# "Hadoop"
## A Distributed Architecture, FileSystem, & MapReduce

**H. Andrew Schwartz**

CSE545
Spring 2023

(freesvg.org/1534373472)

# Big Data Analytics, The Class

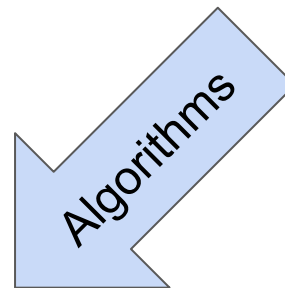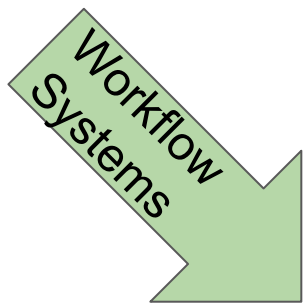**Goal:** Generalizations
A *model* or *summarization* of the data.

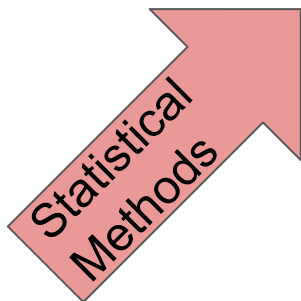*Data Workflow Frameworks*

*Analytics and Algorithms*

**Hadoop File System**
Spark
Streaming
**MapReduce**
Deep Learning Frameworks

Similarity Search
Hypothesis Testing
Transformers/Self-Supervision
Recommendation Systems
Link Analysis

