



Mux-Kmeans: Multiplex Kmeans for Clustering Large-scale Data Set

Chen Li, Yanfeng Zhang, Minghai Jiao, Ge Yu
Northeastern University, China



東北大學
Northeastern University

BIG DATA

IN A SINGLE DAY ONLINE

ENOUGH INFORMATION IS CONSUMED TO FILL

168 MILLION DVDS

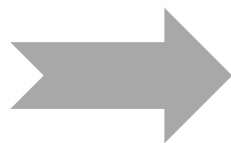
294bn E-MAILS
ARE SENT

MINUTES SPENT
ON FACEBOOK **4.7M**

2 MILLION BLOG POSTS
ARE WRITTEN

VIDEO UPLOADED TO
YOUTUBE **864,000 HRS**

MORE IPHONES
ARE SOLD THAN BABIES BORN



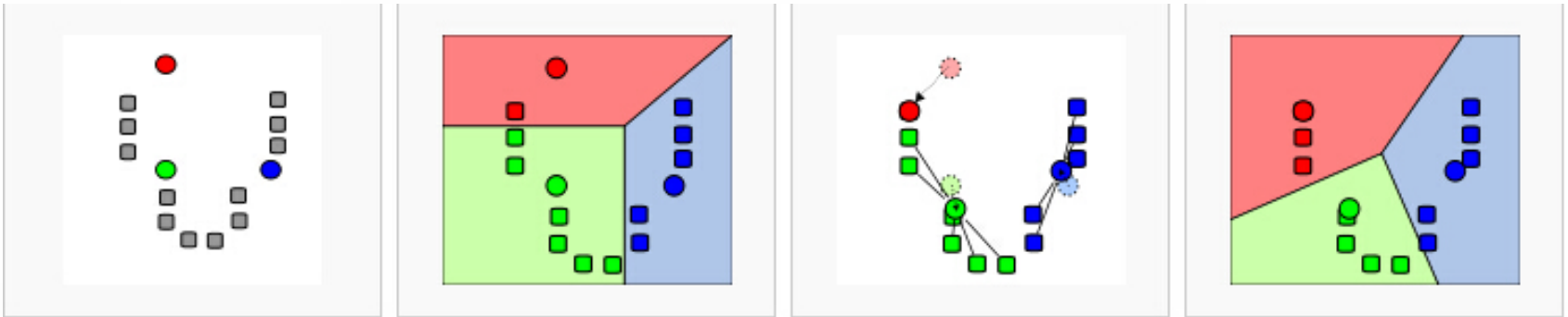
Clustering



Kmeans

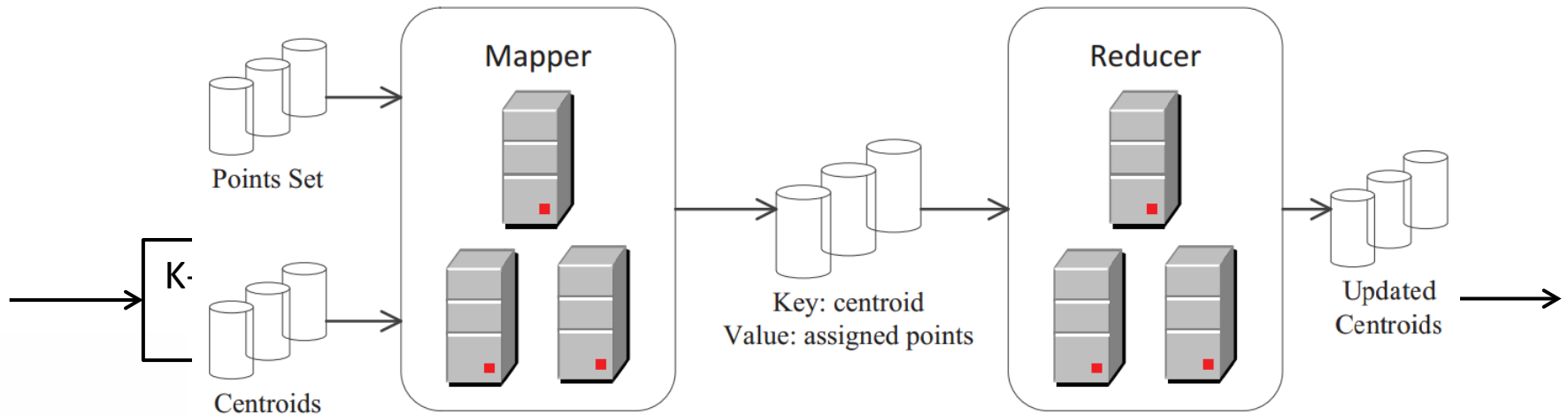
Kmeans

- Kmeans: accept K center patterns and a data set, divide the set into K clusters
- Goal:
 - ✓ 1. similar – data patterns in same cluster;
 - ✓ 2. dissimilar – data patterns in different clusters.



