# 基于Hadoop的
# SNS统计平台和聚类推荐

白伯纯 张叶银

renren.com

# 人人网

- 2.2亿用户
- 平均190好友
- 月40亿照片访问

- 一成付费用户
- 五成用户每天使用
- 八成有真实的资料
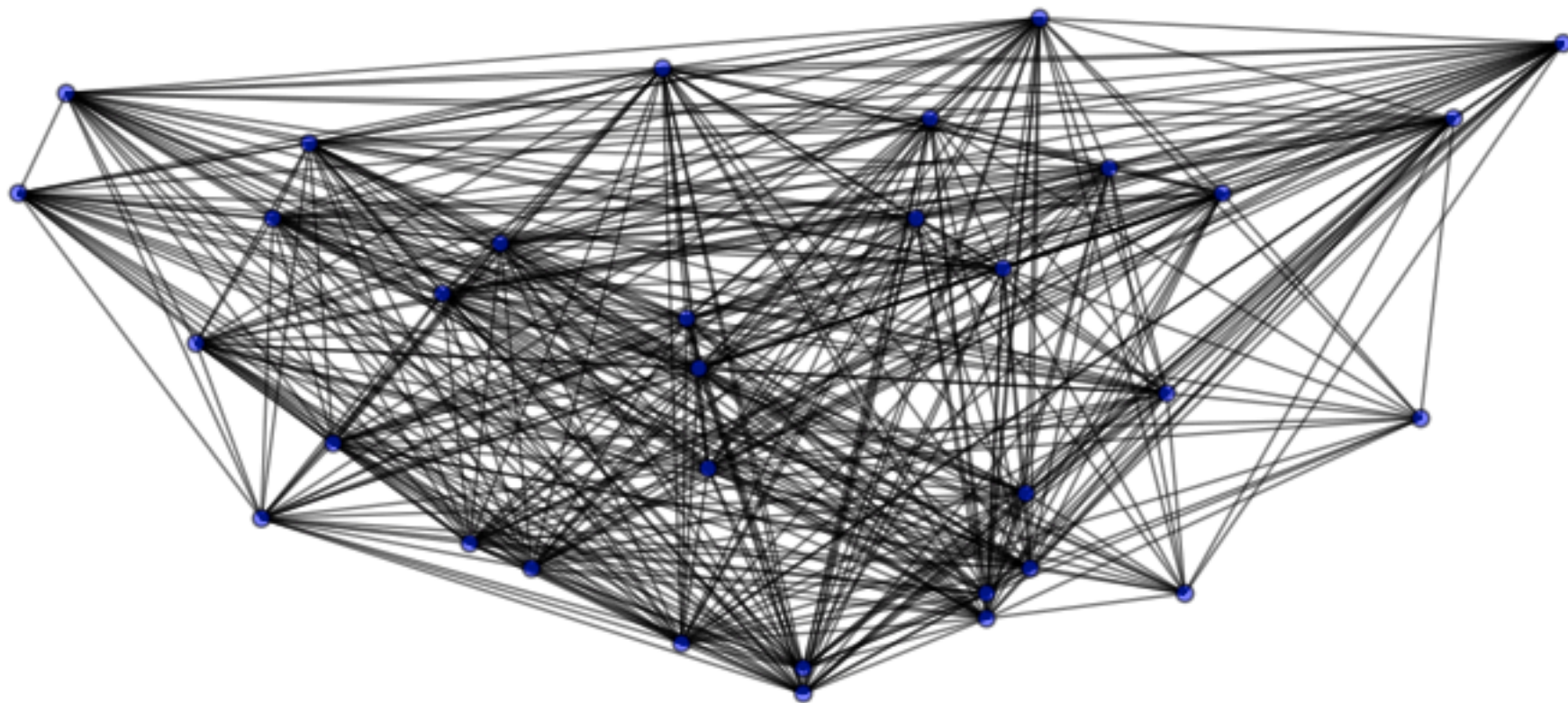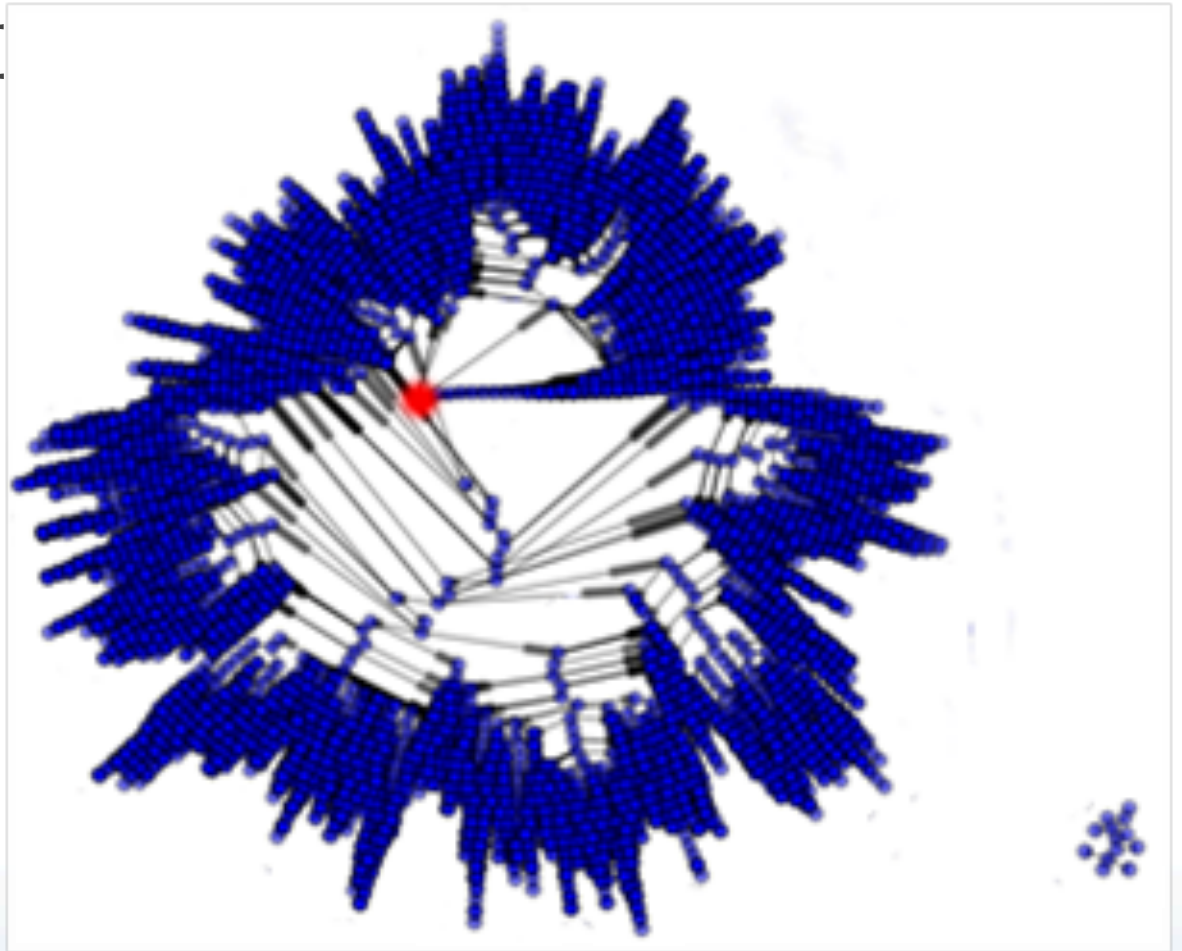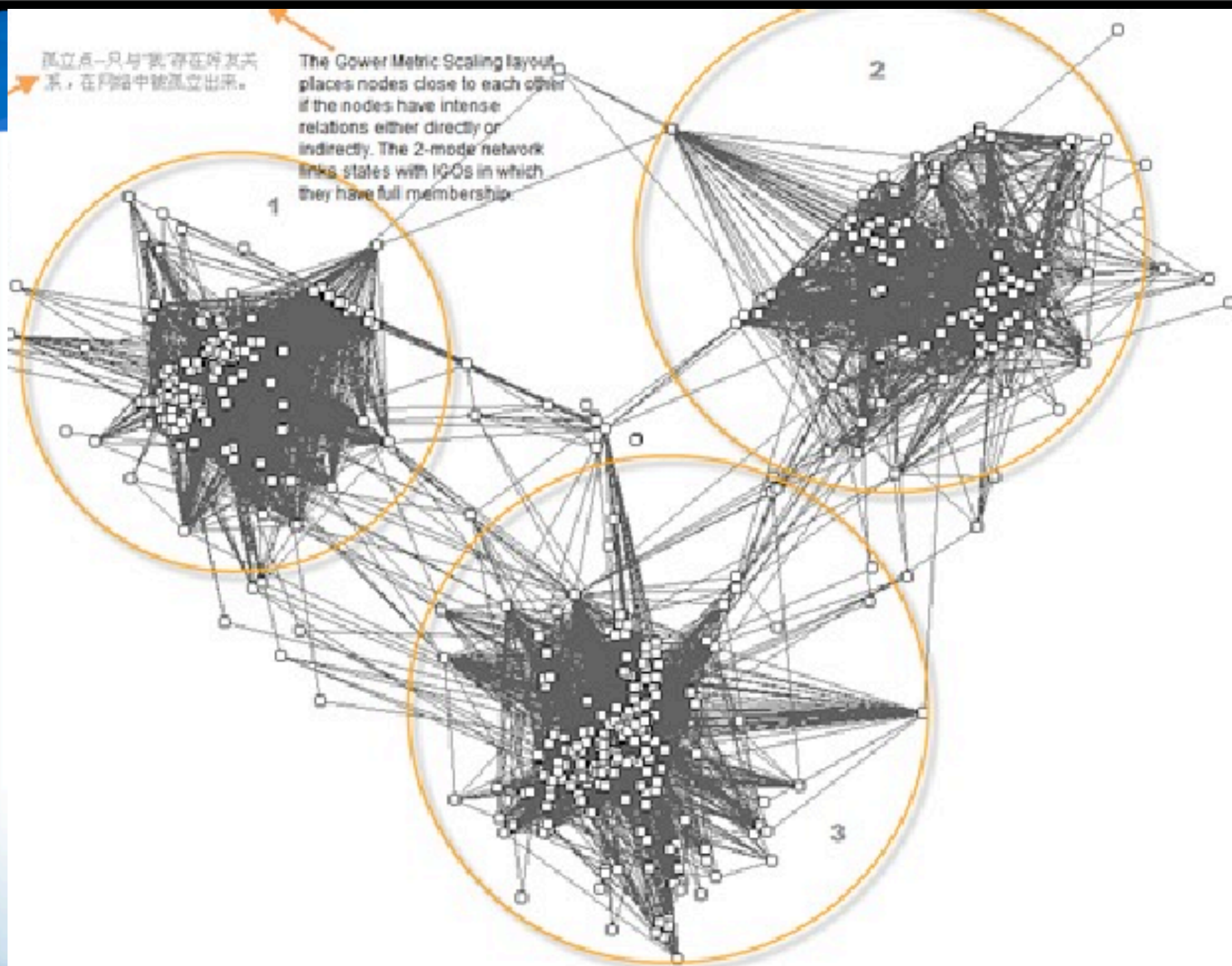
人人网 RENREN.com

