

# DisCo

A Case Study on Distributed Mining

Spiros Papadimitriou  
Jimeng Sun

IBM T.J. Watson Research Center

## Data processing and mining

Mining process

```

    graph LR
      A[Gathering] --> B[Pre-processing]
      B --> C[Analysis]
      C --> D[Post-processing]
  
```

## Data processing and mining

Mining process

```

    graph LR
      A[Gathering] --> B[Pre-processing]
      B --> C[Analysis]
      C --> D[Post-processing]
  
```

Feature extraction, Aggregation, ...

Model estimation, Validation, ...

Scalable infrastructures for end-to-end mining?

"It's a dirty job..."

BUT:  
Both are equally important and equally time consuming

## Data processing and mining

Mining process

```

    graph LR
      A[Gathering] --> B[Pre-processing]
      B --> C[Analysis]
      C --> D[Post-processing]
  
```

BigTable, HBase, ...

Data access APIs

Processing APIs

Sawzall, Pig, Cascading, ...

Google FS, HDFS, KFS, ...

Distributed File System

Shared computation and storage nodes

Distributed framework

## Overview

Mining process

```

    graph LR
      A[Gathering] --> B[Pre-processing]
      B --> C[Analysis]
      C --> D[Post-processing]
  
```

- Background **MapReduce**
- Feature extraction **Aggregates & edgelists**
- Mining **Co-clustering**
- Conclusion

## Overview

Mining process

```

    graph LR
      A[Gathering] --> B[Pre-processing]
      B --> C[Analysis]
      C --> D[Post-processing]
  
```

- Background: MapReduce
- Feature extraction
- Mining
- Conclusion

### Example – Programming model

```
employees.txt
# LAST FIRST SALARY
Smith John $90,000
Brown David $70,000
Johnson George $95,000
Yates John $80,000
Miller Bill $65,000
Moore Jack $85,000
Taylor Fred $75,000
Smith David $80,000
Harris John $90,000
... ..
... ..
```

```
mapper
def getName(line):
    return line.split('\t')[1]
reducer
def addCounts(hist, name):
    hist[name] = \
    hist.get(name, default=0) + 1
    return hist

input = open('employees.txt', 'r')
intermediate = map(getName, input)
result = reduce(addCounts, \
                intermediate, {})
```

Q: "What is the frequency of each first name?"

Three (fixed) "streams":  
key-value iterators

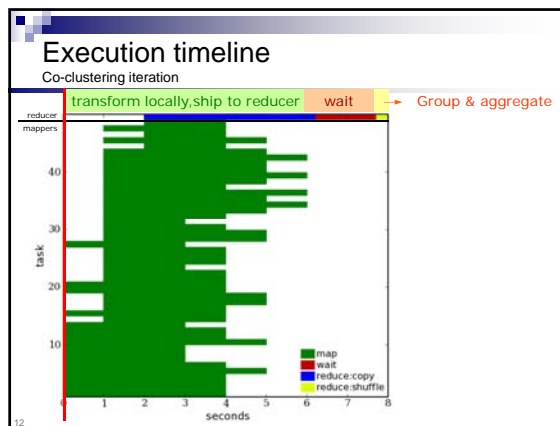
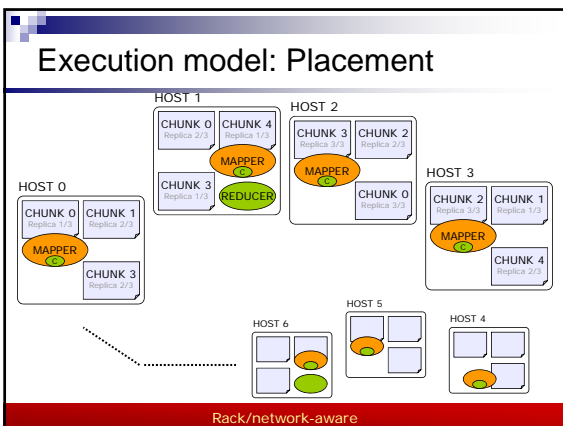
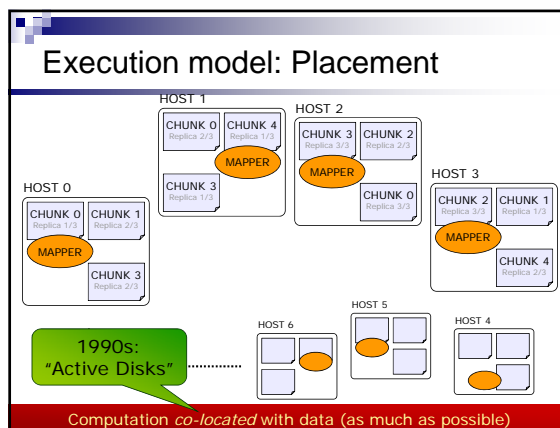
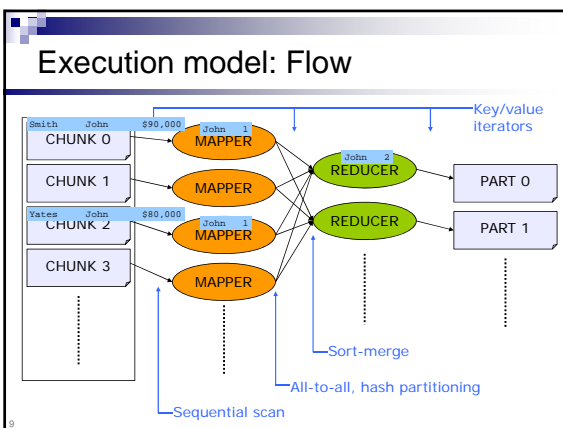
### Comparison