http://en.wikipedia.org/wiki/K-means_clustering

Kmeans on MapReduce

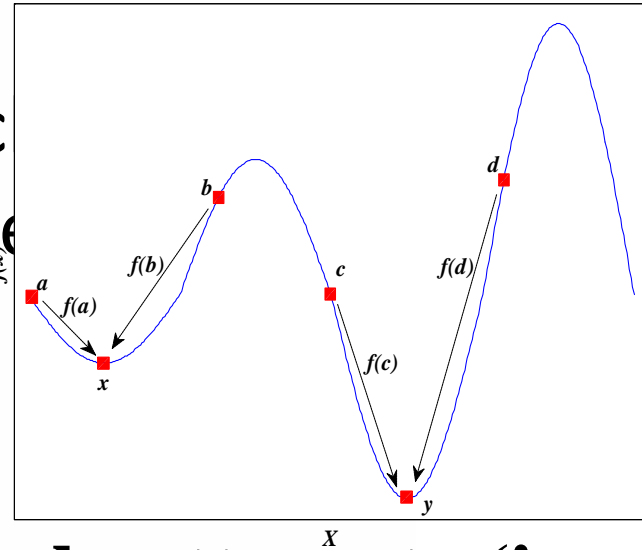
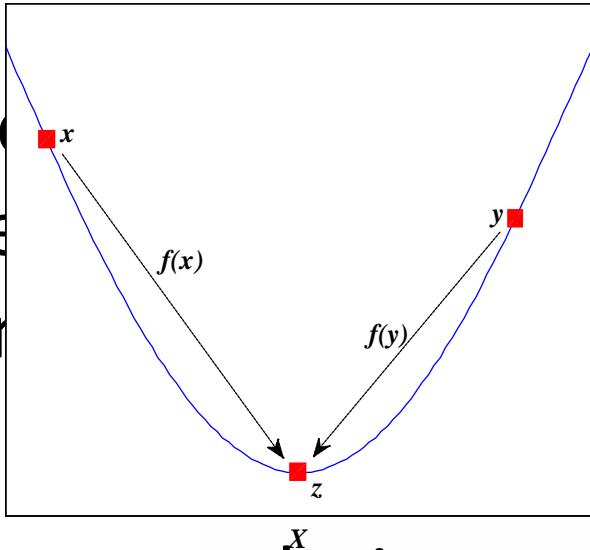


```
procedure MAPPHASE( $\langle i, x_i \rangle$ )  
  for  $c_j, j \in \{1, \dots, k\}$  do  
     $j \leftarrow$  find the closest centroid  $c_j$  to  $x_i$ ,  
     $j \in \{1, \dots, k\}$   
  end for  
  output  $\langle j, x_i \rangle$   
end procedure
```

```
procedure REDUCEPHASE( $\langle j, [x_i] \rangle$ )  
   $c_j \leftarrow$  average  $[x_i]$   
  output  $\langle j, c_j \rangle$   
end procedure
```

Shortcoming of Kmeans

- The current solution is not the global optimum, but only a local optimum, determined by initial centroids.
- Current solution: **multiple attempts (in series)**



- Start from multiple groups of initial centroids
- Execute multiple kmeans processes, obtain multiple local optima
- Pick the one with the lowest cost function

