# Introduction to Hadoop

**ID2210**

**Jim Dowling**

# Large Scale Distributed Computing

☐ In #Nodes
- BitTorrent (millions)
- Peer-to-Peer

☐ In #Instructions/sec
- Teraflops, Petaflops, Exascale
- Super-Computing

☐ In #Bytes stored
- Facebook: 300+ Petabytes (April 2014)*
- Hadoop

☐ In #Bytes processed/time
- Google processed 24 petabytes of data per day  in 2013
- Colossus, Spanner, BigQuery, BigTable, Borg, Omega, ..

*http://www.adweek.com/socialtimes/orcfile/434041

# Where does Big Data Come From?

- On-line services                     PBs per day

- Scientific instruments          PBs per minute

- Whole genome sequencing    250 GB per person

- Internet-of-Things                 Will be lots!
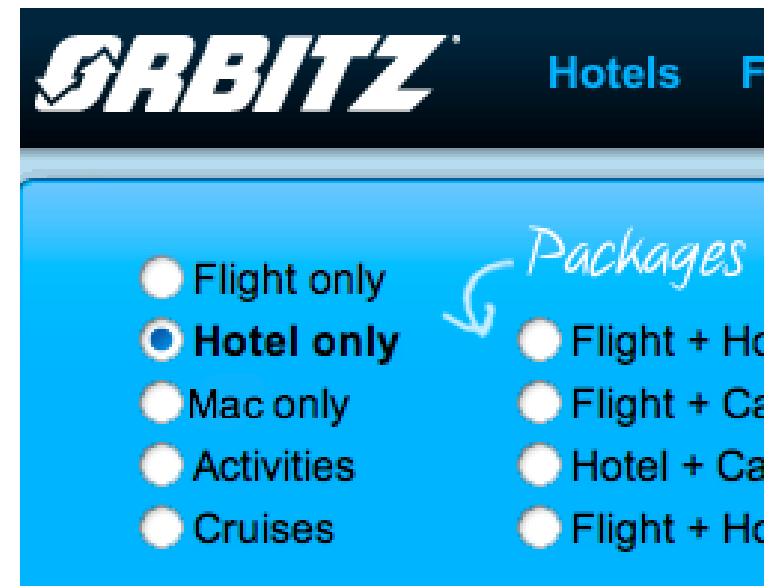
# What is Big Data?



**Small Data**

**Big Data**

# Why is Big Data "hot"?

- Companies like Google and Facebook have shown how to **extract value from Big Data**

Orbitz looks for higher prices from Safari users [WSJ'12]

# Why is Big Data "hot"?

- Big Data helped Obama win the 2012 election through **data-driven decision making**\*



Data said: middle-aged females like contests, dinners and celebrity

\*http://swampland.time.com/2012/11/07/inside-the-secret-world-of-quants-and-data-crunchers-who-helped-obama-win/

KTH
VETENSKAP
OCH KONST
ROYAL INSTITUTE
OF TECHNOLOGY

# Why is Big Data Important in Science?

- In a wide array of academic fields, the ability to effectively process data is superseding other more classical modes of research.


**"More data trumps better algorithms"***

*"The Unreasonable Effectiveness of Data" [Halevey et al 09]

# 4 Vs of Big Data

- Volume

- Velocity

- Variety

- Veracity/Variability/Value
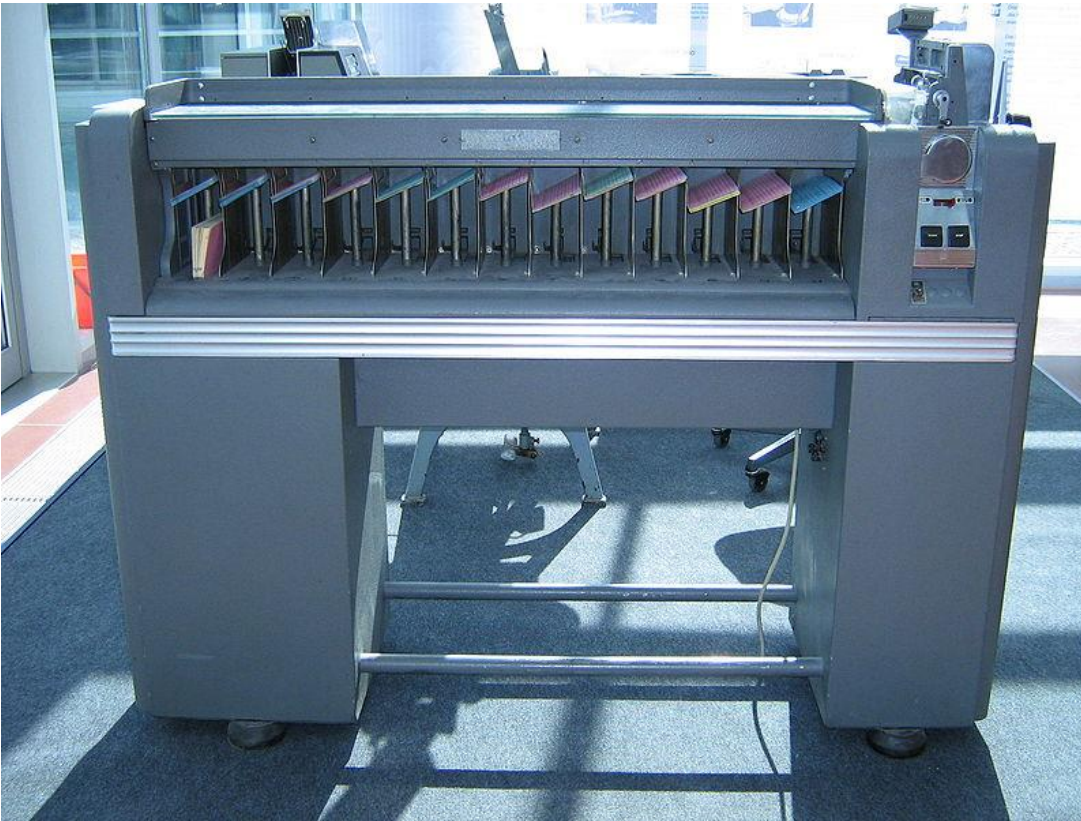
# A quick historical tour of data systems