# 机遇

- 唯一标识

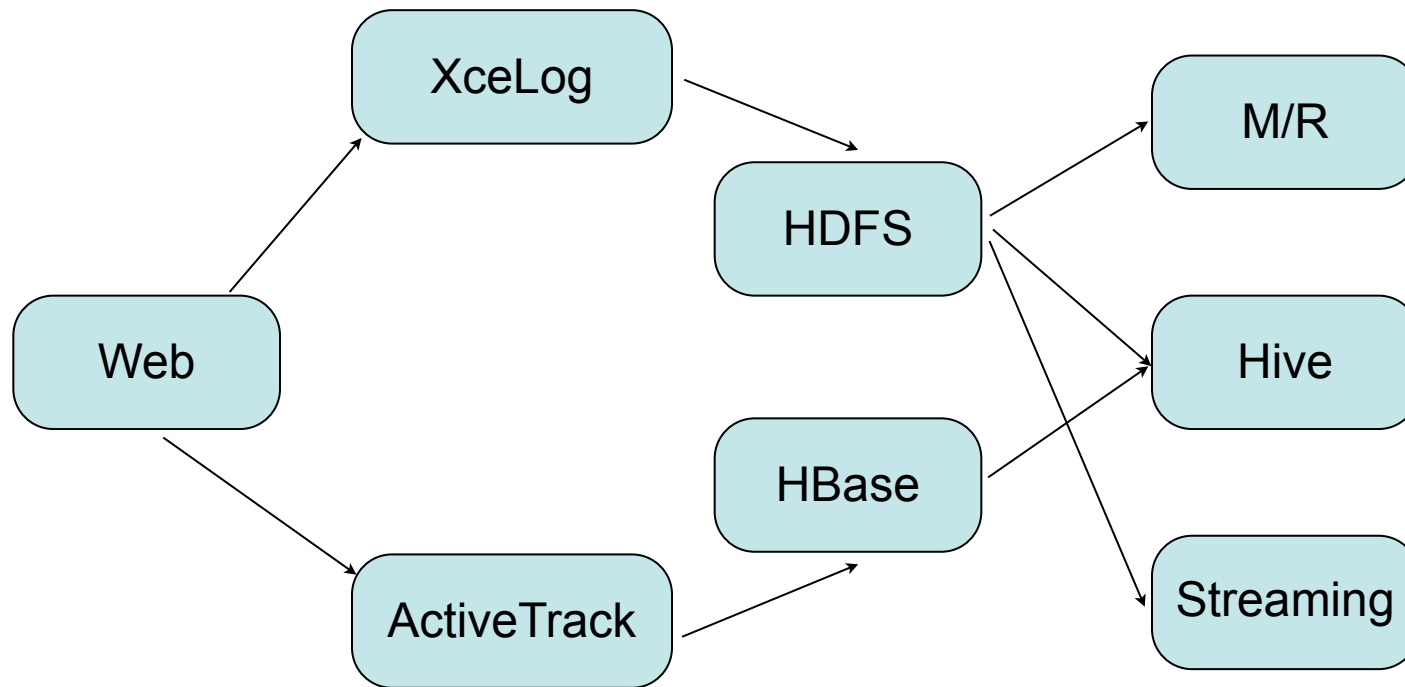# 机遇

- 唯一标识

# 机遇

- 结构化数据

- 高复杂度计算

The Gower Metric Scaling layout places nodes close to each other if the nodes have intense relations either directly or indirectly. The 2-mode network links states with IGOs in which they have full membership.
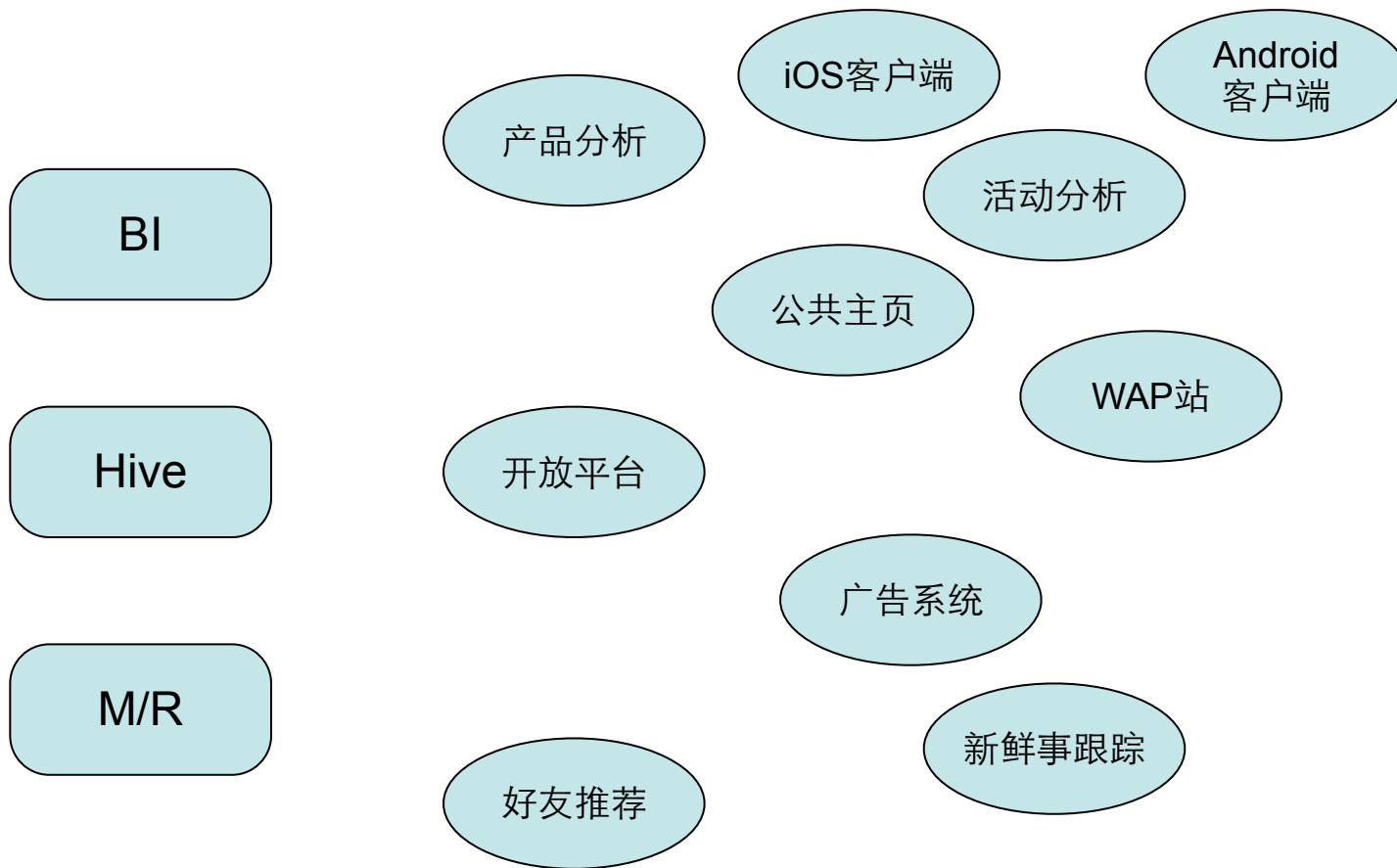
孤立点—只与"孤"存在时发关系，在网络中能孤立出来。

# 部署

- 200台
  - Hadoop 0.21.0
  - 4k+任务/天
  - 700TB Used/1.2PB Total
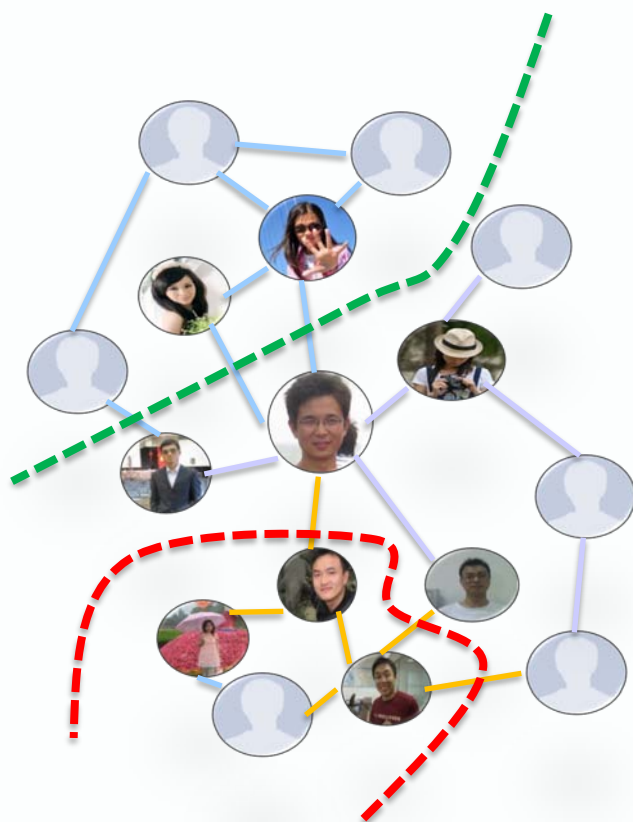  - Hive/HBase/Streaming
- 30台
  - Hadoop 0.20.3
  - HBase only