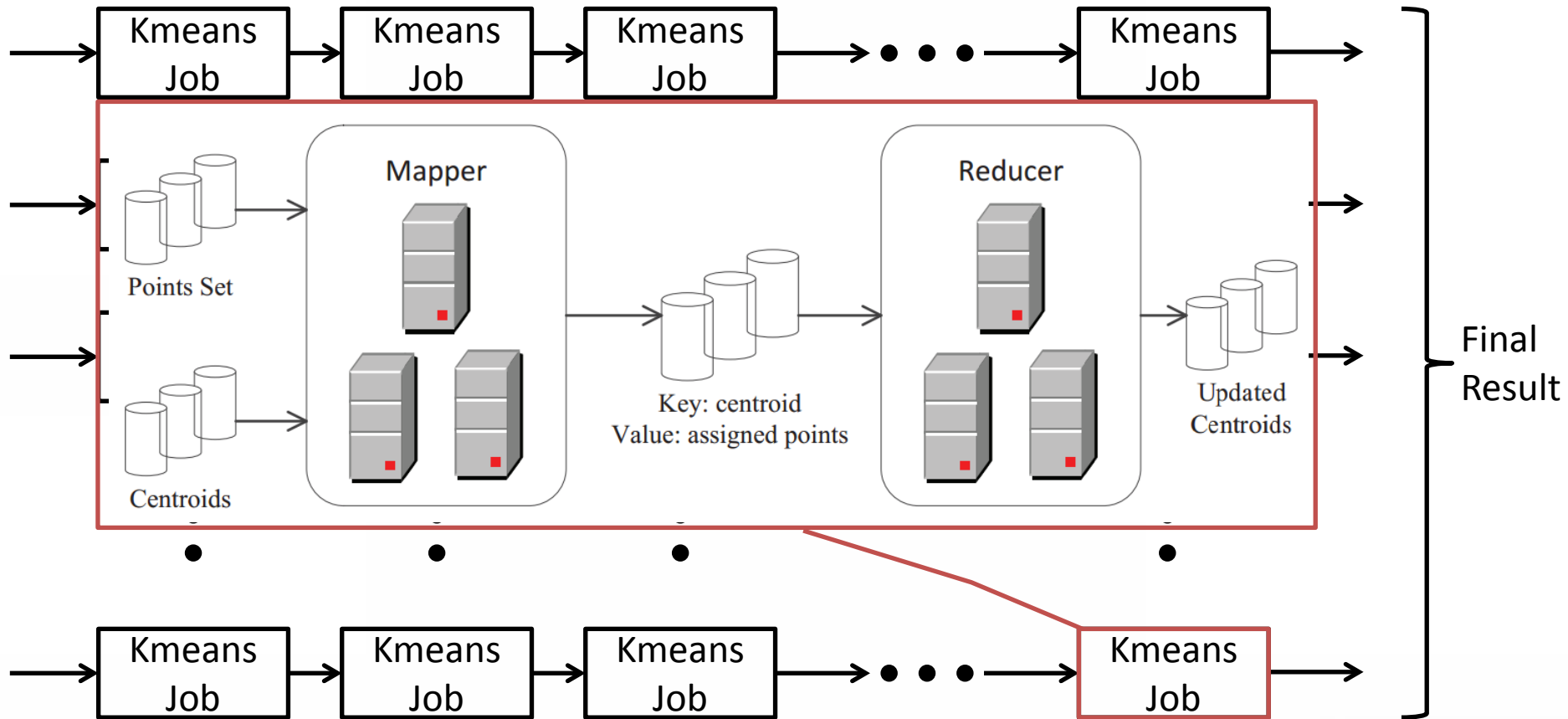
Efficiency Problem



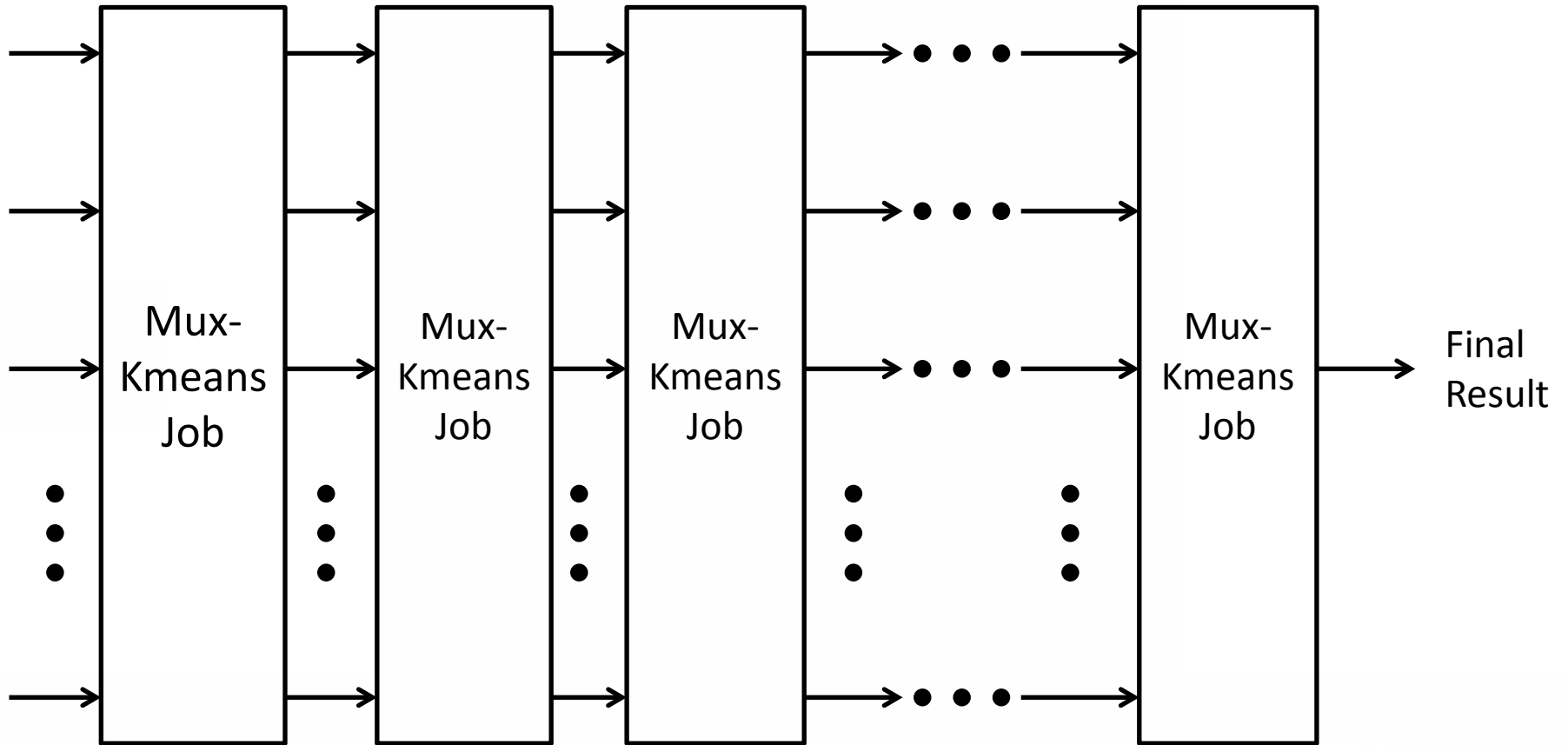
Mux-Kmeans

- Idea:
 - Execute multiple Kmeans attempts in parallel
 - Share states across different Kmeans processes
 - Terminate “hopeless” attempts in early stage
 - Expand searching scale and try more attempts
- Goal:
 - Guarantee the clustering quality
 - Decrease the runtime