In the
Beginning

# Batch Sequential Processing

Scan → Sort



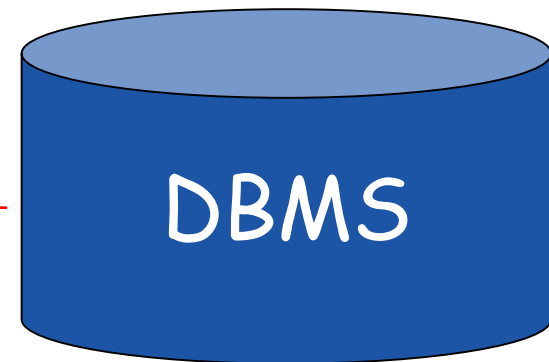IBM 082 Punch Card Sorter



No Fault Tolerance ☺

# 1960s

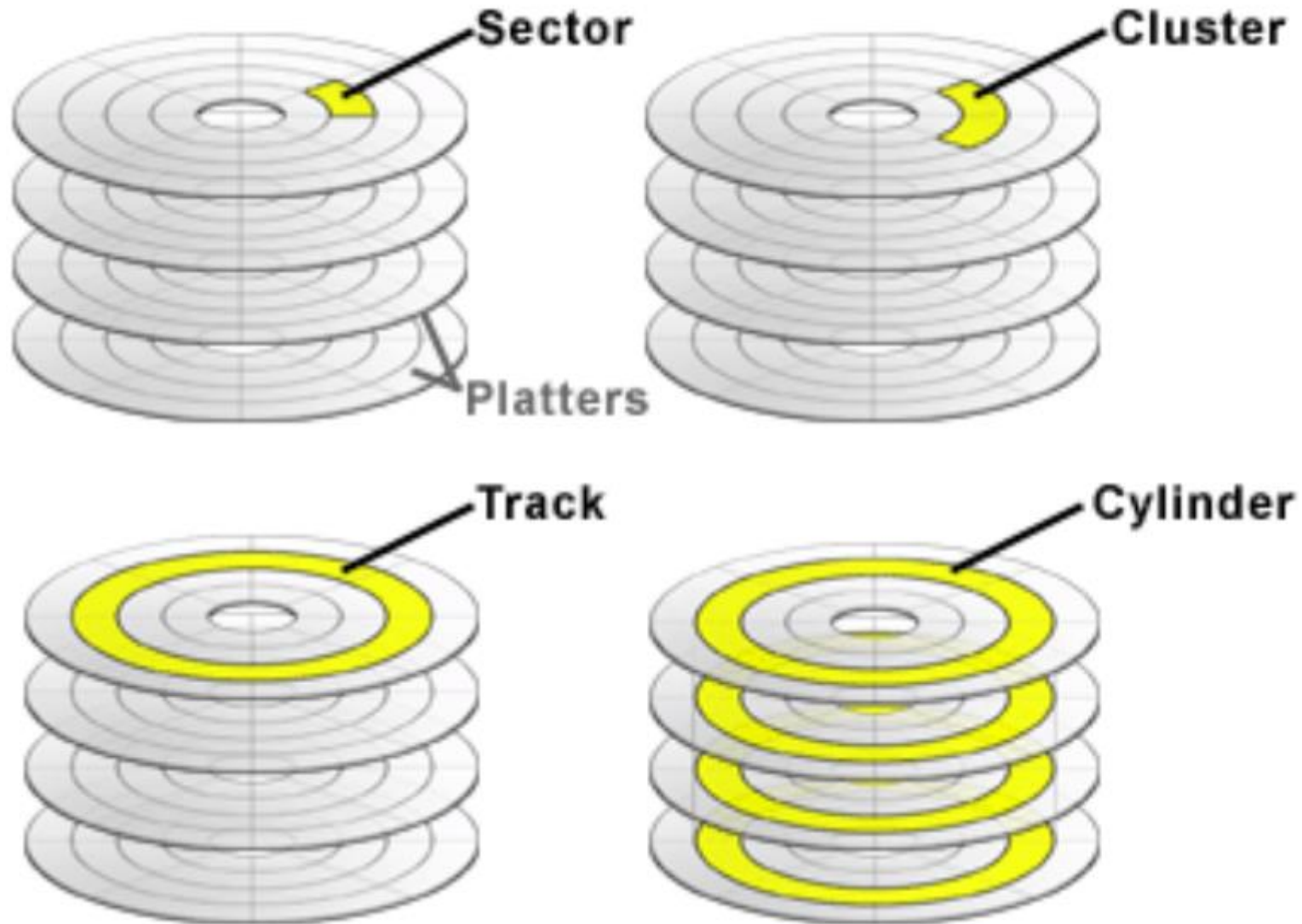# First Database Management Systems

COBOL



DBMS

# Hierarchical and Network Database Management Systems

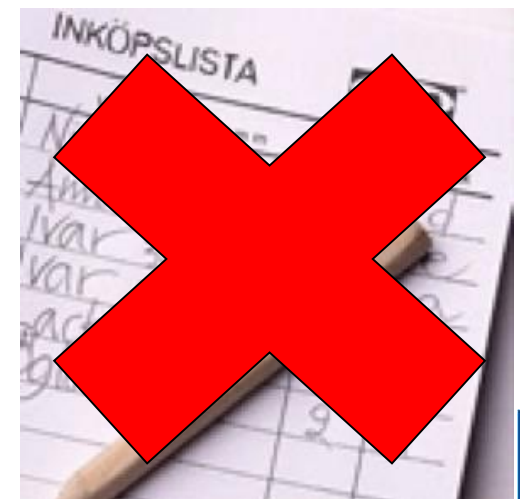# You had to know what data you want, and how to find it

# Early DBMS' were Disk-Aware

# Codd's Relational Model

Just tell me
the data you want,
the system will
find it.