# 体系结构

BI

Hive

M/R

产品分析

iOS客户端

Android 客户端

活动分析

公共主页

WAP站

开放平台

广告系统

新鲜事跟踪

好友推荐

# Social computing at Renren
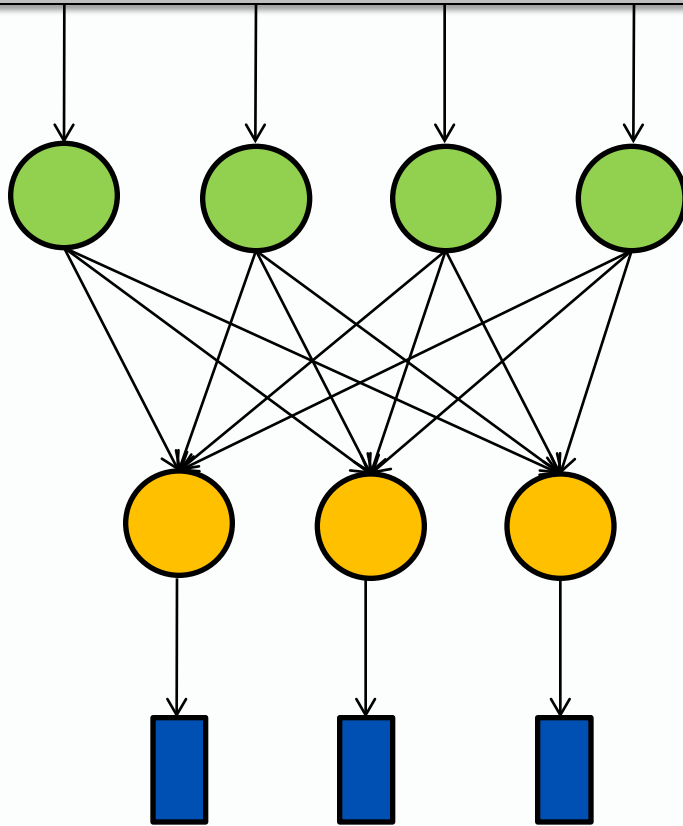


## Data driven applications

# Distribute

- TB of daily log data to be analyzed
- Millions of blogs, videos to be recommended
- Hundreds of millions of friends to be recommended

The most Computational-Intensive applications
with highly structured big data

# MapReduce

Data Points
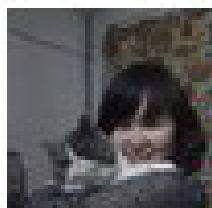
Adjacent list data structure
- A sparse representation facilitate to pass graph structures from iteration to the next iteration

- Parallel Breadth-first search
  - To find shortest path in the graph

- Google page rank
  - Impact index passing through graph links

- Distributed k-means clustering
  - Clustering large data into pre-defined number of groups

# Case I: Friend recommendation by agglomerative hierarchical clustering

- Primary problem of friend recommendation
  - User familiarity
    - Common friends
    - User profile
    - User access
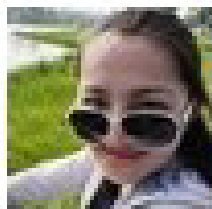    - User interest

# People you may know

- Friends' friends



$$similarity(user1, user2) =$$
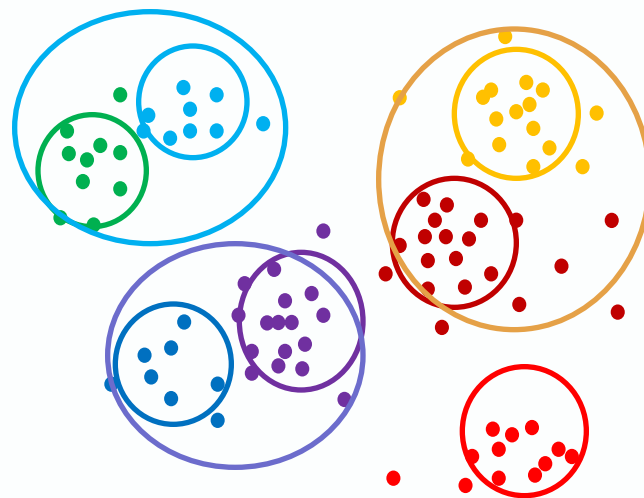$$|friendset1 \cap friendset2|$$

# Hierarchy

- Clustering to find communities in social network
  - All in one community share some properties.
  - These overlapping communities reveal some social relationship of different levels.
  - They help to building new friendships in the social network.
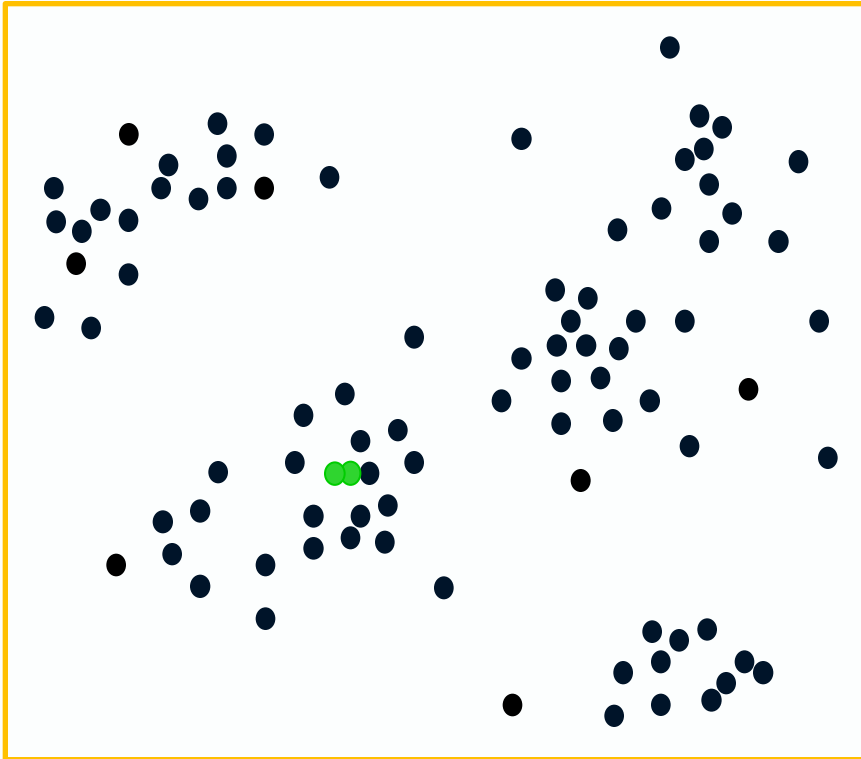
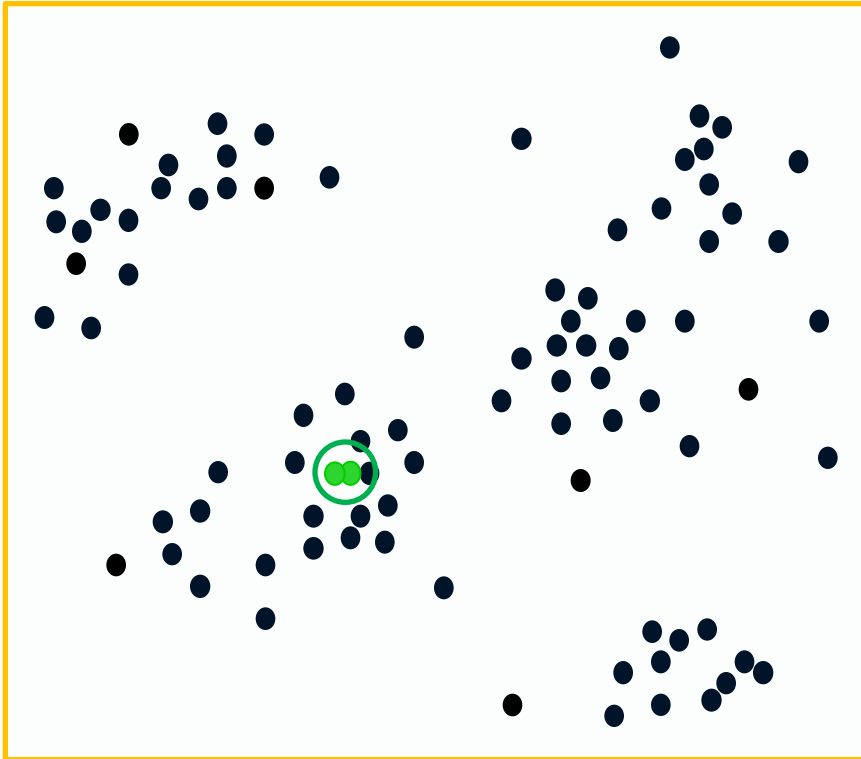# Clustering: unsupervised learning
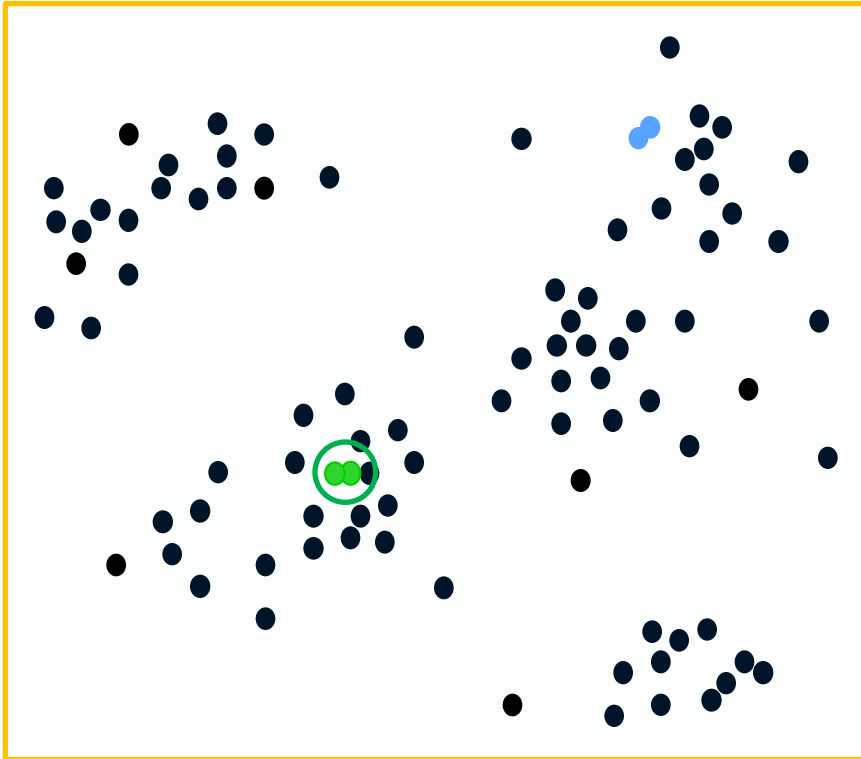