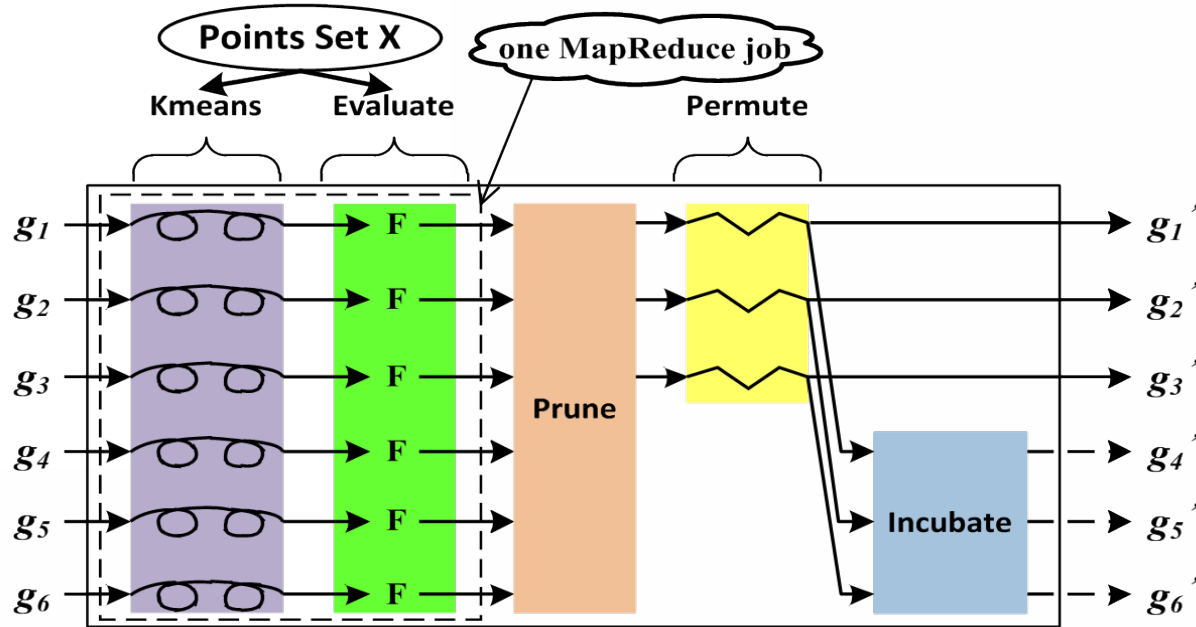
Naive Kmeans VS. Mux-Kmeans



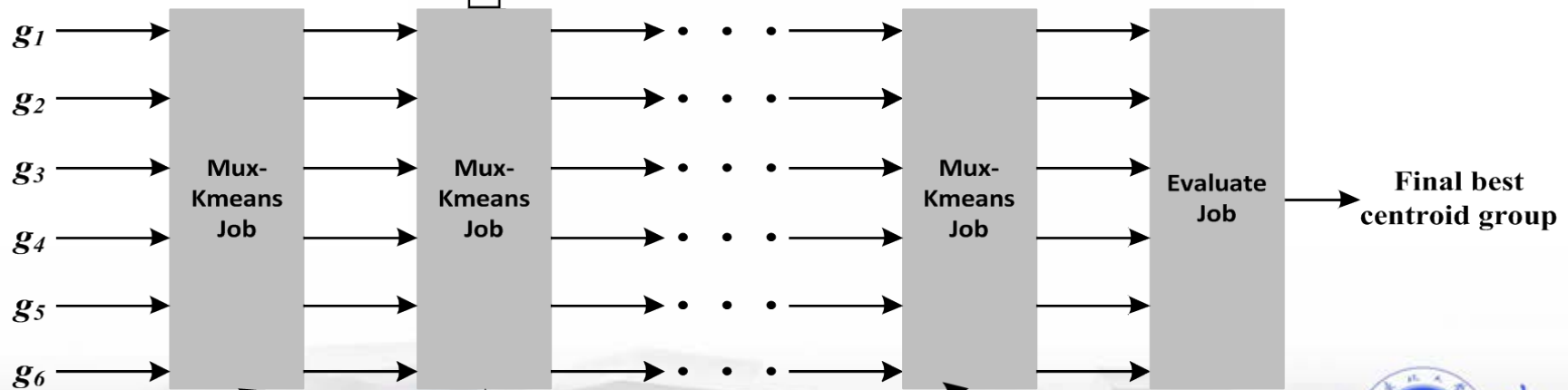
Naive Kmeans VS. **Mux-Kmeans**



Inside Mux-Kmeans Job

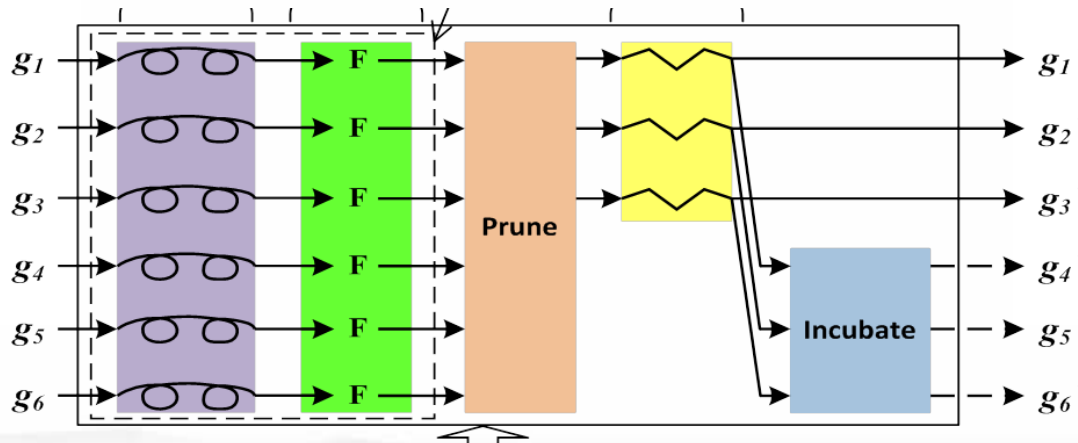


1. **Kmeans** clustering
2. **Evaluate** intermediate result's quality
3. **Prune** some Kmeans attempts with low quality
4. **Incubate** new groups of centroids for next iteration



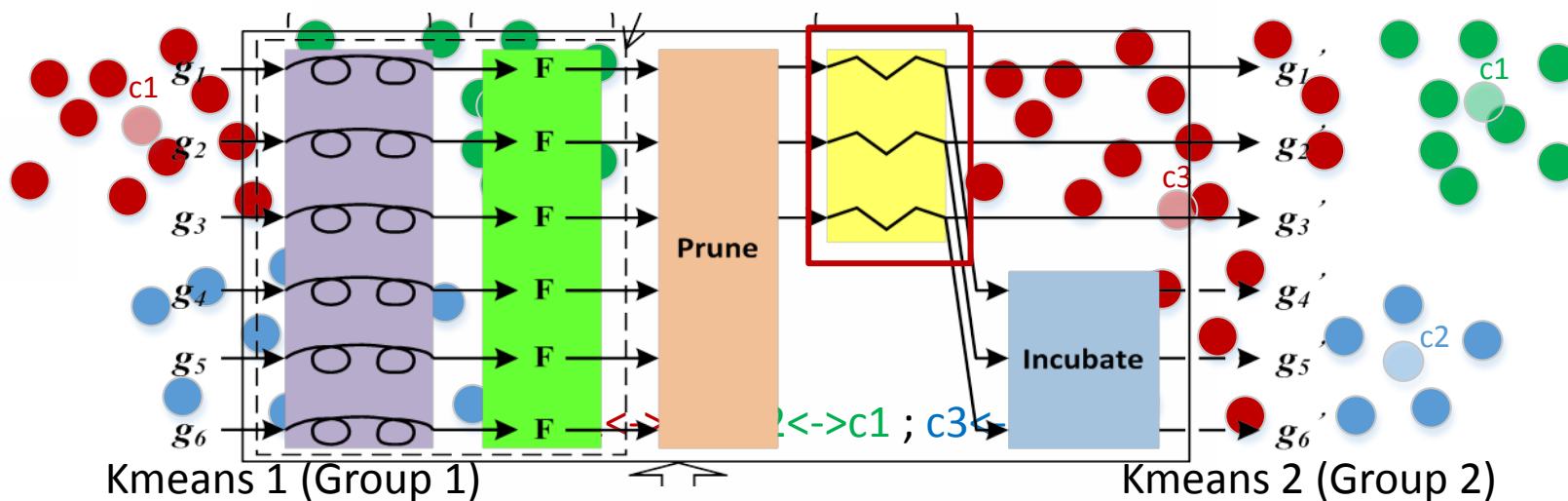
Kmeans, Evaluate & Prune

- Use Kmeans algorithm to do clustering
 - Get multiple updated centroid groups
- Use Total Within-Cluster Variation (TWCV) to evaluate different centroid groups' quality.
- Prune x% centroid groups with relatively low quality