# SystemR

```
CREATE TABLE Students(
id INT PRIMARY_KEY,
firstname VARCHAR(96),
lastname VARCHAR(96)
);
```

```
SELECT * FROM Students
WHERE id > 10;
```

?

Views

Relations

Structured Query
Language

Indexes

Disk Access
Methods

Disk

# Finding the Data using a Query Optimizer

**Each color represents a program in this plan diagram**



- Each program produces the same result for the Query.
- Each program has *different performance characteristics* depending on changes in the data characteristics

- Data Integrity using Transactions*



ACID

Atomicity Consistency Isolation Durability

*Jim Gray, "The Transaction Concept: Virtues and Limitation"

In the 1990s
Data Read Rates Increased Dramatically

# Distribute within a Data Center

**Master-Slave Replication**



Data-location awareness is back:
Clients read from slaves, write to master.
Possibility of reading stale data.

In the 2000s
Data Write Rates Increased Dramatically

# Unstructured Data explodes



Source: IDC whitepaper. As the Economy contracts, the Digital Universe Explodes. 2009

Key-Value stores don't do Big Data yet. Existing Big Data systems currently only work for a single Data Centre.*

*The usual Google Exception applies

Storage and Processing of Big Data

# What is Apache Hadoop?

- Huge data sets and large files
  - Gigabytes files, petabyte data sets
  - Scales to thousands of nodes on commodity hardware

- No Schema Required
  - Data can be just copied in, extract required columns later

- Fault tolerant

- Network topology-aware, Data Location-Aware

- Optimized for analytics: high-throughput file access

# Hadoop (version 1)

Application

MapReduce

Hadoop Filesystem

# HDFS: Hadoop Filesystem
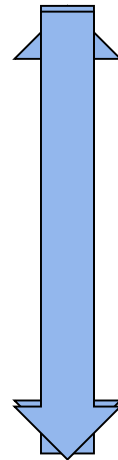
write "/crawler/bot/jd.io/1"

Under-replicated blocks

Name node

Heartbeats

Rebalance
Re-replicate
blocks

| 1 | 2 | 5 |
| 4 | 5 | 6 |

1

3

| 1 | 3 | 4 |
| 6 | | |

| 1 | 2 | 3 |
| 5 | | |

| 2 | 4 | 5 |
| | | 6 |

2

Data nodes

Data nodes

# HDFS v2 Architecture



Active-Standby Replication of NN Log
Agreement on the Active NameNode
Faster Recovery - Cut the NN Log

Journal Nodes

Zookeeper Nodes

NameNode

Standby NameNode

Snapshot Node

HDFS Client

DataNodes

30

# HopsFS Architecture



NDB

Load Balancer

HopsFS Client

HDFS Client

NameNodes

Leader

DataNodes

31

# Processing Big Data
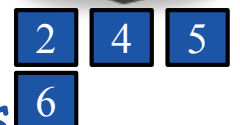
# Big Data Processing with No Data Locality

submit

Job("/crawler/bot/jd.io/1")

Workflow Manager

Compute Grid Node

Job

This doesn't scale.
Bandwidth is the bottleneck

1 2 3

2 5 6

4 3 6

3 5 6

1 2 4

1 4 5

# MapReduce – Data Locality

submit

Job("/crawler/bot/jd.io/1")

Job Tracker

| Task Tracker | Task Tracker | Task Tracker | Task Tracker | Task Tracker | Task Tracker |
|---|---|---|---|---|---|
| Job | Job | Job | Job | Job | Job |
| DN | DN | DN | DN | DN | DN |
| 1 2 3 | 2 5 6 | 4 3 6 | 3 5 6 | 1 2 4 | 1 4 5 |
| | R | | R | R | |

R = resultFile(s)

# MapReduce*

1. Programming Paradigm

2. Processing Pipeline (moving computation to data)

*Dean et al, OSDI'04

```
map(record) ->
        {(key_i, value_i), .., (key_l, value_l)}
```

```
reduce((key_i, {value_k, .., value_y}) -> output
```

# MapReduce Programming Paradigm

- Also found in:

   Functional programming languages

   MongoDB

   Cassandra

```
map(url, doc) ->
        {(term_i, url),(term_m, url)}


reduce((term,{url_k,..,url_y}) ->
        (term, (posting list of url, count))
```

# Example: Building a Web Search Index

```
map( ("jd.io", "A hipster website with news"))
->
{
    emit("a", "jd.io"),
    emit("hipster", "jd.io"),
    emit("website", "jd.io"),
    emit("with", "jd.io"),
    emit("news", "jd.io")
}
```

# Example: Building a Web Search Index

```
map( ("hn.io", "Hacker hipster news"))
->
{

    emit("hacker", "hn.io"),

    emit("hipster", "hn.io"),

    emit("news", "hn.io")

}
```

```
reduce( "hipster", { "jd.io", "hn.io" }) ->
      ( "hipster", (["jd.io", "hn.io"], 2))
```

```
reduce( "website", { "jd.io"}) ->
      ( "website", (["jd.io"], 1))
```

```
reduce( "news", { "jd.io", "hn.io" }) ->
        ( "news", (["jd.io", "hn.io"], 2))
```