flat clustering

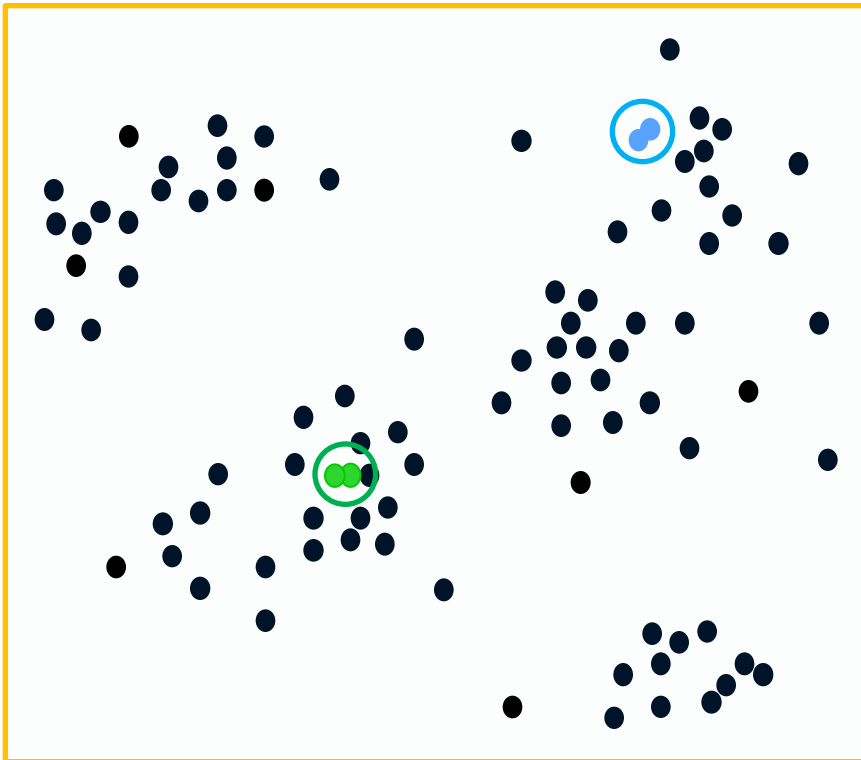hierarchical clustering

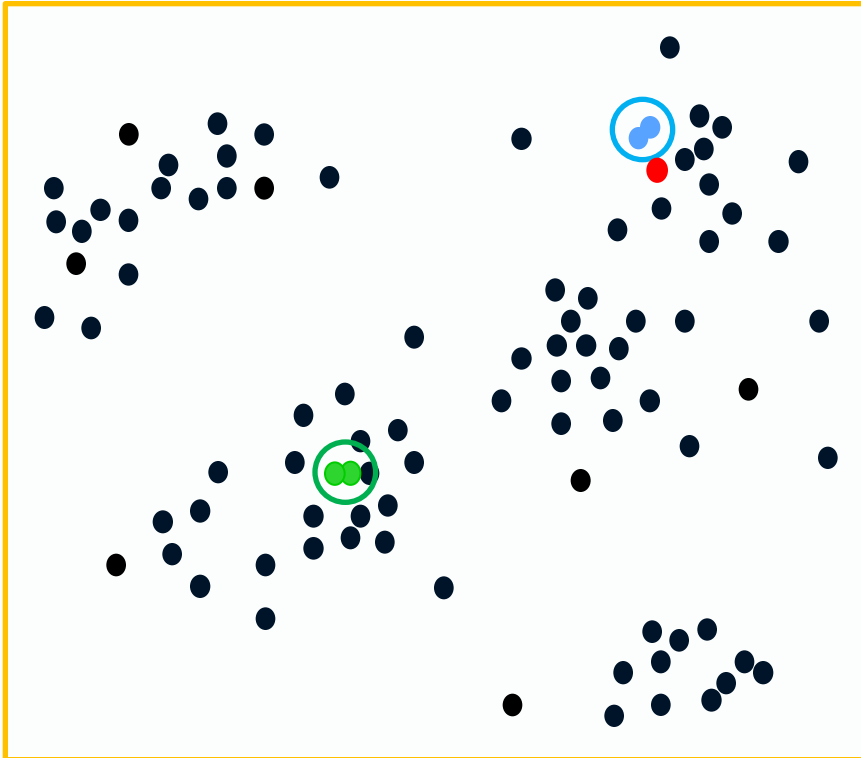# Hierarchical clustering

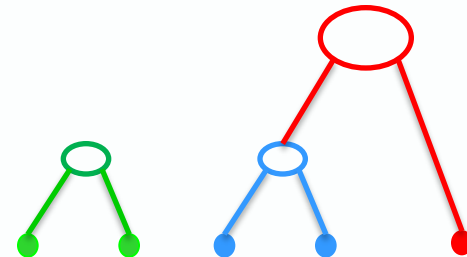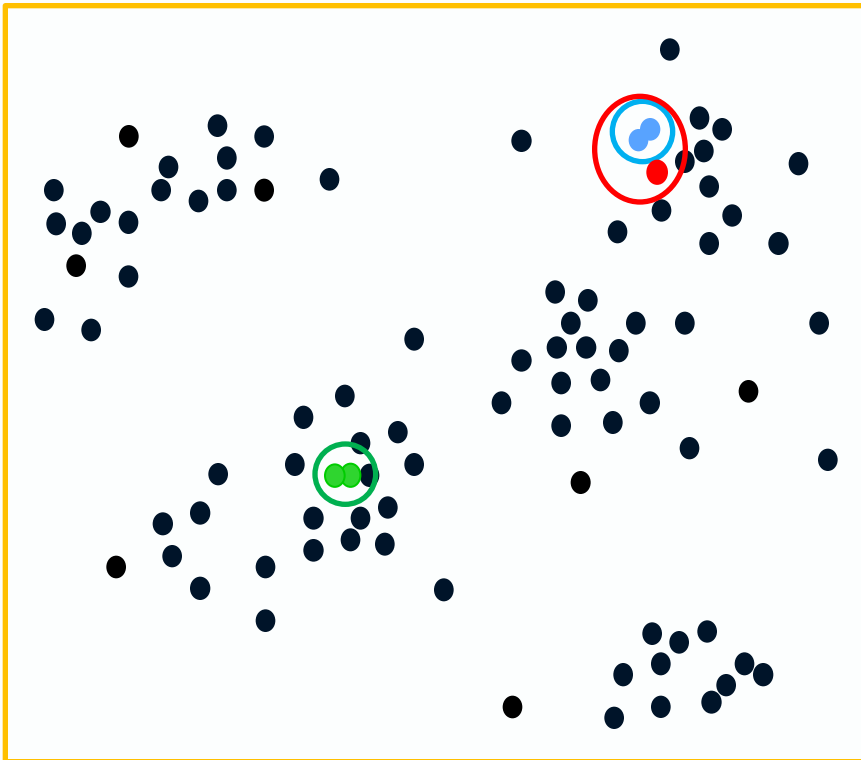# Hierarchical clustering

# Hierarchical clustering

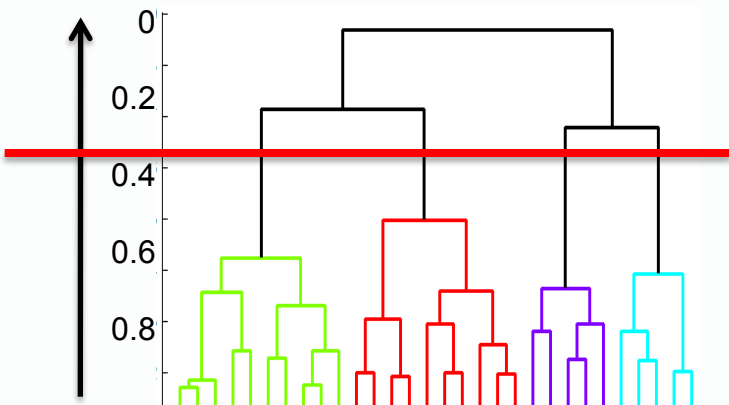# Hierarchical clustering

# Hierarchical clustering

# Hierarchical clustering

# Hierarchical agglomerative clustering

Monotonic



*Method:* **Merge the nearest clusters until a single cluster is left**

*Procedure HAC* (N points, stop criterion)
{
(1) Initialize n points as n cluster centers;
(2) Iterate over centers until stop criterion is satisfied:

      a. Compute pair-wise similarity between any two centers $sim(c_i, c_j)$

      b. Find the nearest pair of centers
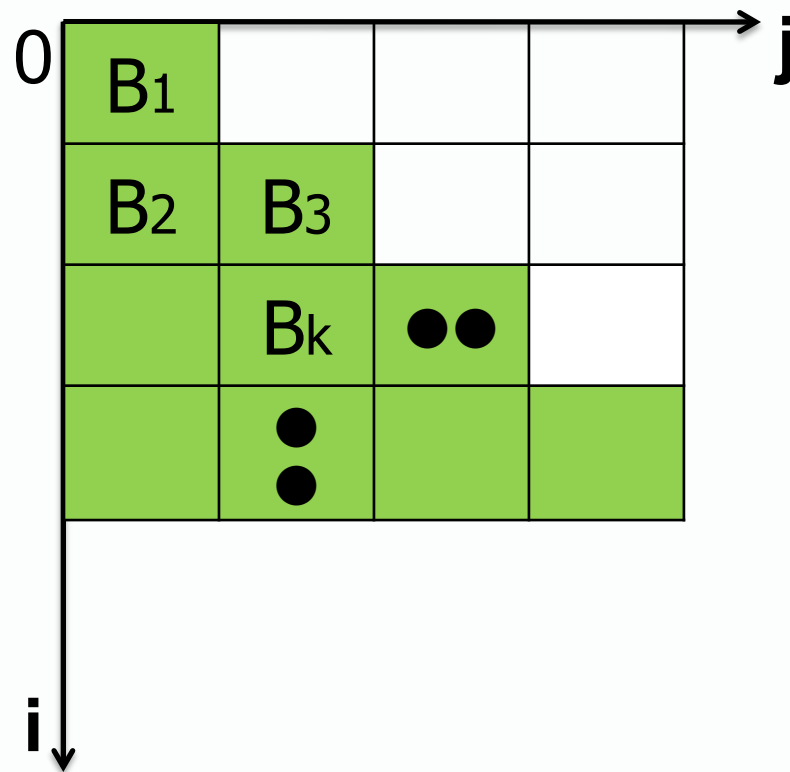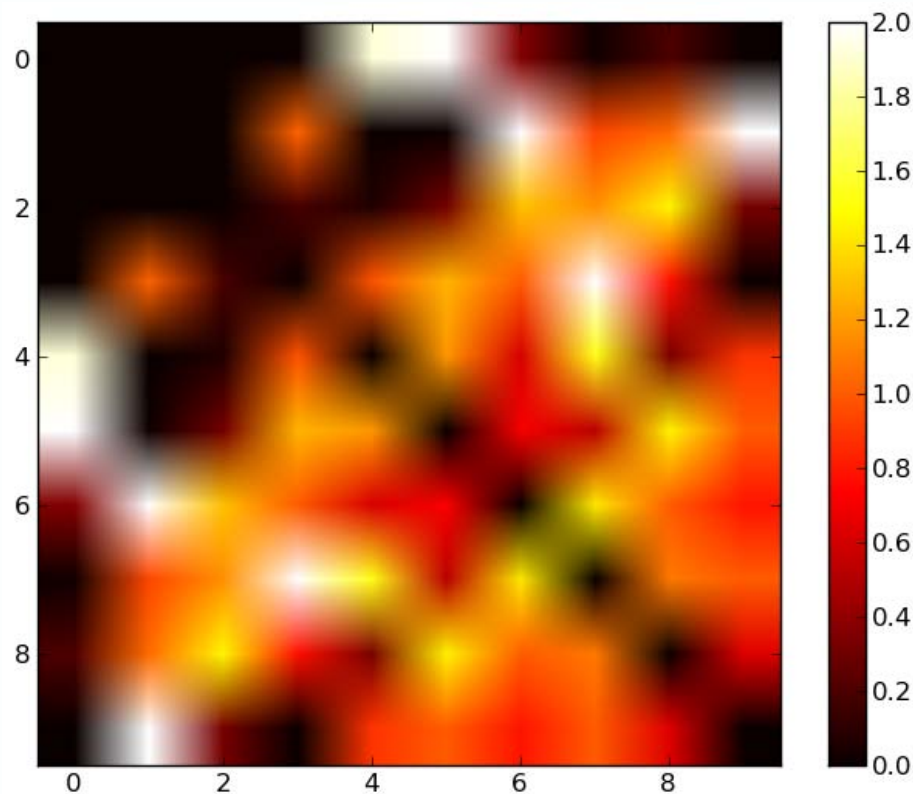
      c. Merge the two centers

$$<i, j> \leftarrow \arg\max_{i,j} sim(c_i, c_j)$$

(3) Output the hierarchical clusters.
}

# Pair-wise distances



Symmetric similarity matrix

0 | B₁ | | | → j
B₂ | B₃ | | |
| Bₖ | ●● | |
| ● ● | | |

# Iterative map/reduce



Data Points

Assignment table

Find the global nearest centers and merge them

**Iterative map/reduce**

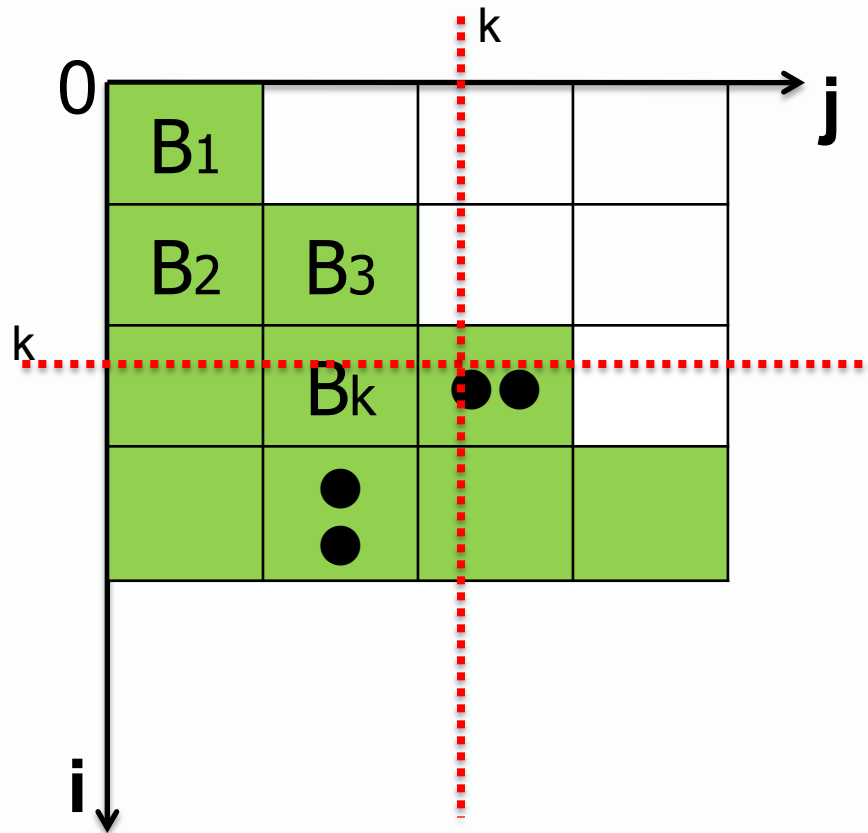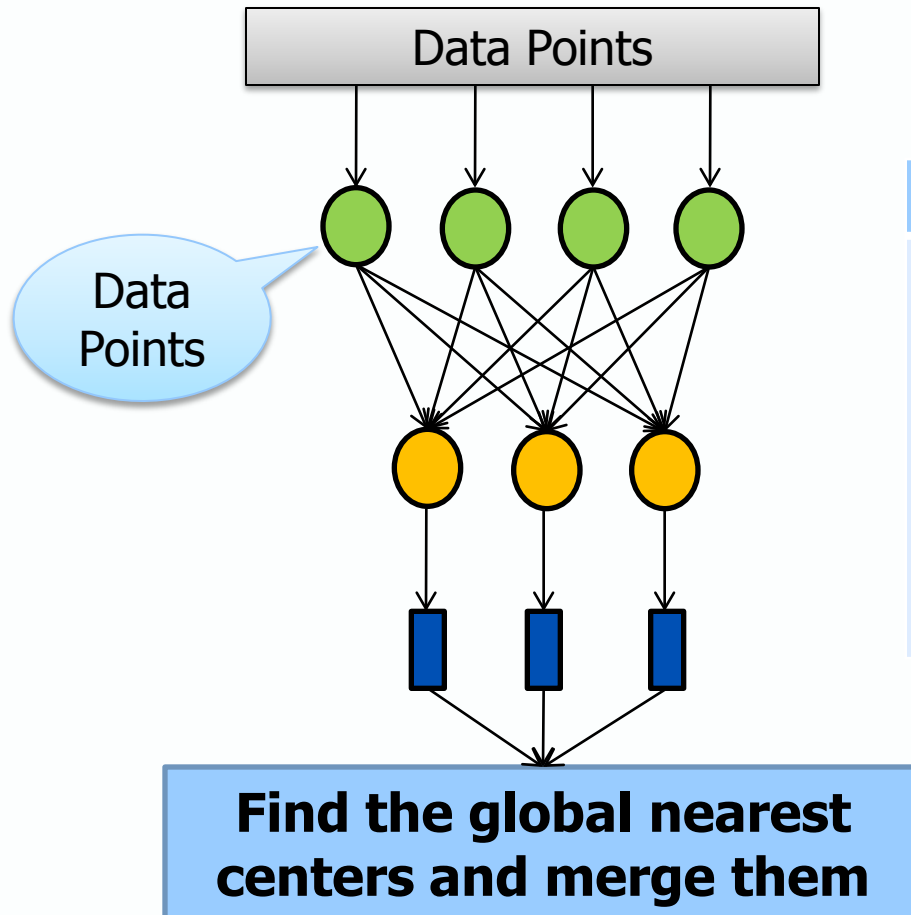*Mapper:*
     assign Block IDs to each data point.
*Reducer:*
     each reducer is responsible for find the local nearest pair of centers.