|                                    |  |
|------------------------------------|--|
| Quick-n-dirty script               | vs. MapReduce (Hadoop)                           |
| ~5 lines of (non-boilerplate) code |  |
| Single machine,<br>local drive     | Up to <i>thousands</i> of<br>machines and drives |

What is hidden to achieve this:

- Data partitioning, placement and replication
- Computation placement (and replication)
- Number of nodes (mappers / reducers)

As a user, you don't *need* to know  
what I'm about to show you next...



## Hadoop

- Open-source implementation (ASF/Yahoo!)
- Main Hadoop components:
  - HDFS: Clone of Google FS
  - MapReduce
  - HTable: Clone of BigTable
  - Hadoop On Demand (HOD)

## Overview

Mining process

```

    graph LR
      A[Gathering] --> B[Pre-processing]
      B --> C[Analysis]
      C --> D[Post-processing]
    
```

- Background
- Feature extraction
- Mining
- Conclusion

## Histogram benchmark

- 347.5GB raw data input (text)
- ~30KB total reducer output
- Up to 39 quad-core nodes
- ...see paper for setup details

## Scalability

Simple benchmark

| Number of nodes | Aggregate bandwidth (Mbps) |
|-----------------|----------------------------|
| 1               | 113                        |
| 25              | 3766                       |
| 40              | 6844                       |

~170Mbps/node

7 hours (at 1 node), 7 minutes (at 40 nodes)

Fibre (at ~1500 Mbps), Single drive (at ~500 Mbps)

Caveat: cluster is running a single job

## Parameter tuning study

- Process pool size
- Number of reducers
- Input file split size
- ...see paper

## Pre-processing vs. co-clustering

Edgelist extraction, 39 nodes

← Pre-processing | Co-clustering →

ISS: Pre-processing ~0.5, Co-clustering ~0.5

TREC: Pre-processing ~0.6, Co-clustering ~0.4

... both are equally important and equally time consuming

### Overview

Mining process

```

    graph LR
    A[Gathering] --> B[Pre-processing]
    B --> C[Analysis]
    C --> D[Post-processing]
  
```

- Background
- Feature extraction
- Mining: (co-)clustering**
- Conclusion

### Co-clustering Summary

**Split:** Increase  $k$  or  $l$

**Shuffle:** Rearrange rows and cols

### Search for solution Shuffles

Similarity (KL-divergence) of row fragments to blocks of a row group

Assign to second row-group

### (Co-)clustering with MapReduce

| KEY | VAL          |
|-----|--------------|
| 1   | 5 7 13       |
| 2   | 3 9 11 19 27 |
| 3   | 6 12         |
| ... | ...          |
| m   | 98           |

Sequence file in HDFS

### (Co-)clustering with MapReduce

Sequence file in HDFS

Cluster statistics

Broadcast job parameters

### (Co-)clustering with MapReduce

Sequence file in HDFS

Cluster statistics

Broadcast job parameters

### (Co-)clustering with MapReduce

25

### (Co-)clustering with MapReduce

26

### (Co-)clustering with MapReduce

27

### (Co-)clustering with MapReduce

28

### (Co-)clustering with MapReduce

Sequence file in HDFS

Broadcast job parameters

Scales w.r.t. number of edges – what about nodes ?