# Big Data Analytics, The Class
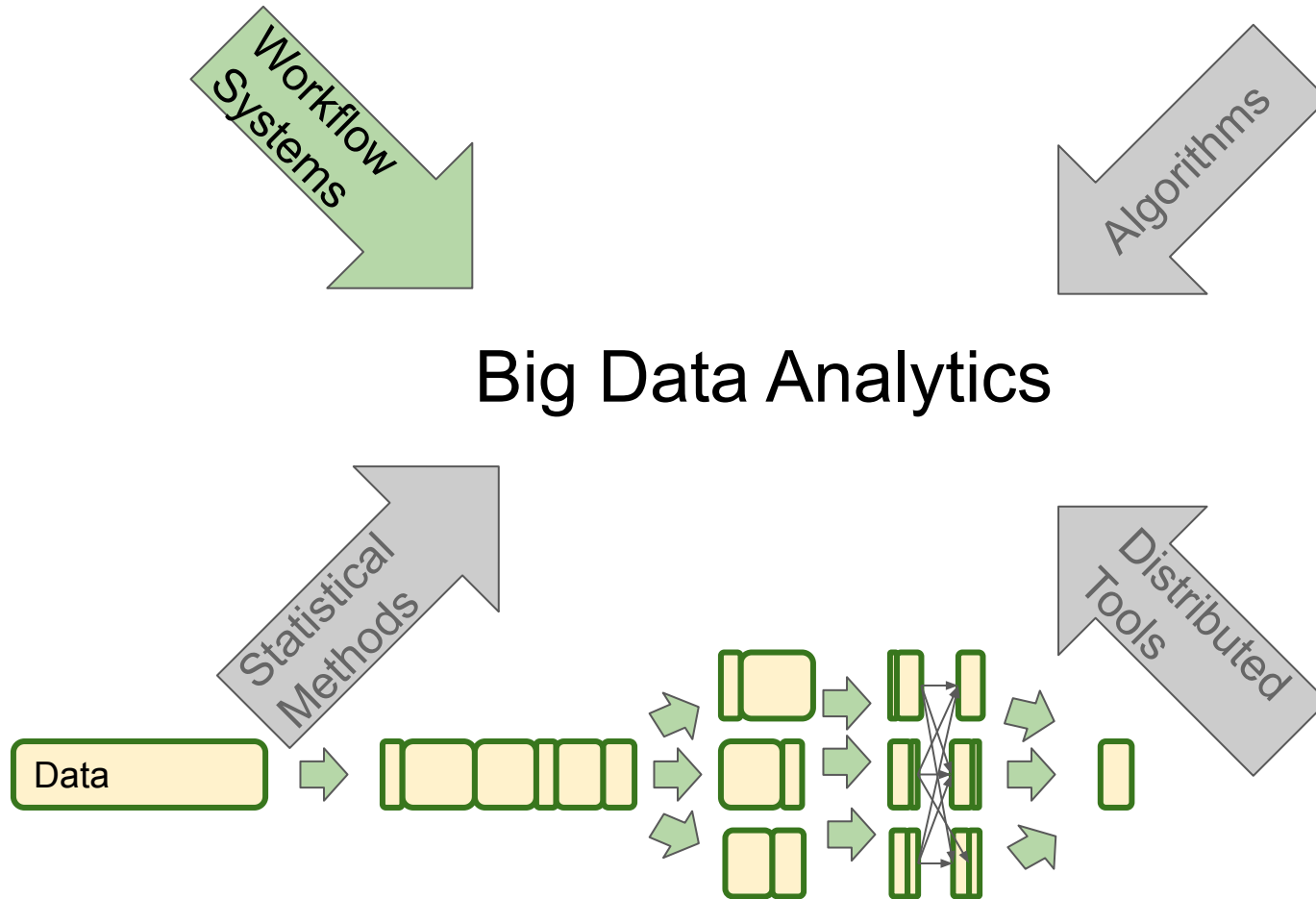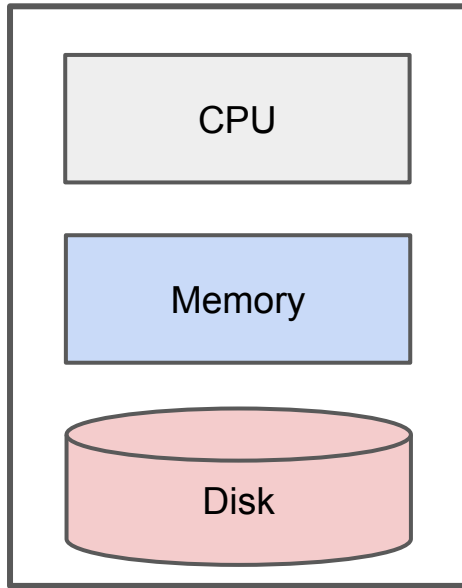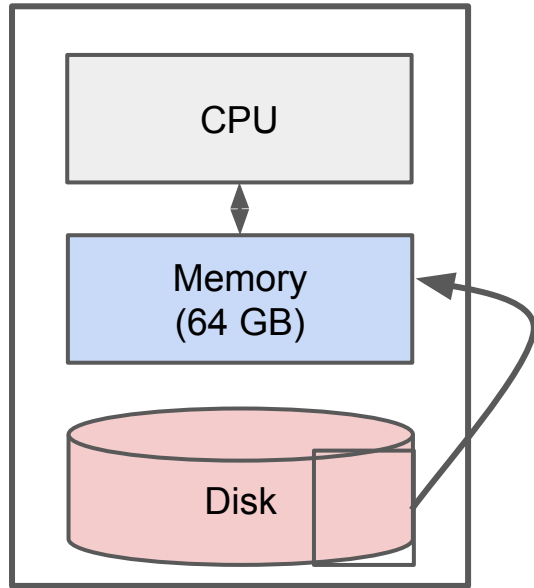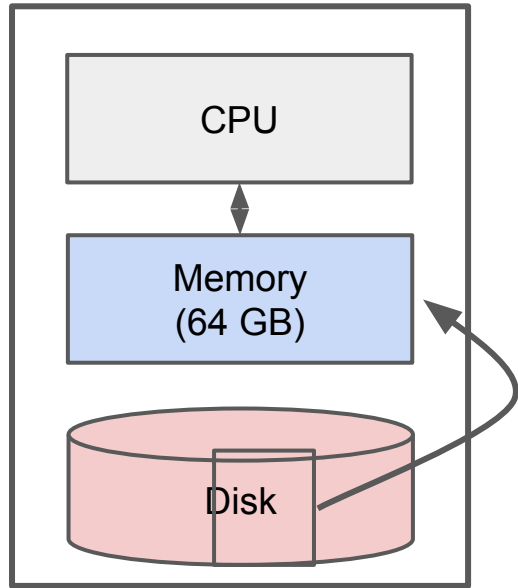
# Classical Data Analytics

# Classical Data Analytics

# Classical Data Analytics

# Classical Data Analytics

# IO Bounded

Reading a word from disk versus main memory: $10^5$ slower!

Reading many contiguously stored words
is faster per word, but fast modern disks
still only reach ~1GB/s for sequential reads.

# IO Bounded

Reading a word from disk versus main memory: $10^5$ slower!

Reading many contiguously stored words
is faster per word, but fast modern disks
still only reach ~1GB/s for sequential reads.



IO Bound: biggest performance bottleneck is reading / writing to disk.

starts around 500 GBs: >10 minutes just to read

500 TBs: ~8,600 minutes = ~6 days

# Classical Big Data



**Classical focus:** efficient use of disk.
e.g. Apache Lucene / Solr

**Classical limitation:** Still bounded when needing to process all of a large file.

# Classical Big Data

*How to solve?*

**Classical limitation:** Still bounded when needing to process all of a large file.

# Distributed Architecture

# Distributed Architecture

In reality, modern setups often have multiple cpus and disks per server, but we will model as if one machine per cpu-disk pair.