# Assignment table

# Iterative map/reduce

Data Points

Data Points

**Find the global nearest centers and merge them**
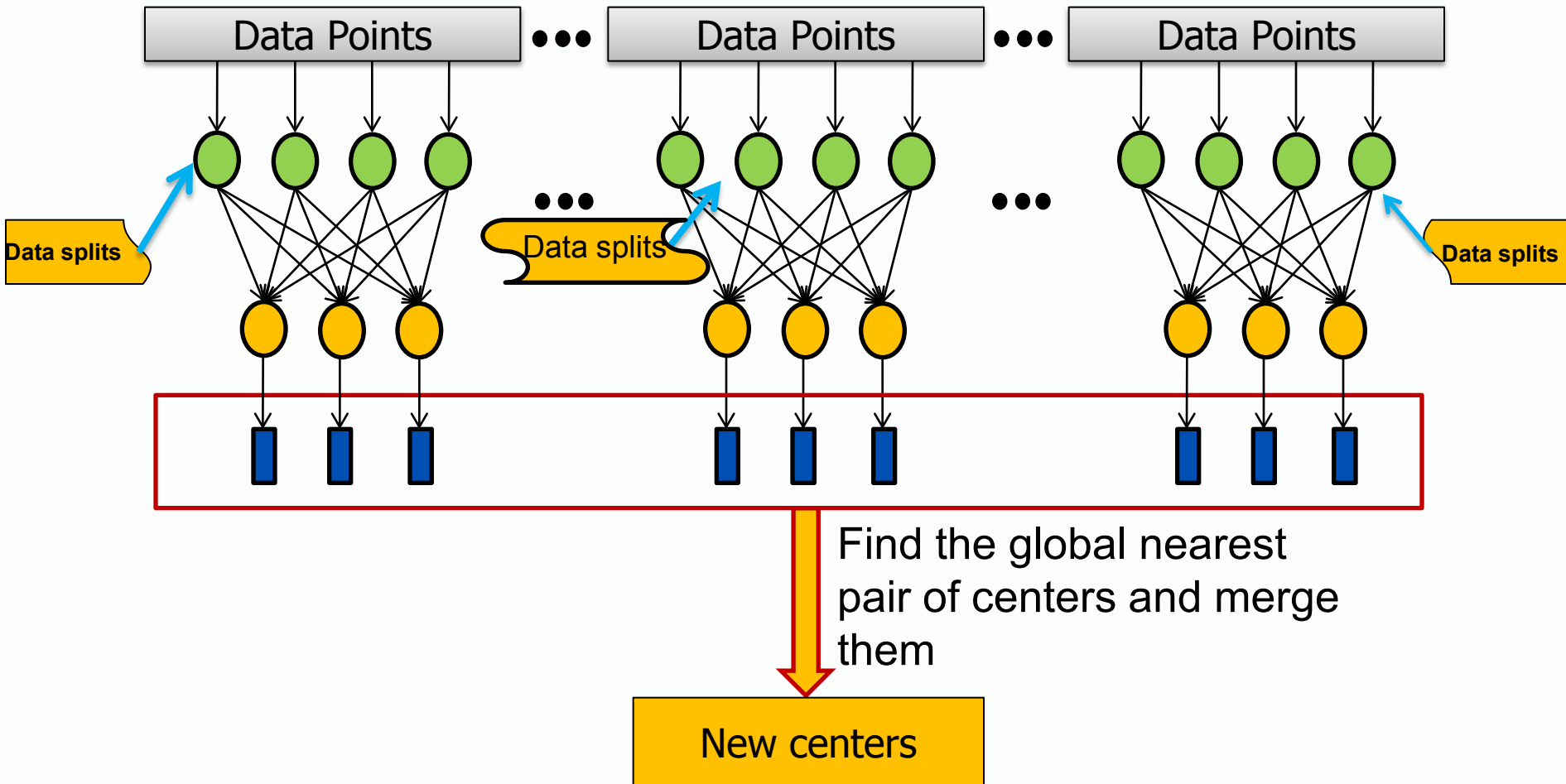
## Iterative map/reduce

*Mapper:*
   preload centers in memory and compute distances from each center to all the centers coming into mappers.
*Reducer:*
   each reducer is responsible for find the local nearest pair of centers.

# Block map/reduce



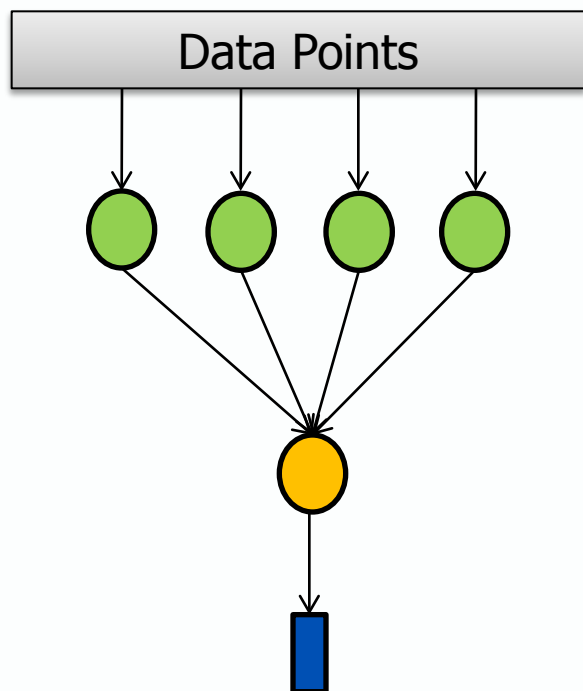Find the global nearest pair of centers and merge them

# Partition

- 100,000,000 users have to be partitioned to blocks before clustering.

- User profile helps to partition users into overlapping blocks.

- There are millions of blocks and each block contains several thousands of users on average.

# Speed-up

- For small blocks, iterative map/reduce is not efficient for the overhead of start and end of a job.
- Only few of elements of the similarity matrix need to be updated.

# One-off map/reduce

Data Points

One-off map/reduce

*Mapper:*
    Passing friend list to reducers.
*Reducer:*
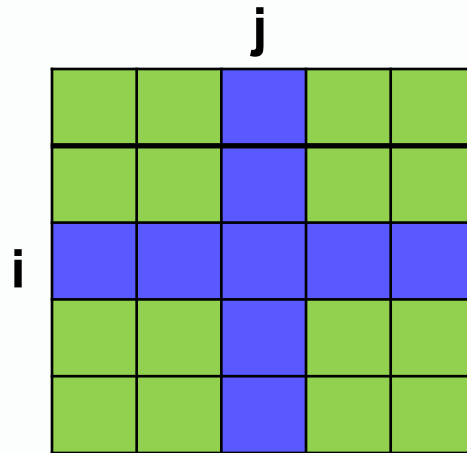    Agglomeration until clustering stops and  output clustering results.

# Scalability

- Only suitable for m*n matrix where m,n $< 10^4$

# Distance caching

- Avoid re-calculating pair-wise distance between centers not for agglomeration from iteration to the next.

# Compressed storage for lower triangle

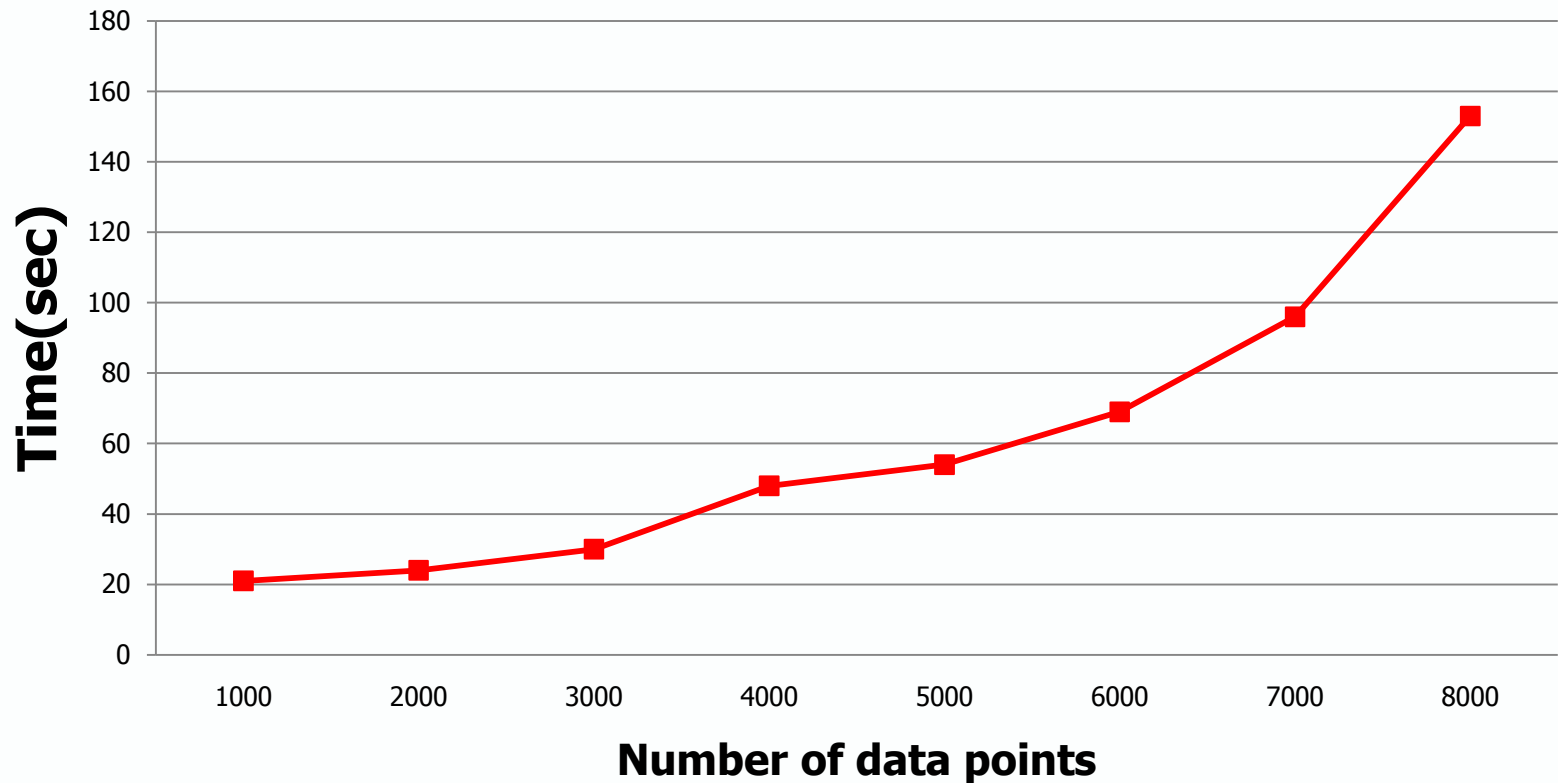- Choose a row compressed storage mode to keep pair-wise distance between centers in memory.



$$K=(i+1)*i/2+j$$

# Performance


Iterative map/reduce — Time(sec) vs Number of data points

# Performance



Overall efficiency

Time(m) vs Num of data points(million)

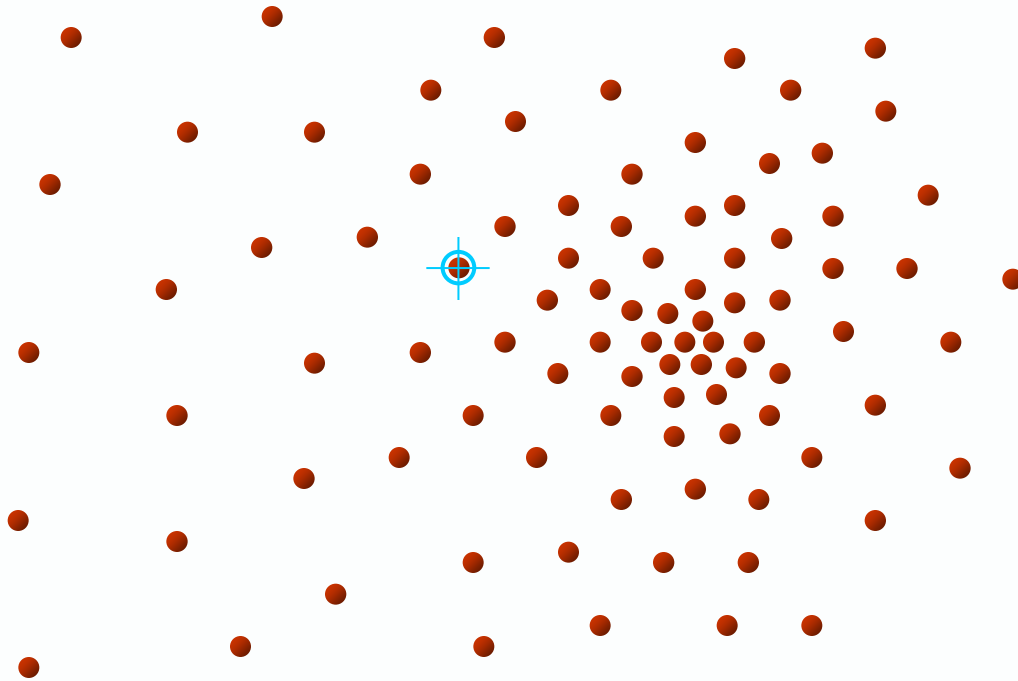Legend: Iterative, one-off, Cached distances