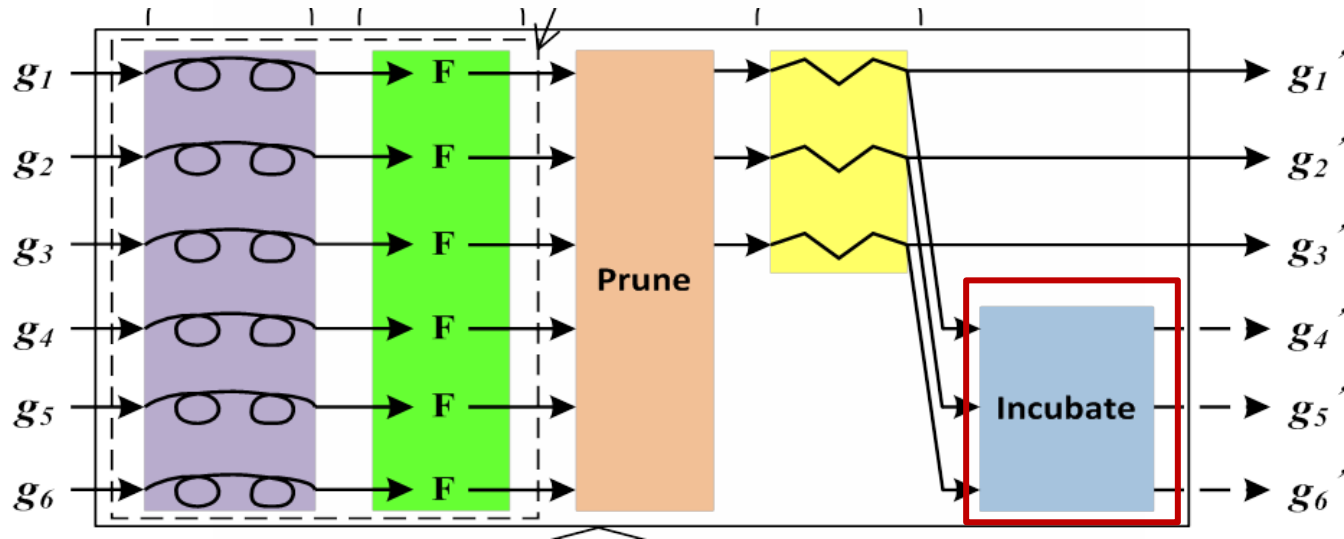
Permute

- A preprocessing step before the incubate step
- Aim to find the related centroids between different groups

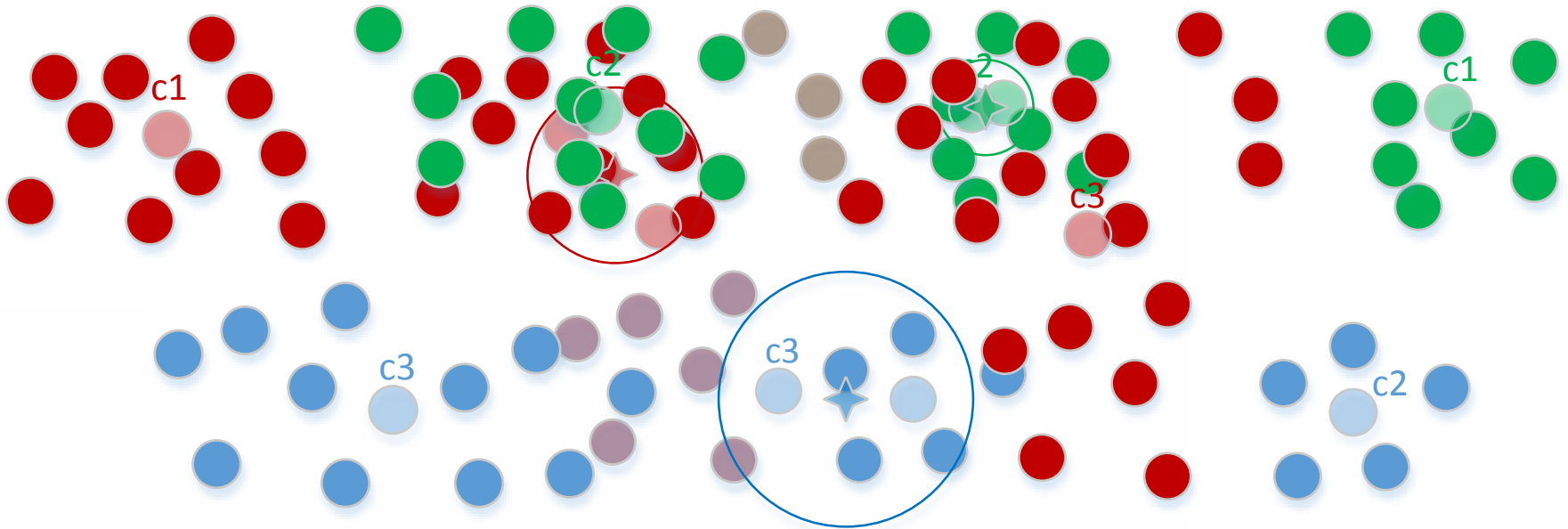


- $\text{sim}(c_i, c_j) = \frac{1}{1 + \text{distance}(c_i, c_j)}$, $\text{sim}(c_i, c_j) \in [0, 1]$