# Map Phase

MapReduce

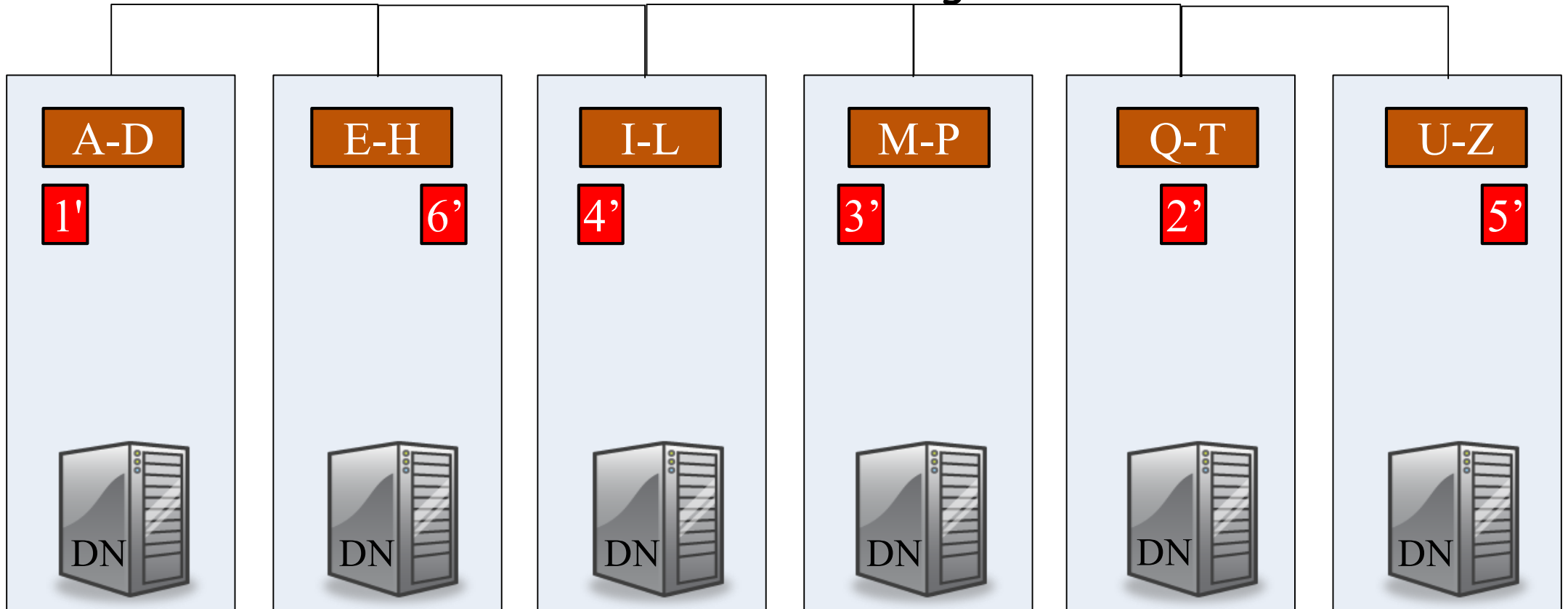`map(url, doc) -> {(term_i, url),(term_l, url)}`
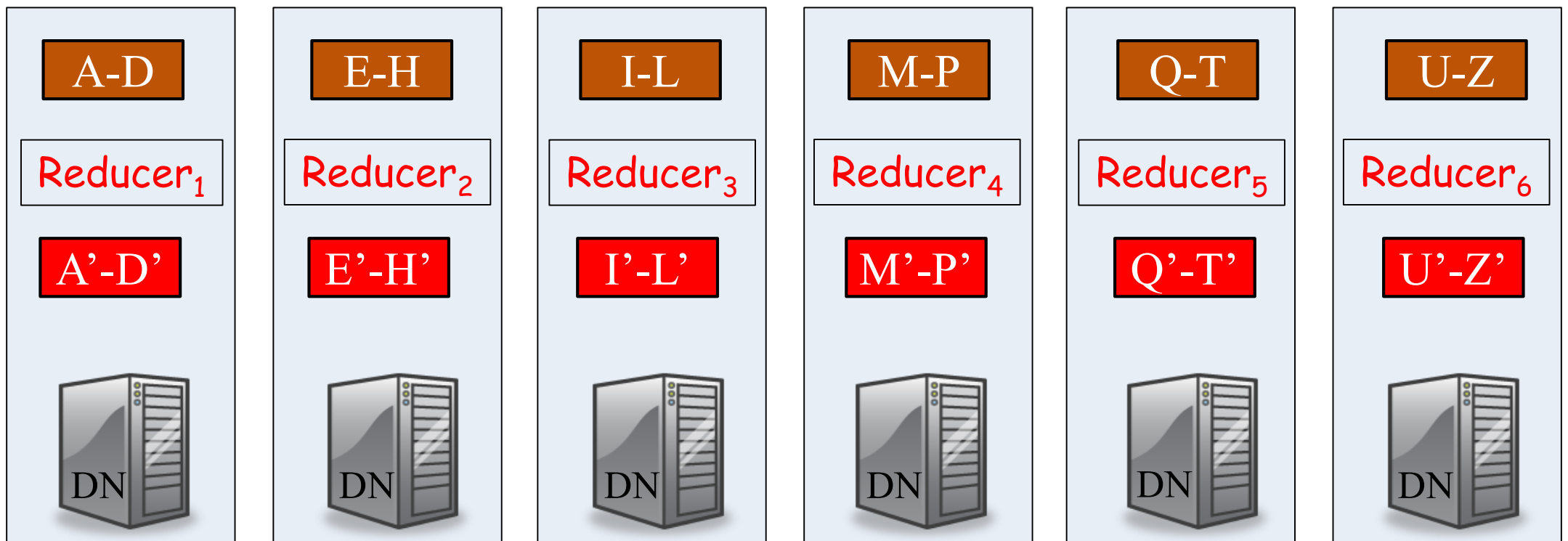
# Shuffle Phase

MapReduce

group by term

Shuffle over the Network using a Partitioner

| A-D | E-H | I-L | M-P | Q-T | U-Z |
|---|---|---|---|---|---|
| 1' | 6' | 4' | 3' | 2' | 5' |
| DN | DN | DN | DN | DN | DN |