# Distributed Architecture (Cluster)

# Distributed Architecture (Cluster)

Challenges for IO Cluster Computing

1. Nodes fail

   1 in 1000 nodes fail a day

2. Network is a bottleneck

   Typically 1-10 Gb/s throughput

3. Traditional distributed programming is
   often ad-hoc and complicated

# Distributed Architecture (Cluster)

Challenges for IO Cluster Computing

1. Nodes fail

   1 in 1000 nodes fail a day

   Duplicate Data

2. Network is a bottleneck

   Typically 1-10 Gb/s throughput

   Bring computation to nodes, rather than data to nodes.

3. Traditional distributed programming is often ad-hoc and complicated

   Stipulate a programming system that can easily be distributed

# Distributed Architecture (Cluster)

Challenges for IO Cluster Computing

1. Nodes fail

   1 in 1000 nodes fail a day

   Duplicate Data

2. Network is a bottleneck

   Typically 1-10 Gb/s throughput

   Bring computation to nodes, rather than data to nodes.

3. Traditional distributed programming is often ad-hoc and complicated

   Stipulate a programming system that can easily be distributed

*HDFS with MapReduce accomplishes all!*

# Distributed Filesystem

The effectiveness of MapReduce, Spark, and other distributed processing systems is in part simply due to use of a <u>distributed filesystem!</u>

# Distributed Filesystem

Characteristics for Big Data Tasks

Large files (i.e. >100 GB to TBs)

Reads are most common

No need to update in place
(append preferred)

CPU

Memory

Disk

# Distributed Filesystem

(e.g. Apache HadoopDFS, GoogleFS, EMRFS)

C, D: Two different files

C

D

# Distributed Filesy

(e.g. Apache **Hadoop**DFS, GoogleFS, EM

C, D: Two different files

C

D

# Distributed Filesystem

(e.g. Apache HadoopDFS, GoogleFS, EMRFS)

C, D: Two different files; break into chunks (or "partitions"):

| $C_0$ |
| $C_1$ |
| $C_2$ |
| $C_3$ |
| $C_4$ |
| $C_5$ |

| $D_0$ |
| $D_1$ |
| $D_2$ |
| $D_3$ |
| $D_4$ |
| $D_5$ |

# Distributed Filesystem

(e.g. Apache HadoopDFS, GoogleFS, EMRFS)

C, D: Two different files



chunk server 1     chunk server 2     chunk server 3     chunk server n

# Distributed Filesystem

(e.g. Apache HadoopDFS, GoogleFS, EMRFS)

C, D: Two different files



chunk server 1          chunk server 2          chunk server 3          chunk server n

(Leskovec at al., 2014; http://www.mmds.org/)

# Distributed Filesystem

(e.g. Apache HadoopDFS, GoogleFS, EMRFS)

C, D: Two different files



chunk server 1     chunk server 2     chunk server 3     chunk server n

(Leskovec at al., 2014; http://www.mmds.org/)

# Distributed Filesystem

## Chunk servers (on Data Nodes)

File is split into contiguous chunks

Typically each chunk is 16-64MB

Each chunk replicated (usually 2x or 3x)

Try to keep replicas in different racks

# Components of a Distributed Filesystem

## Chunk servers (on Data Nodes)

File is split into contiguous chunks

Typically each chunk is 16-64MB

Each chunk replicated (usually 2x or 3x)

Try to keep replicas in different racks

## Name node (aka master node)

Stores metadata about where files are stored

Might be replicated or distributed across data nodes.

# Components of a Distributed Filesystem

## Chunk servers (on Data Nodes)

File is split into contiguous chunks

Typically each chunk is 16-64MB

Each chunk replicated (usually 2x or 3x)

Try to keep replicas in different racks

## Name node (aka master node)

Stores metadata about where files are stored

Might be replicated or distributed across data nodes.

## Client library for file access

Talks to master to find chunk servers

Connects directly to chunk servers to access data

(Leskovec at al., 2014; http://www.mmds.org/)

# Distributed Architecture (Cluster)

Challenges for IO Cluster Computing

1. Nodes fail

   1 in 1000 nodes fail a day

   Duplicate Data **(Distributed FS)**

2. Network is a bottleneck

   Typically 1-10 Gb/s throughput

   Bring computation to nodes, rather than data to nodes.

3. Traditional distributed programming is often ad-hoc and complicated

   Stipulate a programming system that can easily be distributed

# What is MapReduce

*noun.1*  -  **A *style of programming***

input chunks => map tasks   |     group_by keys   |     reduce tasks => output

"|"  is the linux "pipe" symbol: passes stdout from first process to stdin of next.

# What is MapReduce

*noun.1*  -  **A *style of programming***

input chunks => map tasks  |    group_by keys  |    reduce tasks => output

"|"  is the linux "pipe" symbol: passes stdout from first process to stdin of next.

