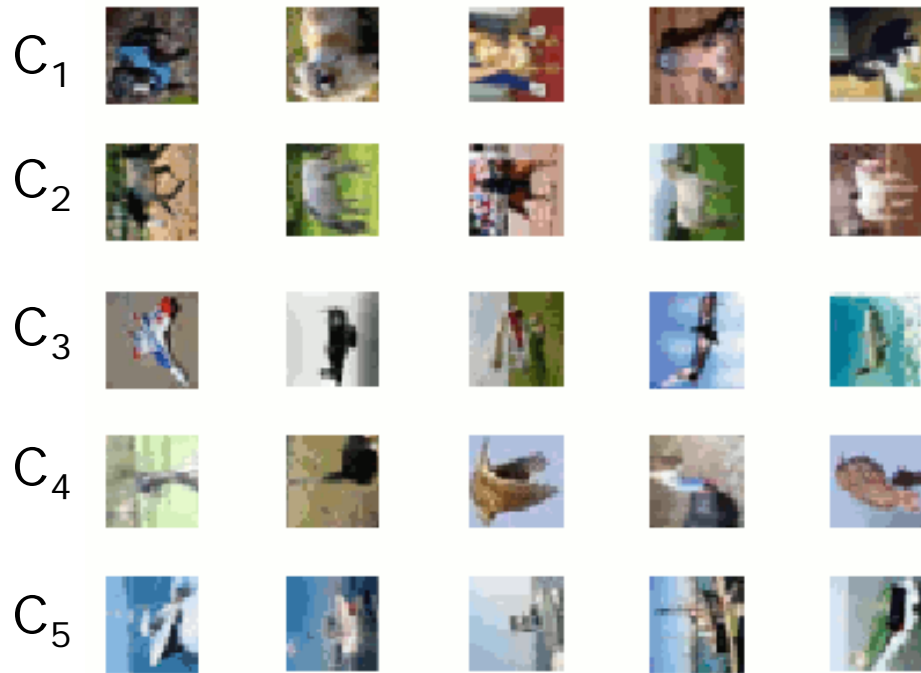
~80 million 32x32 images from ~75K classes (bamboo, fish, mushroom, leaf, mountain,...); image represented by 384-dim. GIST descriptors

# Tiny Image Data set

10-class subset (CIFAR-10): 60K manually annotated images



# Clustering Tiny Images



Example Clusters

Average clustering time (100 clusters)	
Approximate kernel K-means (m=1,000)	8.5 hours
K-means	6 hours