# Reduce Phase

MapReduce

```
reduce((term,{url_k,url_y}) ->
    (term, (posting list of url, count))
```

| A-D | E-H | I-L | M-P | Q-T | U-Z |
|-----|-----|-----|-----|-----|-----|
| Reducer$_1$ | Reducer$_2$ | Reducer$_3$ | Reducer$_4$ | Reducer$_5$ | Reducer$_6$ |
| A'-D' | E'-H' | I'-L' | M'-P' | Q'-T' | U'-Z' |
| DN | DN | DN | DN | DN | DN |

# Hadoop 2.x

**Single Processing Framework**
Batch Apps

**Multiple Processing Frameworks**
Batch, Interactive, Streaming ...

## Hadoop 1.x

MapReduce
(resource mgmt, job scheduler, data processing)

HDFS
(distributed storage)

## Hadoop 2.x

MapReduce
(data processing)

Others
(spark, mpi, giraph, etc)

YARN
(resource mgmt, job scheduler)
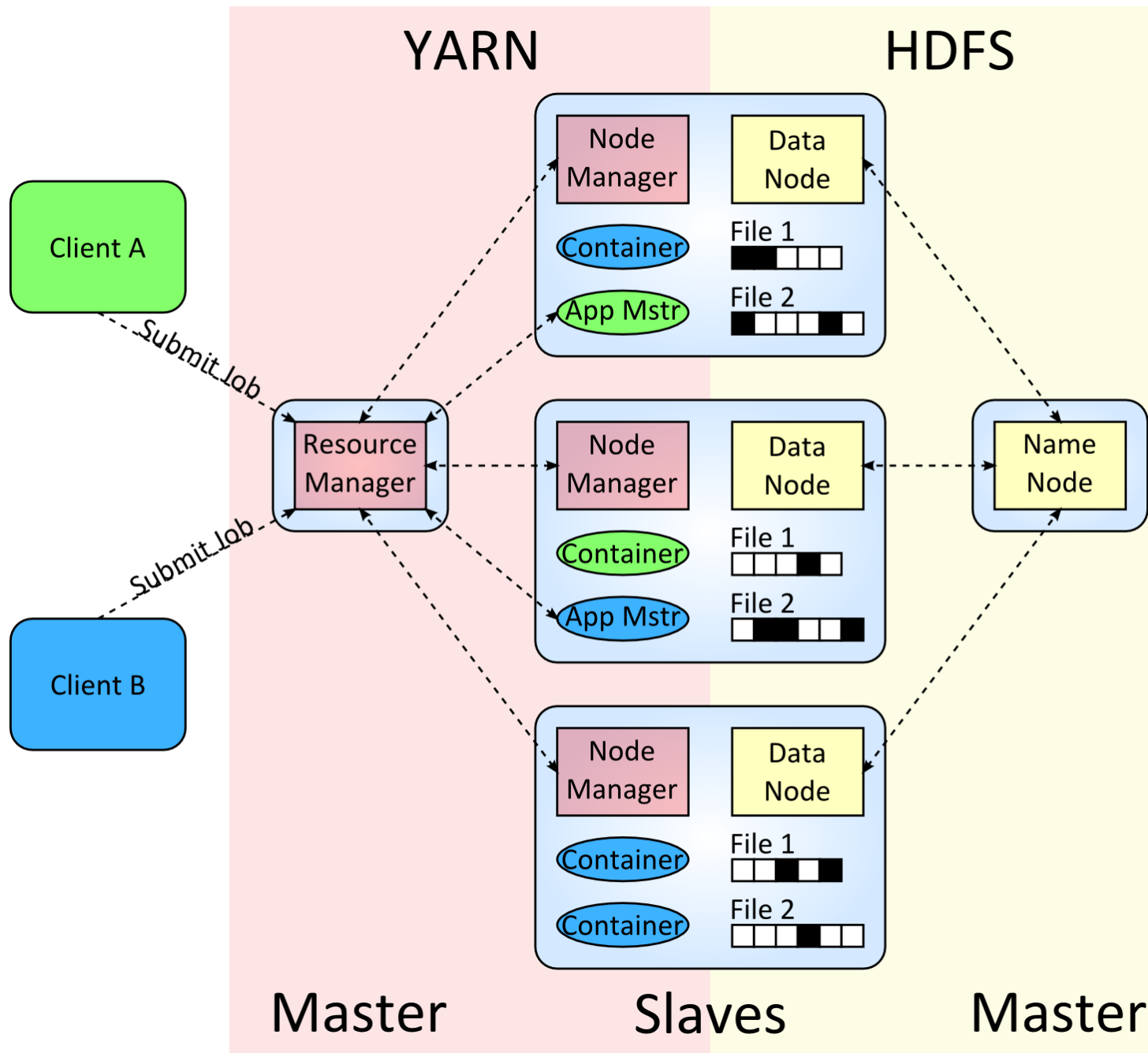
HDFS
(distributed storage)
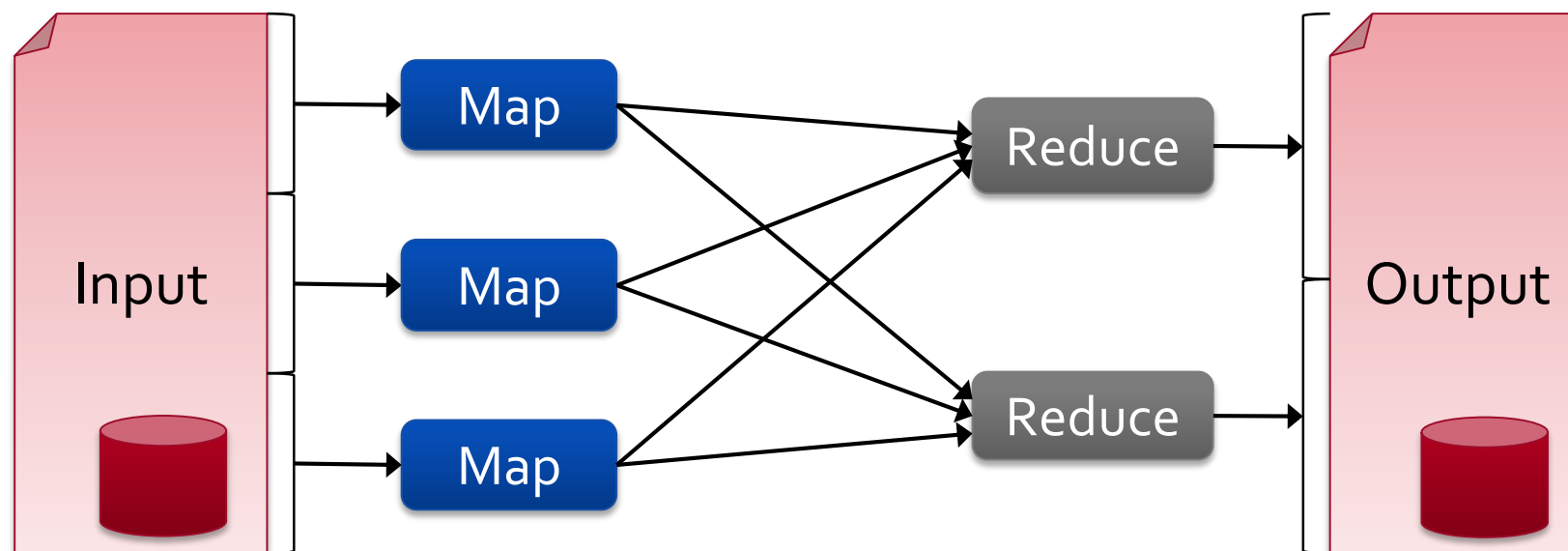
# MapReduce and MPI as YARN Applications