Incubate



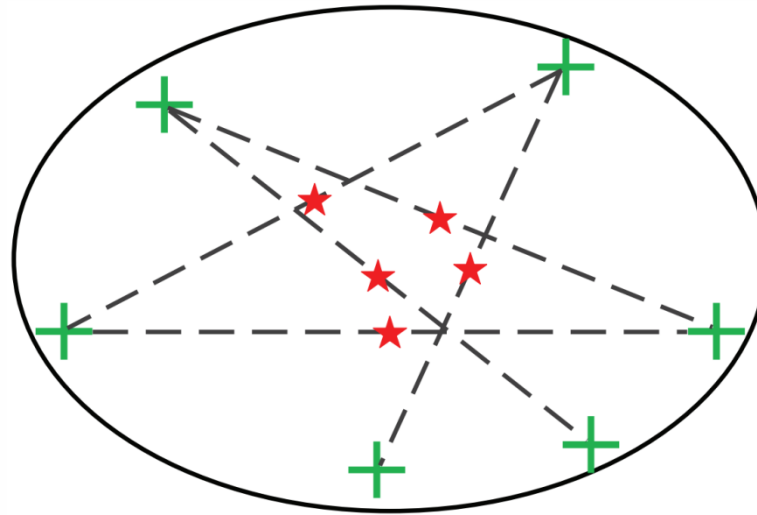
Incubate 1: *RSDS*



Random Search within a Definite Scope

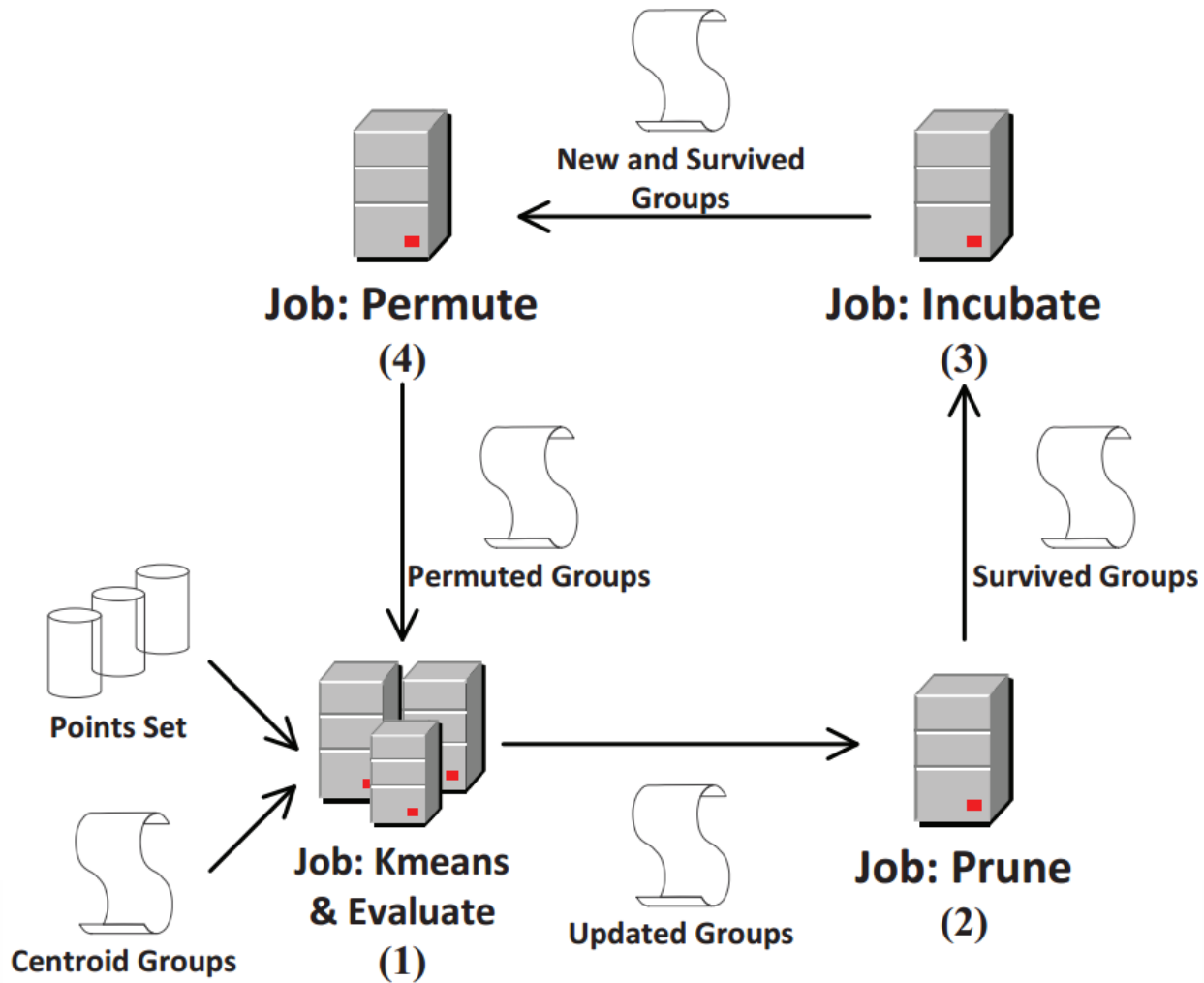
1. Star = compute the middle of two related centroids
2. Random search around the star, radius = $2 * \text{distance}(\text{centroid}, \text{start})$

Incubate 2: *ADGP*



Average of Dissimilar Group Pairs

Implementation



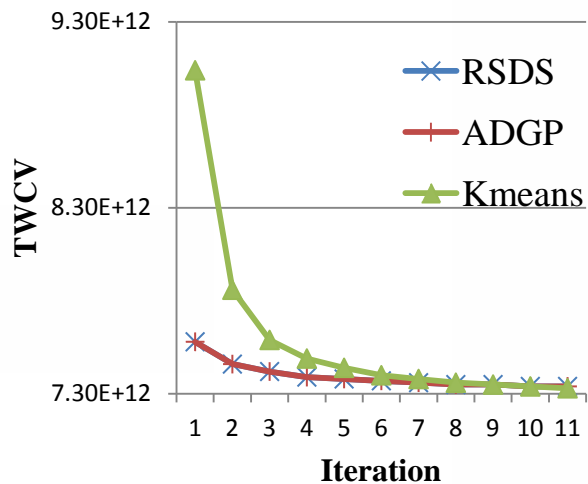
Experiment Setup

- Experiment Environment
 - Amazon EC2: 16 nodes, Ubuntu 12.04, Hadoop 1.4.1
 - Each node: 2 ECUs and 1 CPU; 3.7 GB memory; 410 GB storage; moderate network performance
- Experiment Dataset

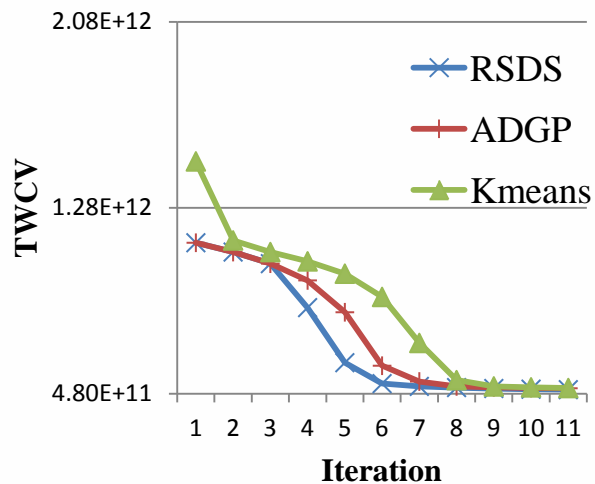
	Points	features
Bio_train	145751	74
Netflix	17770	1000
Lastfm	359330	40