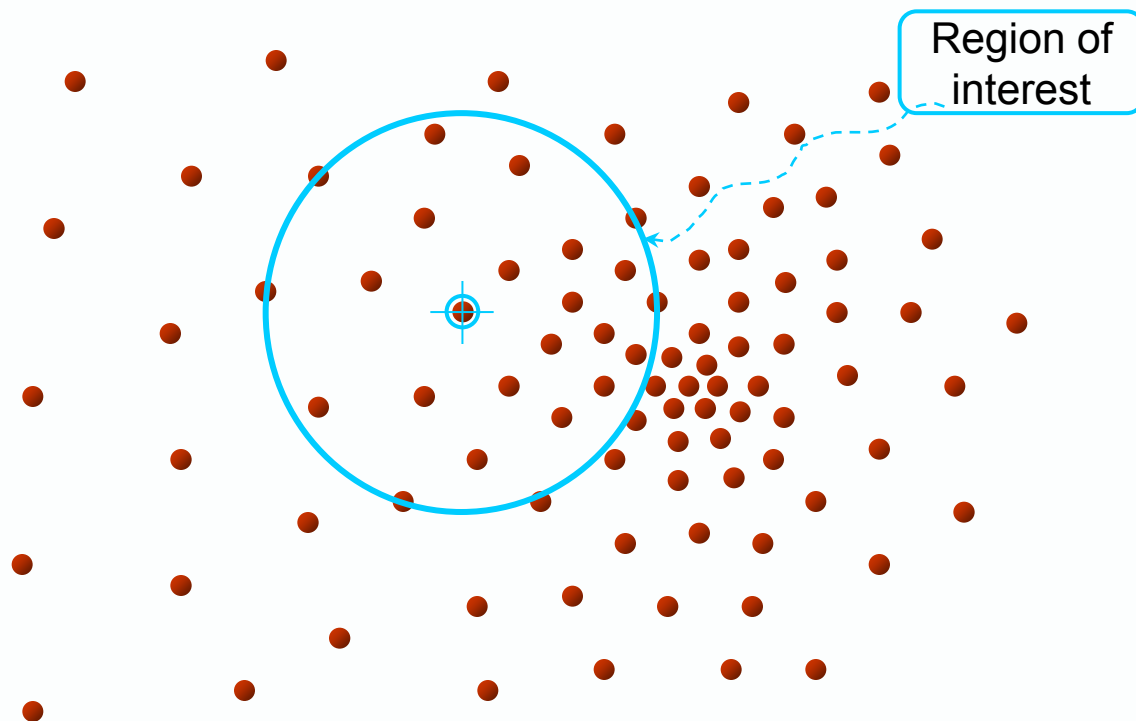# Case II: Topic detection by Mean-shift clustering

- Clustering to find topics in news feed
  - High density indicates hot news.
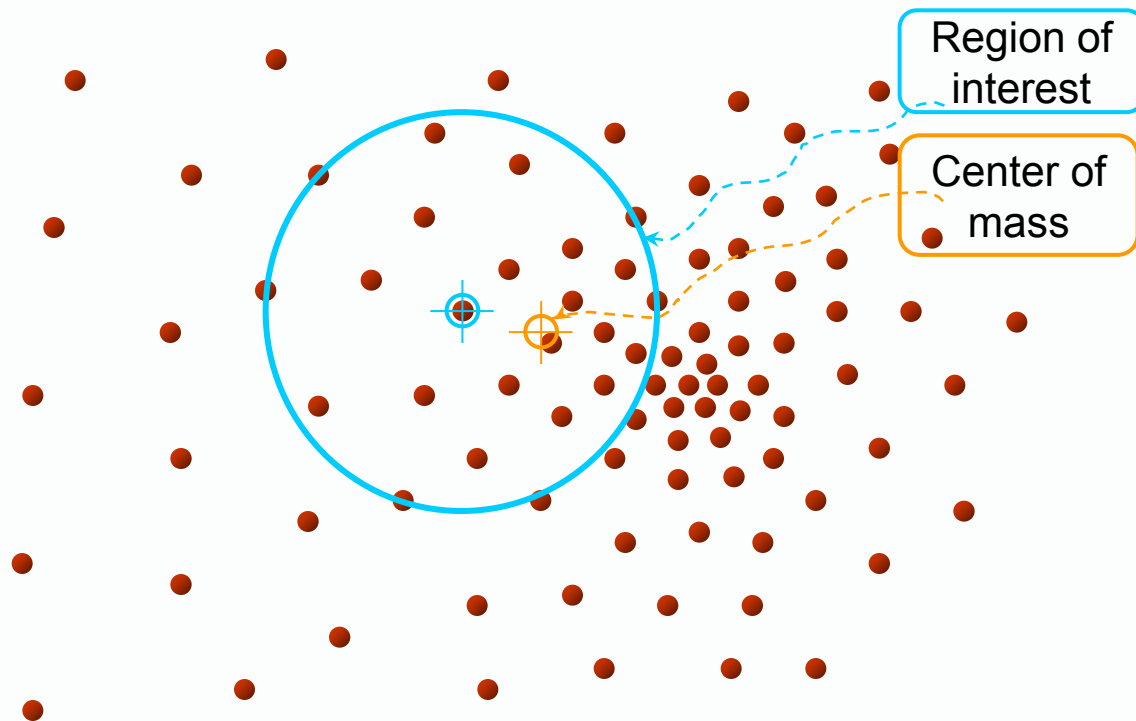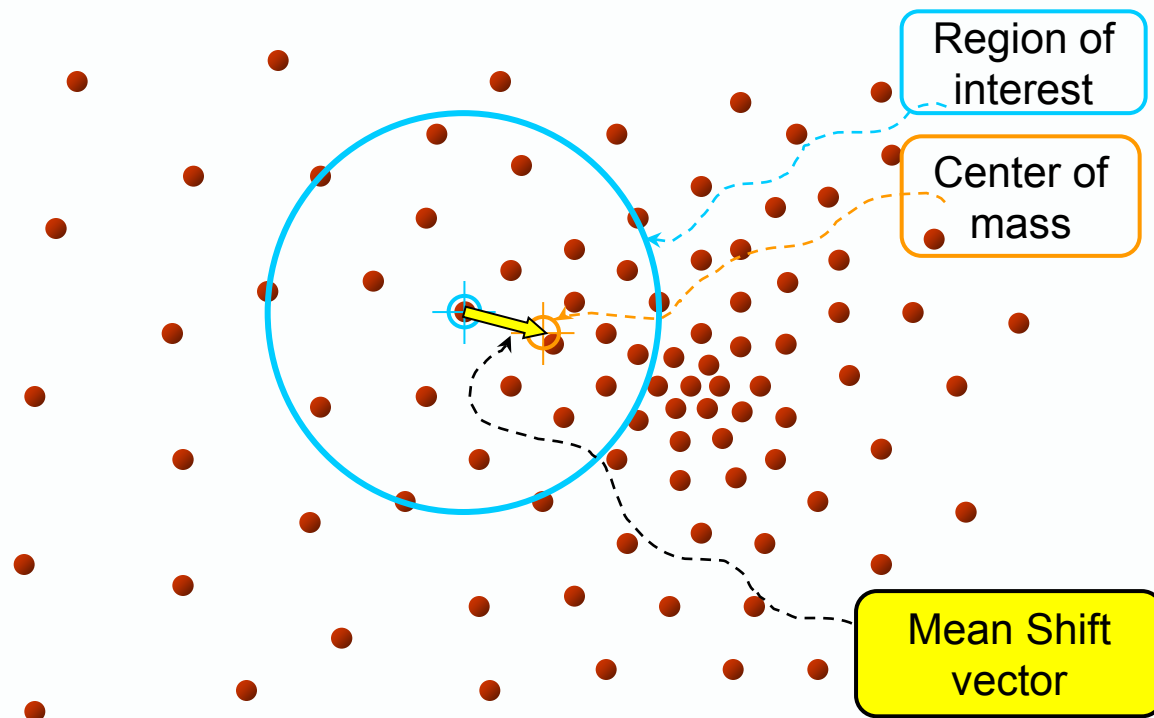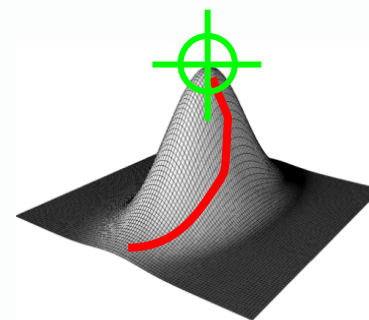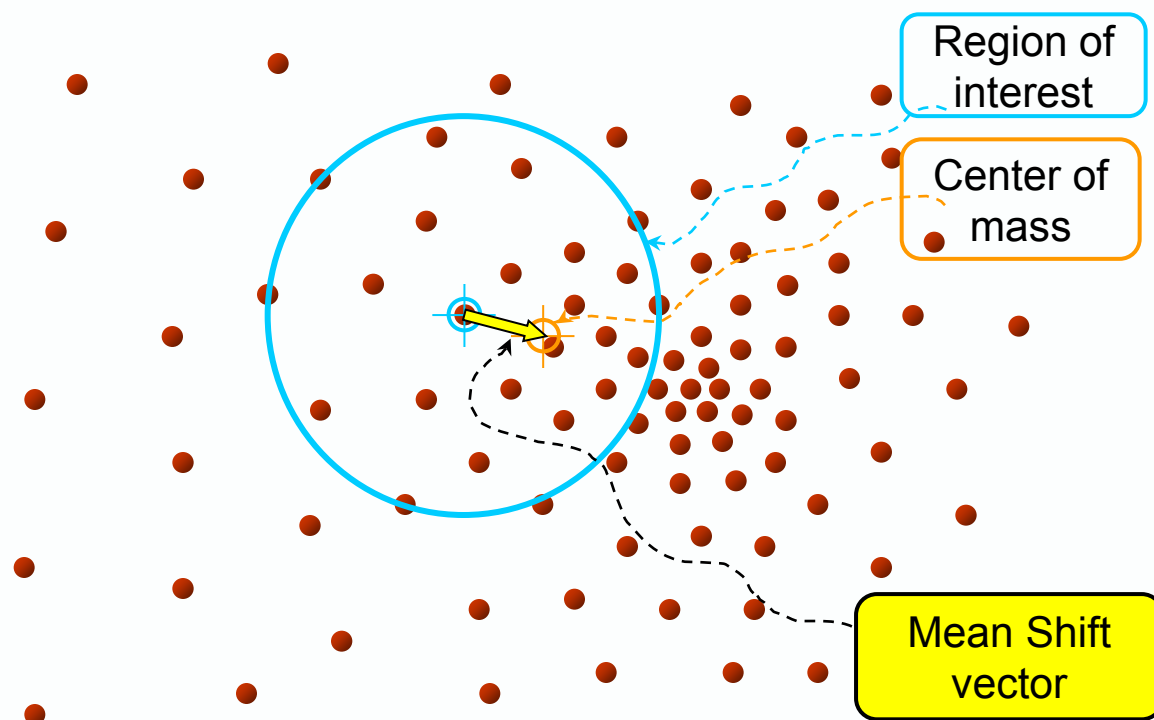  - Without knowing the number of topics.