2.3GHz, 150GB memory

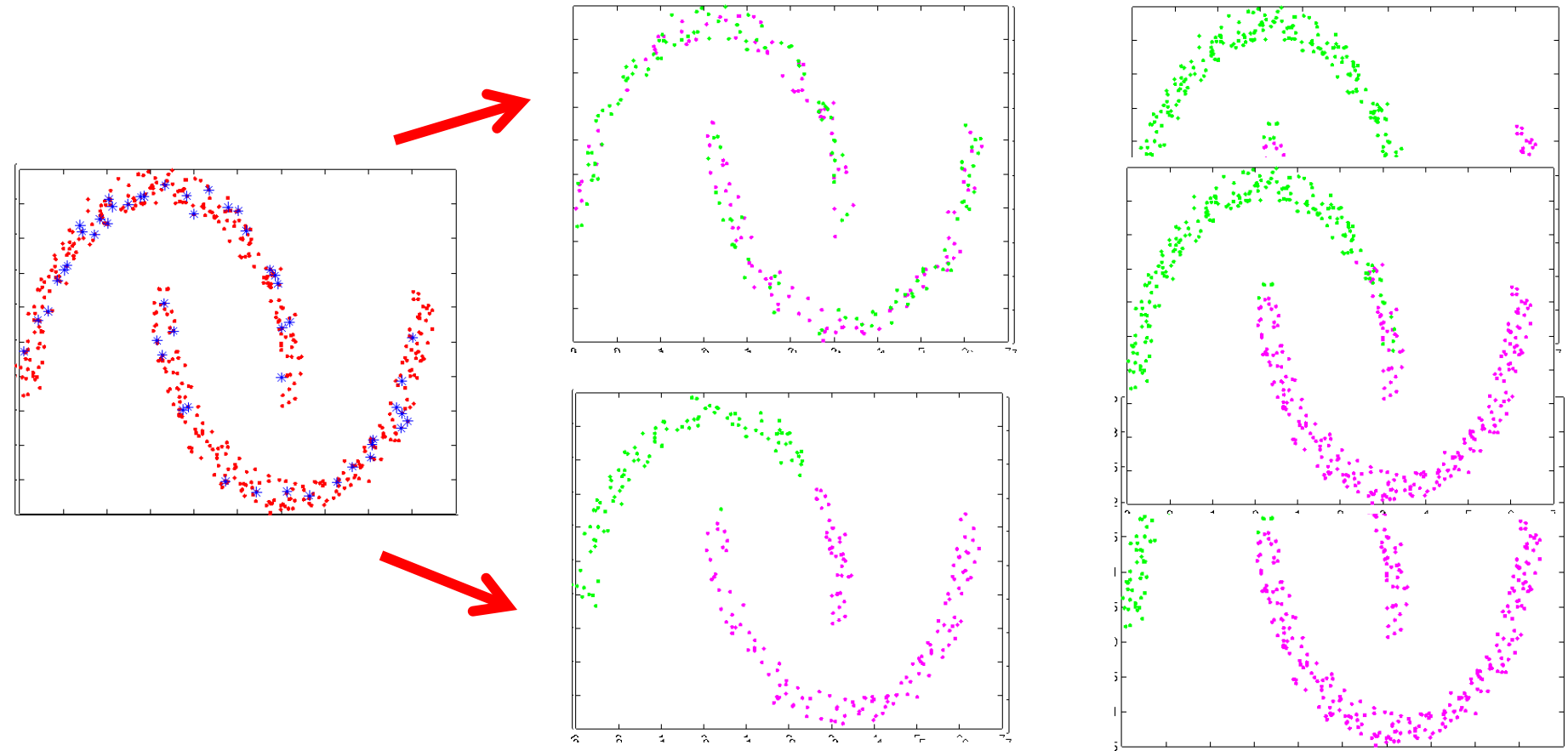
# Clustering Tiny Images

Best Supervised Classification Accuracy on CIFAR-10: 54.7%

Clustering accuracy	
Kernel K-means	29.94%
Approximate kernel K-means (m = 5,000)	29.76%
Spectral clustering	27.09%
K-means	26.70%