[Murthy et. al, Apache Hadoop YARN: Yet Another Resource Negotiator", SOCC'13]

# Data Locality in Hadoop v2

# Limitations of MapReduce [Zaharia'11]

- MapReduce is based on an *acyclic data flow* from stable storage to stable storage.
  - Slow writes data to HDFS at every stage in the pipeline
- Acyclic data flow is inefficient for applications that repeatedly reuse a *working set* of data:
  - **Iterative** algorithms (machine learning, graphs)
  - **Interactive** data mining tools (R, Excel, Python)
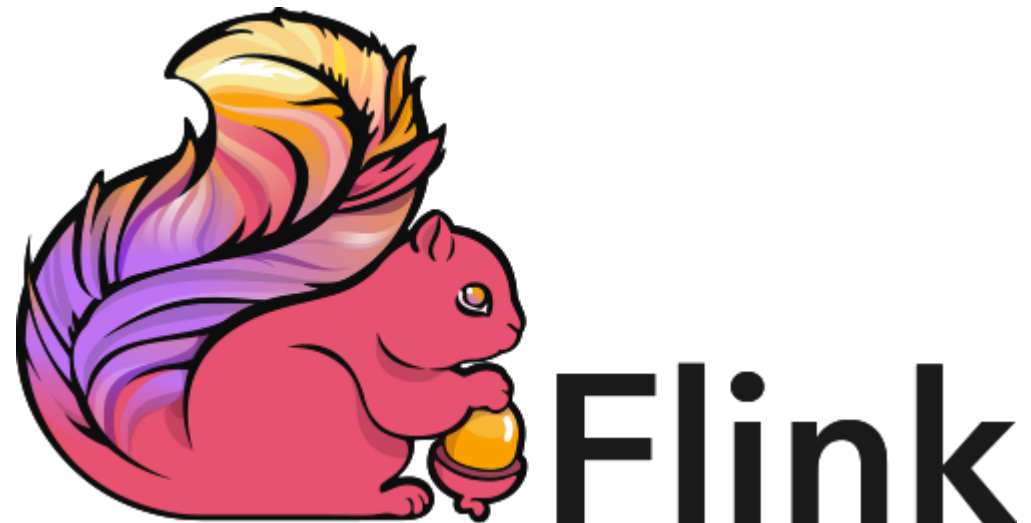
# Iterative Data Processing Frameworks

```scala
val input= TextFile(textInput)

val words = input
    .flatMap
      { line => line.split(" ") }


val counts = words
    .groupBy
        { word => word }
    .count()


val output = counts
    .write (wordsOutput,
        RecordDataSinkFormat() )

val plan = new ScalaPlan(Seq(output))
```

# Spark – Resiliant Distributed Datasets

- Allow apps to keep working sets in memory for efficient reuse

- Retain the attractive properties of MapReduce
  - Fault tolerance, data locality, scalability

Resilient distributed datasets (RDDs)
  - Immutable, partitioned collections of objects
  - Created through parallel *transformations* (map, filter, groupBy, join, …) on data in stable storage
  - Can be *cached* for efficient reuse