Clustering Quality

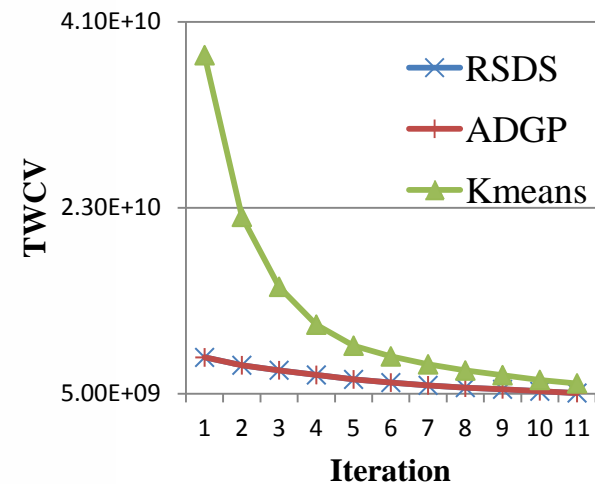
Different data set



Lastfm



Bio_train



Netflix

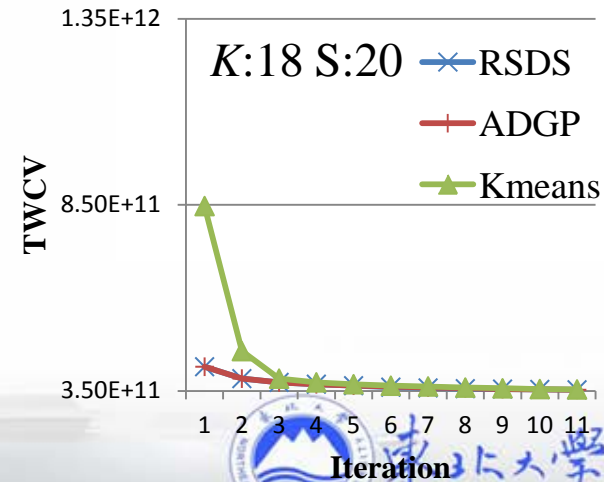
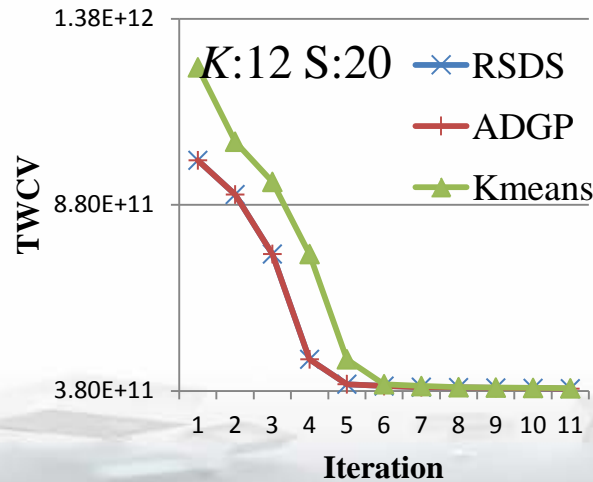
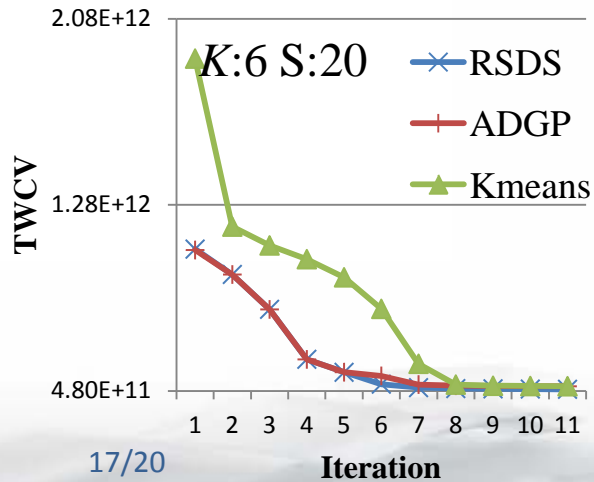
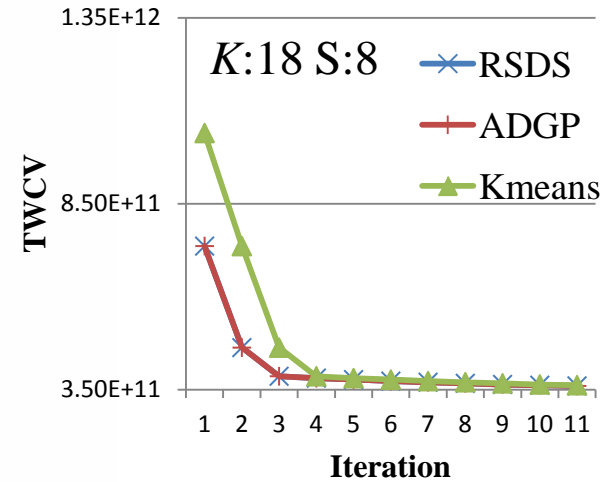
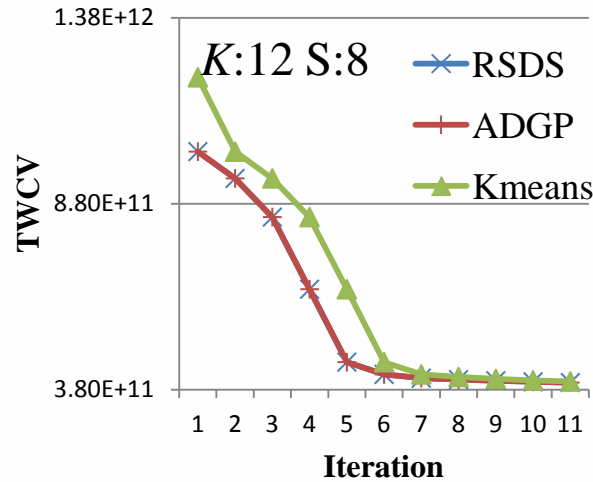
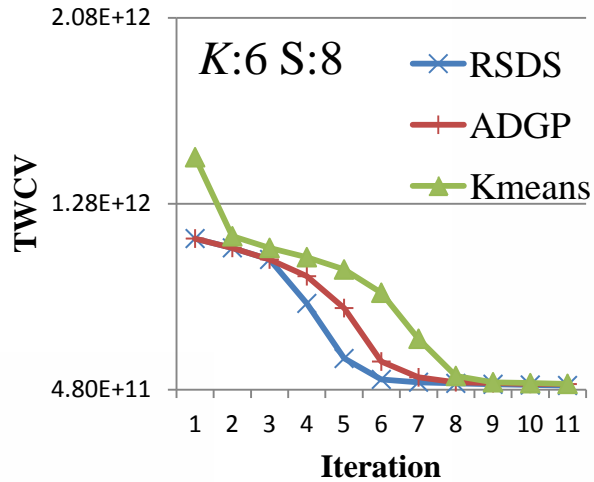
K value:6