# Distributed Approx. Kernel K-means

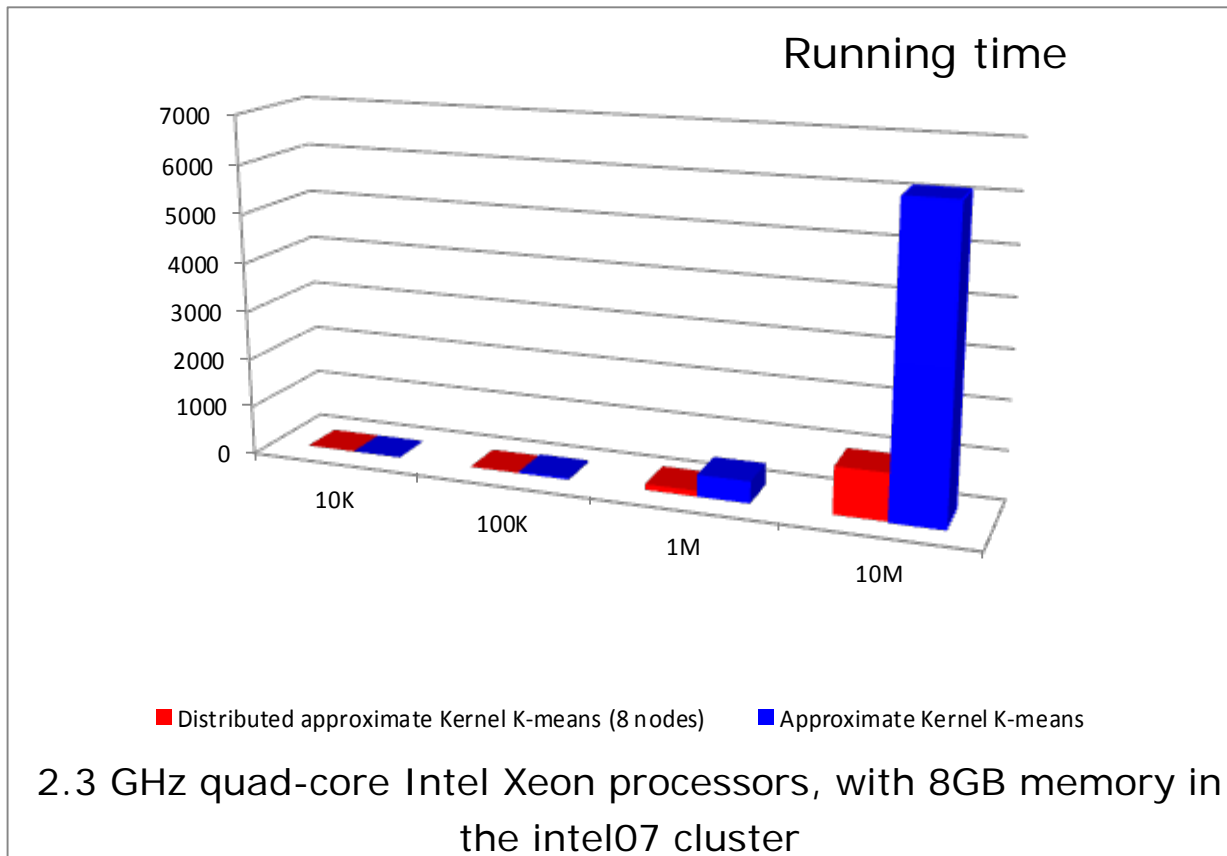
For better scalability and faster clustering



Splitting the data into  $t$  partitions (e.g.,  $t=2$ ) and performing K-means on each partition in parallel. The final result is obtained by aggregating the results from all partitions.

# Distributed Approximate kernel K-means

2-D data set with 2 concentric circles

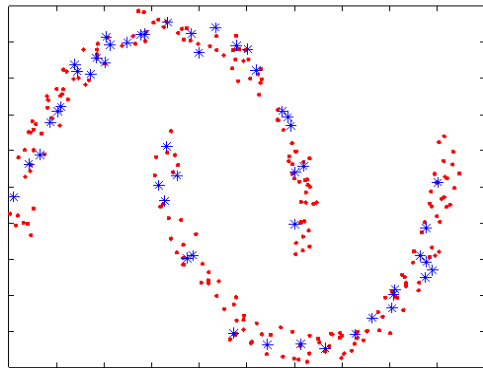
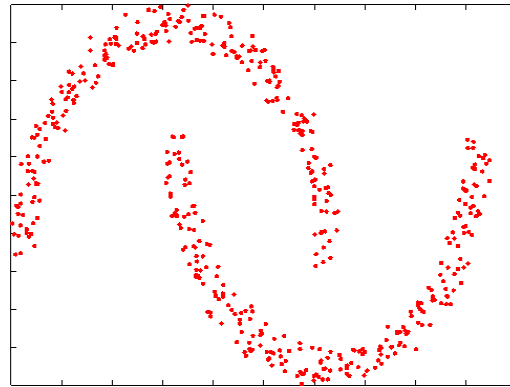


Size of data set	Speedup
10K	3.8
100K	4.8
1M	3.8
10M	6.4

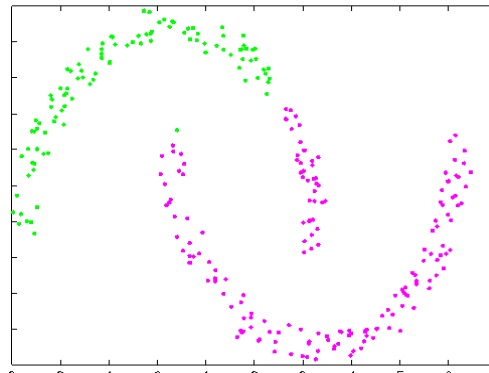
# Limitations of Approx. kernel K-means

Clustering data with more than 10 million points will require terabytes of memory!