29

### Co-clustering with MapReduce

Sequence file

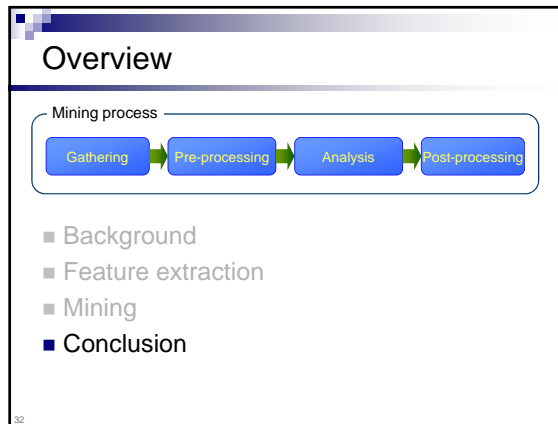
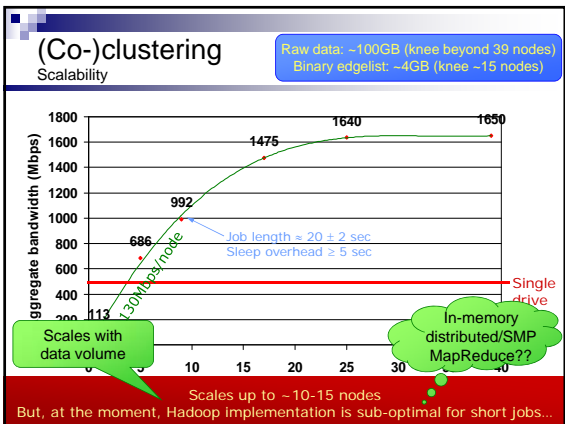
Sequence file in HDFS

Broadcast job parameters

sharing scans ?

map-side joins ?

30

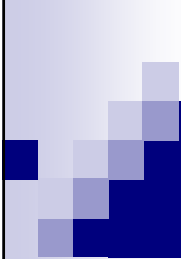


- ### Other work
- Hadoop/MapReduce "ecosystem"
- Lower-level (storage/virtualization):
    - Kosmos FS (KFS), Amazon S3, Parascle VSN, Ceph, Lustre/Gluster, PanFS
    - Amazon EC2
  - Higher-level (data access):
    - BigTable, HBase, Hypertable
  - Higher-level (computation):
    - PIG, Cascading, Sawzall, Dryad/DryadLINQ

- ### Other work
- Distributed/Parallel Mining
- Sector
  - FREERIDE-G
  - MapReduce for ML on multi-core [NIPS06]

- ### Conclusions and lessons
- MapReduce
- Simple programming model
  - Part of a growing, open "ecosystem" of data processing tools
  - Scalable and fault-tolerant
    - Scalability ≠ performance
  - Ideal for (pre-)processing large volumes of data
- [...] MapReduce is the first instruction of the "data center computer" [...]  
 — David Patterson, "The Data Center Is The Computer", CACM, Jan. 2008

- ### Conclusions and lessons
- Distributed mining process
- Data pre-processing *and* mining
    - End-to-end view
    - Both equally important and time-consuming
  - How big the data *really* is depends on the task at hand (cf. [NIPS06])

The logo consists of several overlapping squares of varying shades of blue and white, arranged in a stepped pattern that suggests a staircase or a grid.

**DisCo**  
A Case Study on Distributed Mining

Spiros Papadimitriou  
Jimeng Sun

*IBM T.J. Watson Research Center*

See also <http://www.bitquill.net/trac/wiki/PCC/Start>