centroid group amount:8

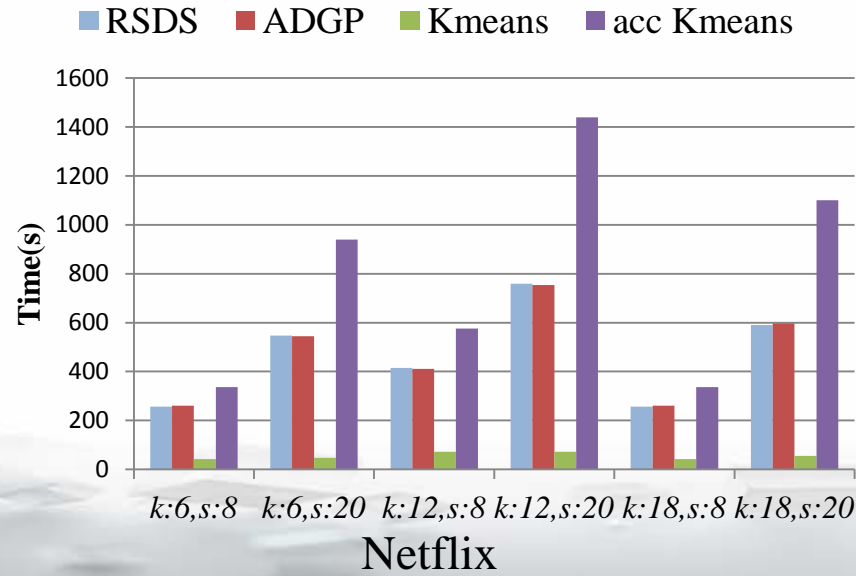
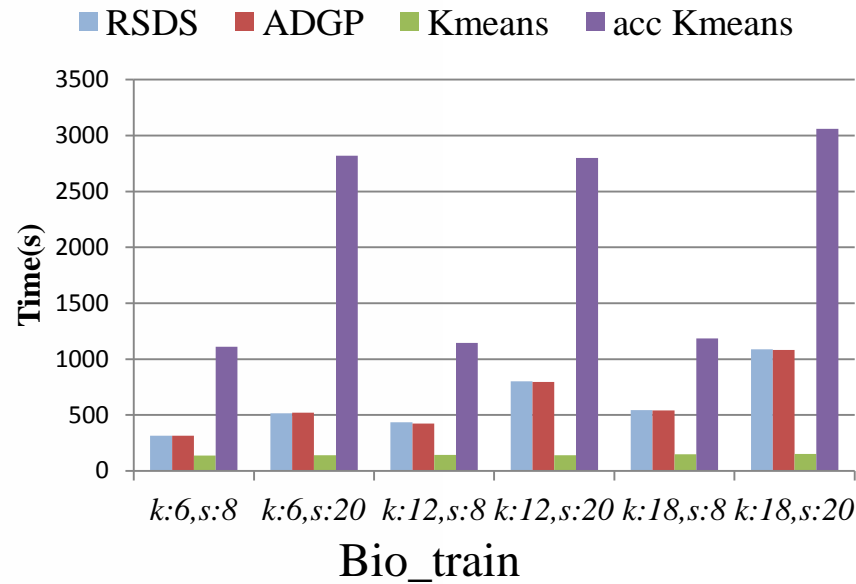
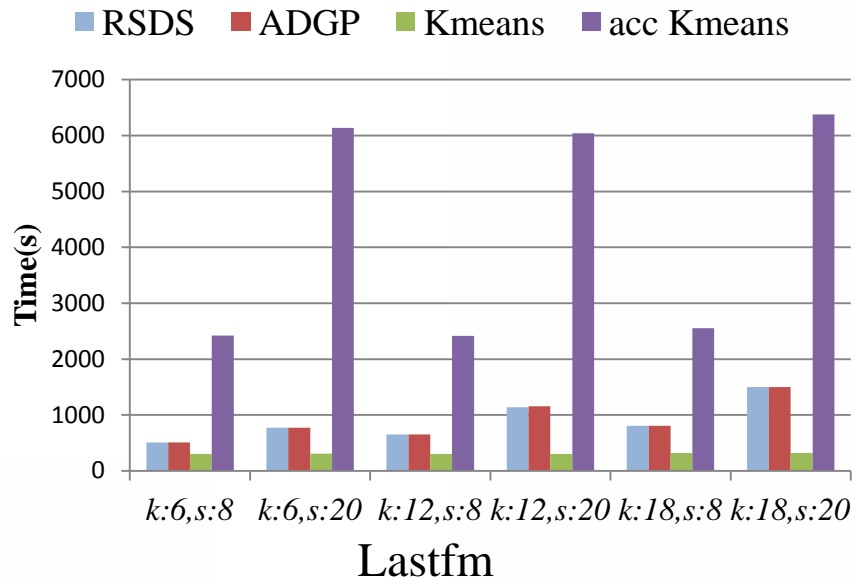
Clustering Quality

Same data set: Bio_train

✖S: centroid groups amount



Elapsed Time



Summary & Future Work

- The Mux-Kmeans algorithm.
 - Idea: execute multiple attempts in parallel, share states across different kmeans processes, terminate “hopeless” attempts in early stage, expand search scale and try more attempts
 - Implementation: deployed on MapReduce
 - Result: better clustering quality and shorter runtime when processing multiple centroid groups
- Future work
 - Different k in different centroid groups
 - Many possible Mux-XXX algorithms (Mux-EM, Mux-FCM, etc.)



Thank you!

QUESTIONS?