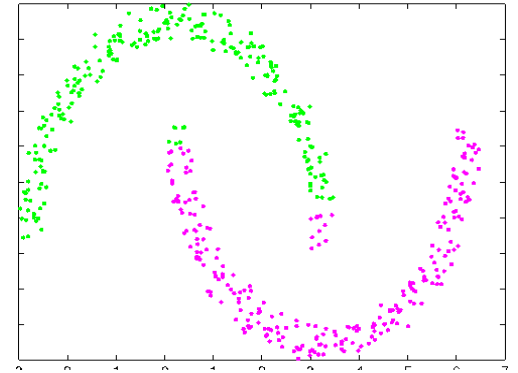
*Sample and Cluster Algorithm (SnC)*



Sample  $s$  points from data



Run approximate kernel K-means on the  $s$  points

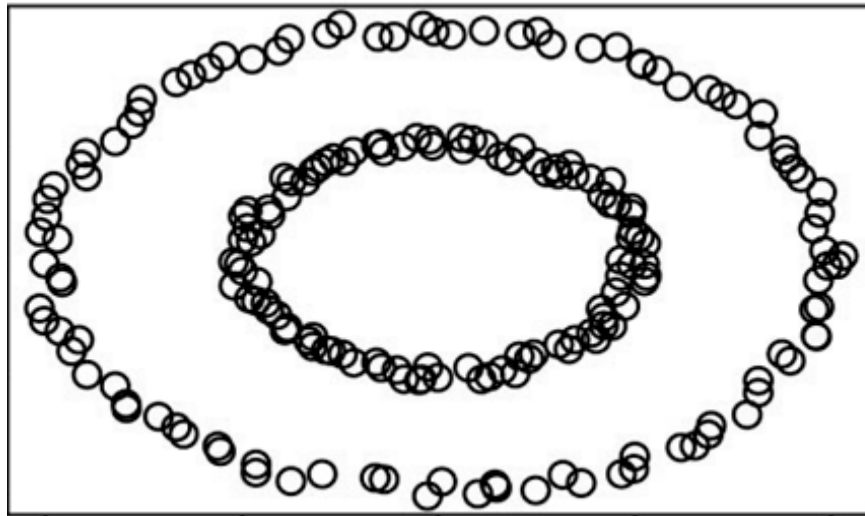


Assign remaining points to the nearest cluster center



# Clustering one billion points

Sample and Cluster ( $s = 1$  million,  $m = 100$ )



Running time			Average Clustering Accuracy	
K-means	SnC	SnC –distributed (8 cores)	K-means	SnC
53 minutes	1.2 hours	45 minutes	50%	85%

# Clustering billions of points

- Work in progress
  - Application to real data sets
  - Yahoo! AltaVista Web Page Hyperlink Connectivity Graph (2002) containing URLs and hyperlinks for over 1.4 billion public web pages
- Challenges
  - Graph Sparsity: Reduce the dimensionality using random projection, PCA
  - Cluster Evaluation: No ground truth available, internal measures such as link density of clusters

# Summary

- Clustering is an exploratory technique; used in every scientific field that collects data
- Choice of clustering algorithm & its parameters is data dependent
- Clustering is essential for “Big Data” problem
- Approximate kernel K-means provides good tradeoff between scalability & clustering accuracy
- Challenges: Scalability, very large no. of clusters, heterogeneous data, streaming data, validity

# Big Data