*Actions* on RDDs
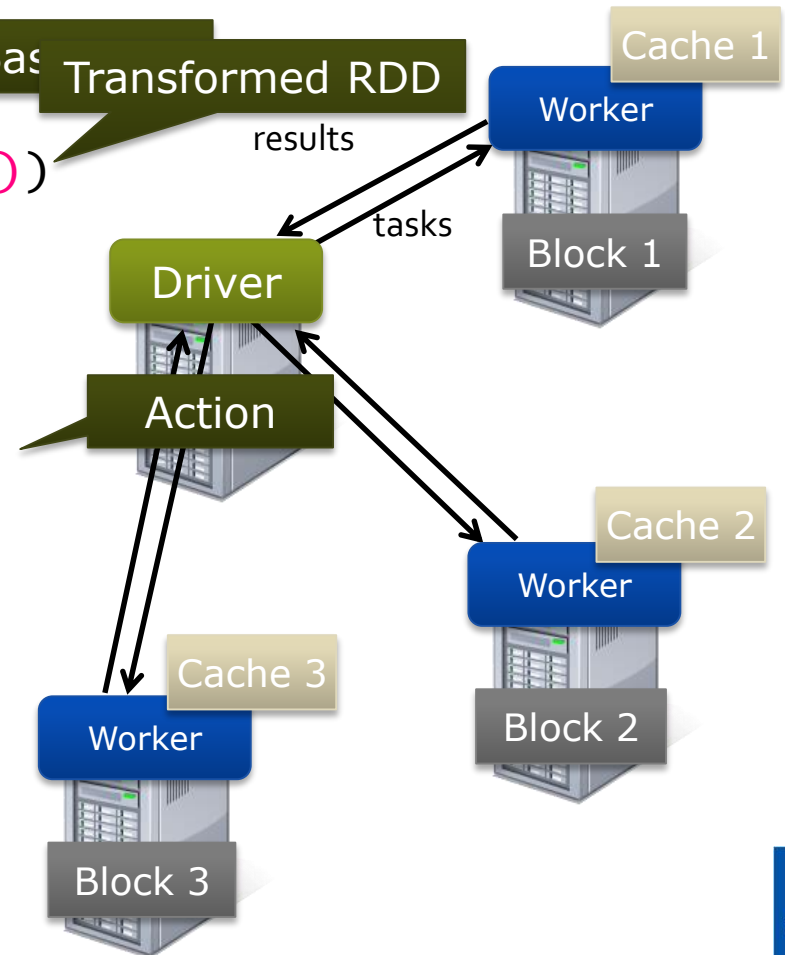  - Count, reduce, collect, save, …

# Example: Log Mining

Load error messages from a log into memory, then interactively search for various patterns

```
lines = spark.textFile("hdfs://...")
errors = lines.filter(_.startsWith("ERROR"))
messages = errors.map(_.split('\t')(2))
cachedMsgs = messages.cache()


cachedMsgs.filter(_.contains("foo")).count
cachedMsgs.filter(_.contains("bar")).count
. . .
```

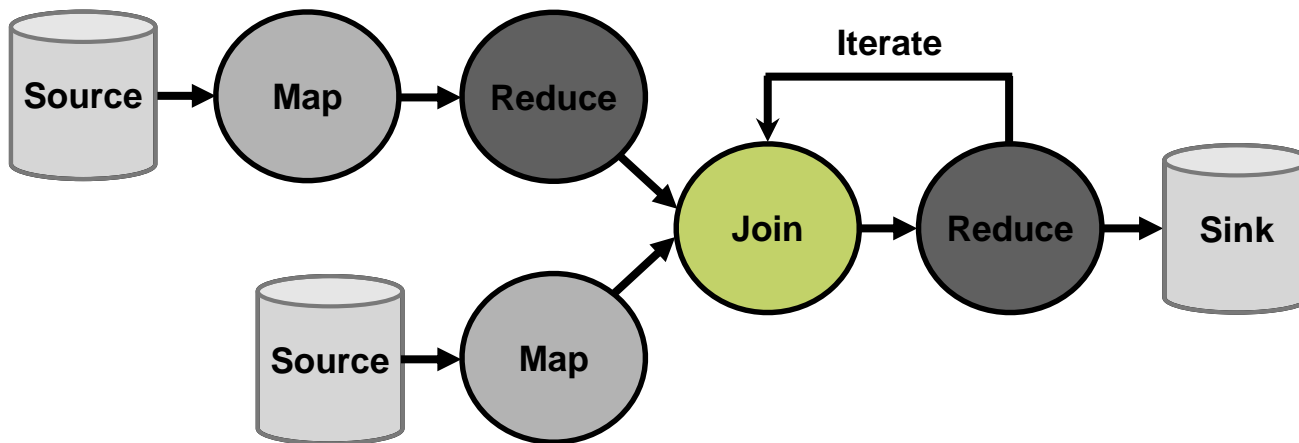**Result:** scaled to 1 TB data in 5-7 sec
(vs 170 sec for on-disk data)

Base RDD
Transformed RDD
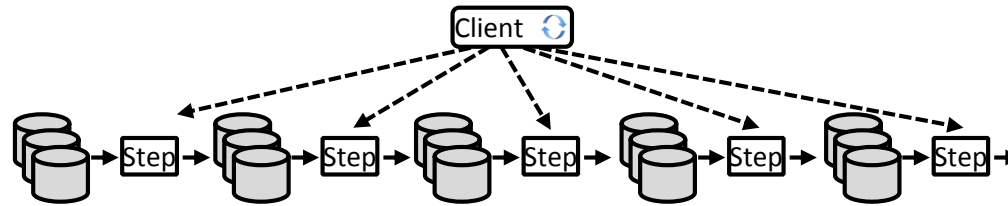results
tasks

Driver

Action

Worker

Cache 1
Block 1

Worker
Cache 2
Block 2

Worker
Cache 3
Block 3

# Apache Flink – DataFlow Operators

Flink

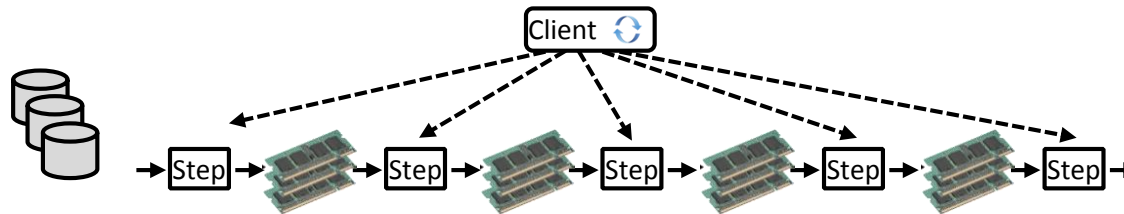| Map | Iterate | Project |
|---|---|---|
| Reduce | Delta Iterate | Aggregate |
| Join | Filter | Distinct |
| CoGroup | FlatMap | Vertex Update |
| Union | GroupReduce | Accumulators |