# Mean shift

# Mean shift



Region of interest

# Mean shift



Region of interest

Center of mass

# Mean shift



Region of interest

Center of mass

Mean Shift vector

# Mean shift



Region of interest

Center of mass

Mean Shift vector

# Mean shift clustering



Compute center of mass give a random point

Generate a new center

iteration

Possible merge between close centers

# Iterative map/reduce



**Data Points**

Assign label to points within bandwidth

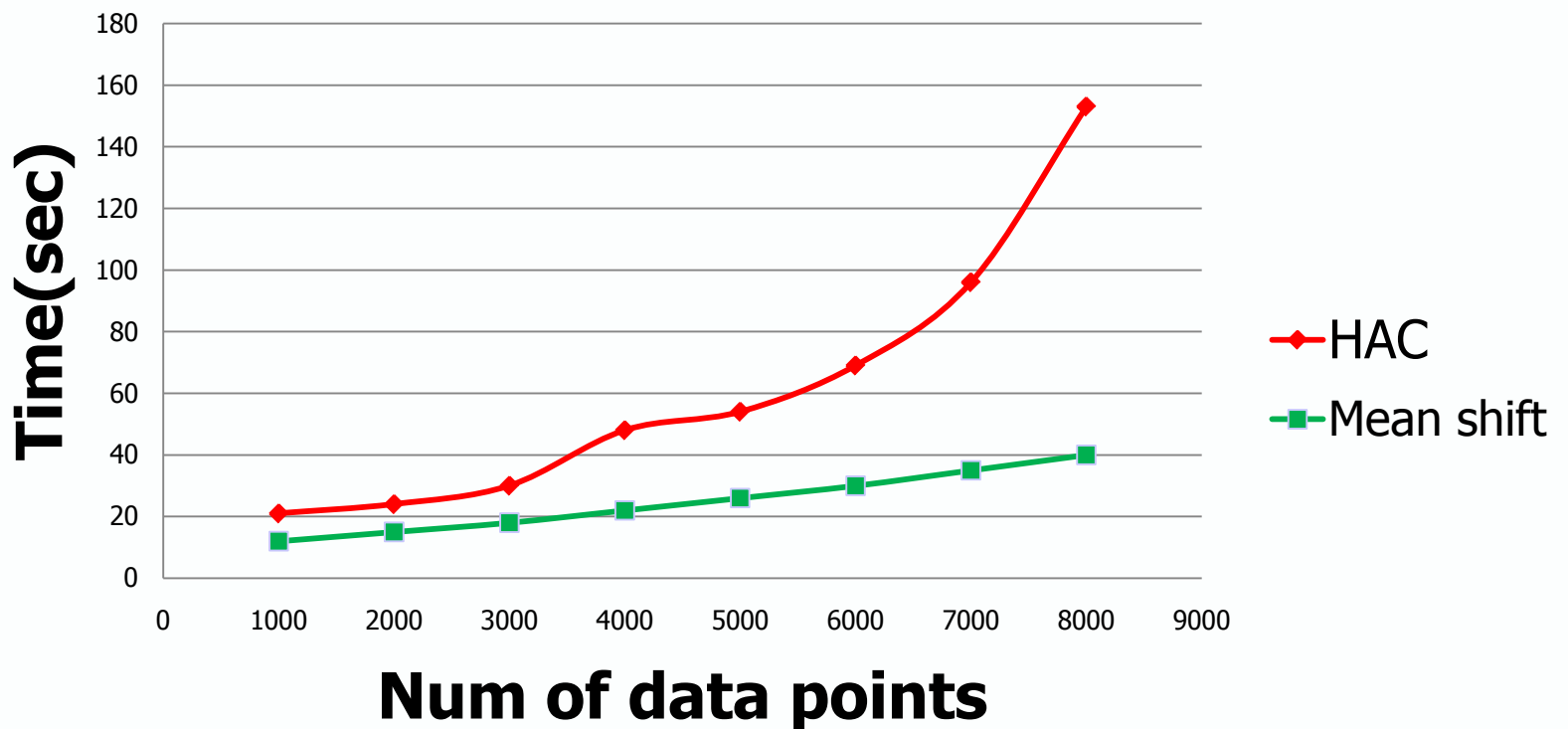**Update centers**

## Iterative map/reduce

*Mapper:*
select a point as the center and compute distances from the center to all the data points, and assign a label to the points within the bandwidth.
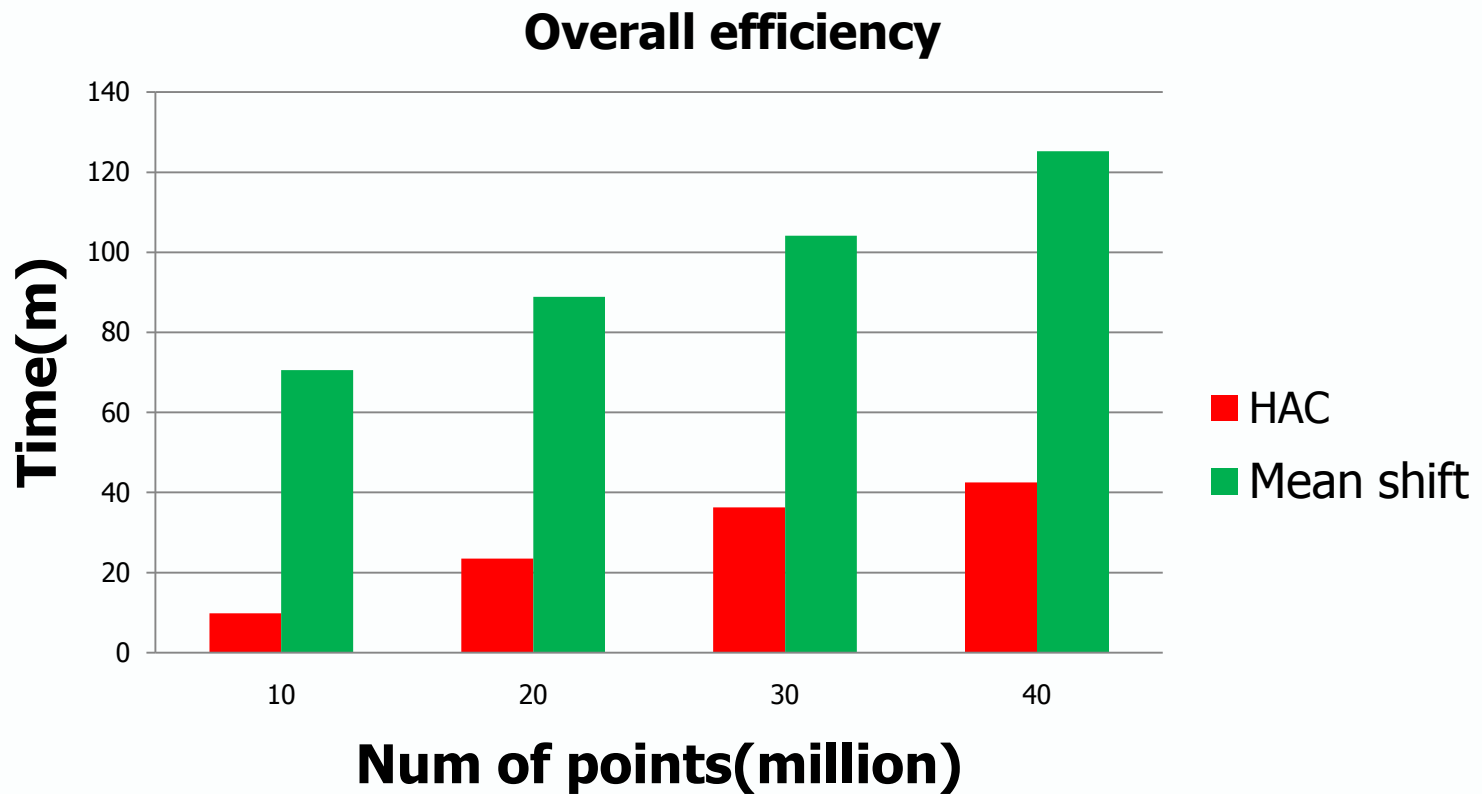
*Reducer:*
collect data points of the same label and compute the center of mass of them.

# Performance

# Performance

# 人人数据平台

白伯纯
bochun.bai@renren-inc.com
http://renren.com/bbc

张叶银
yeyin.zhang@renren-inc.com
http://renren.com/hartebeest