*Alexandrov et al.: "The Stratosphere Platform for Big Data Analytics," VLDB Journal 5/2014
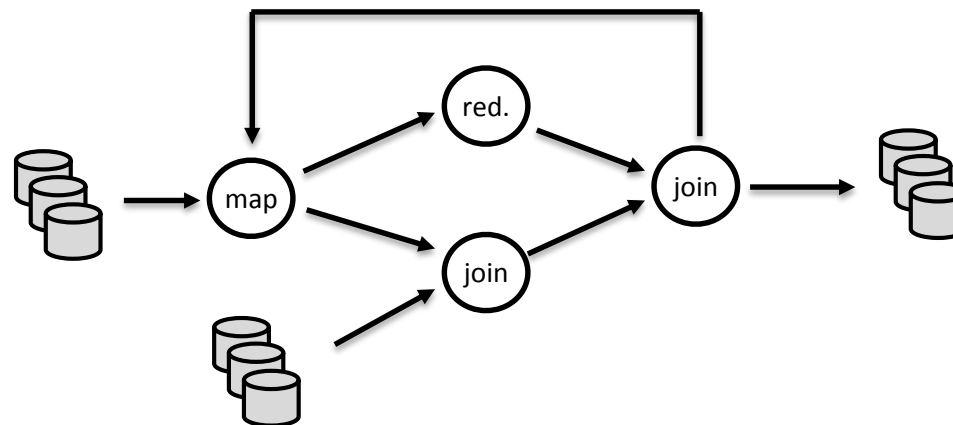
# Built-in vs. driver-based looping

Loop outside the system, in driver program

Iterative program looks like many independent jobs

Dataflows with feedback edges

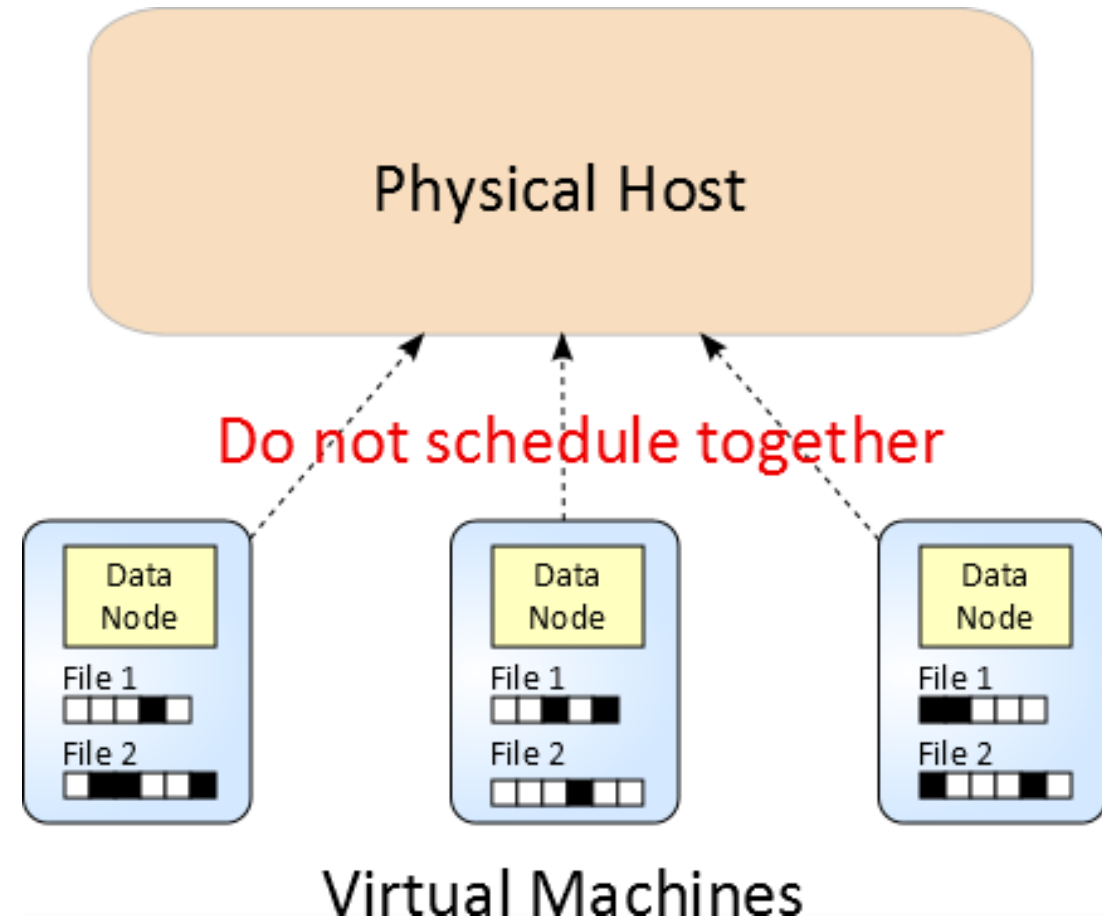System is iteration-aware, can optimize the job

# Hadoop on the Cloud

- Cloud Computing traditionally separates storage and computation.

Amazon Web Services
OpenStack

Nova (Compute) (EC2)

Glance (VM Images) (S3)

Swift (Object Storage)
Elastic Block Storage

# Data Locality for Hadoop on the Cloud

- Cloud hardware configurations should support data locality

- Hadoop's original topology awareness breaks
  - Placement of >1 VM containing block replicas for the same file on the same physical host increases correlated failures

- VMWare introduced a NodeGroup aware topology
  - HADOOP-8468

# Conclusions

- Hadoop is the open-source enabling technology for Big Data

- YARN is rapidly becoming the operating system for the Data Center

- Apache Spark and Flink are in-memory processing frameworks for Hadoop

# References

- Dean et. Al, "MapReduce: Simplified Data Processing on Large Clusters", OSDI'04.

- Schvachko, "HDFS Scalability: The limits to growth", Usenix, :login, April 2010.

- Murthy et al, "Apache Hadoop YARN: Yet Another Resource Negotiator", SOCC'13.

- "Processing a Trillion Cells per Mouse Click", VLDB'12