E.g. counting words:

```
tokenize(document) | sort | uniq -c
```

# What is MapReduce

*noun.1* - A *style of programming*

input chunks  |  map tasks  |  group_by keys  |  reduce tasks => output

"|"  is the linux "pipe" symbol: passes output from first process to input of next.

E.g. counting words:

```
cat file.txt | tr -s '[[:space:]]' '\n' | sort | uniq -c
```

*noun.2* - A *system* that distributes MapReduce style programs across a distributed file-system.

(e.g. Google's internal "MapReduce" or apache.hadoop.mapreduce with hdfs)

# What is MapReduce

*noun.1*  -  **A *style of programming***

input chunks => map tasks  |     group_by keys   |     reduce tasks => output

"|"  is the linux "pipe" symbol: passes output from first process to input of next.

E.g. counting words:

```
tokenize(document) | sort | uniq -c
```

*noun.2*  - **A *system* that distributes MapReduce style programs across a distributed file-system.**

(e.g. Google's internal "MapReduce" or apache.hadoop.mapreduce with hdfs)

# What is MapReduce

# What is MapReduce



Input chunks → Map → Key−value pairs (k,v) → Group by keys → Keys with all their values (k, [v, w,...]) → Reduce tasks → Combined output

extract what you care about.

line => (k, v)

Map

# What is MapReduce

# What is MapReduce

# What is MapReduce

Easy as 1, 2, 3!

Step 1: **Map**        Step 2: **Sort / Group by**        Step 3: **Reduce**

# What is MapReduce

Easy as 1, 2, 3!

Step 1: **Map**     Step 2: **Sort / Group by**     Step 3: **Reduce**

# (1) The *Map* Step



(Leskovec at al., 2014; http://www.mmds.org/)

# (2) The *Sort / Group-by* Step

# (3) The *Reduce* Step



(Leskovec at al., 2014; http://www.mmds.org/)

# What is MapReduce

Easy as 1, 2, 3!

Step 1: **Map**    Step 2: **Sort / Group by**    Step 3: **Reduce**

# What is MapReduce

Map: $(k,v)$ -> $(k', v')$*
   (Written by programmer)

Group by key: $(k_1', v_1')$, $(k_2', v_2')$, ... -> $(k_1', (v_1', v', …),$
   (system handles)                              $(k_2', (v_1', v', …), …$

Reduce: $(k', (v_1', v', …))$ -> $(k', v'')$*
   (Written by programmer)

# Example: Word Count

```
tokenize(document) | sort | uniq -c
```

# Example: Word Count

```
tokenize(document) | sort | uniq -c
```

Map: extract what you care about.

sort and shuffle

Reduce: aggregate, summarize

# Example: Word Count



The crew of the space shuttle Endeavor recently returned to Earth as ambassadors, harbingers of a new era of space exploration. Scientists at NASA are saying that the recent assembly of the Dextre bot is the first step in a long-term space-based man/mache partnership. '"The work we're doing now -- the robotics we're doing - - is what we're going to need ………………..

**Big document**

(Leskovec at al., 2014; http://www.mmds.org/)

**Provided by the programmer**

**MAP:**
Read input and produces a set of key-value pairs

The crew of the space shuttle Endeavor recently returned to Earth as ambassadors, harbingers of a new era of space exploration. Scientists at NASA are saying that the recent assembly of the Dextre bot is the first step in a long-term space-based man/mache partnership. '"The work we're doing now -- the robotics we're doing - - is what we're going to need ………………..

(The, 1)
(crew, 1)
(of, 1)
(the, 1)
(space, 1)
(shuttle, 1)
(Endeavor, 1)
(recently, 1)
….

**Big document**      **(key, value)**

## Provided by the programmer

**MAP:**
Read input and produces a set of key-value pairs

**Group by key:**
Collect all pairs with same key

The crew of the space shuttle Endeavor recently returned to Earth as ambassadors, harbingers of a new era of space exploration. Scientists at NASA are saying that the recent assembly of the Dextre bot is the first step in a long-term space-based man/mache partnership. "'The work we're doing now -- the robotics we're doing - - is what we're going to need ...........................

(The, 1)
(crew, 1)
(of, 1)
(the, 1)
(space, 1)
(shuttle, 1)
(Endeavor, 1)
(recently, 1)

....

(crew, 1)
(crew, 1)
(space, 1)
(the, 1)
(the, 1)
(the, 1)
(shuttle, 1)
(recently, 1)

...

**Big document**      **(key, value)**      **(key, value)**

(Leskovec at al., 2014; http://www.mmds.org/)

Chunks

**Provided by the programmer**

**MAP:** Read input and produces a set of key-value pairs

**Group by key:** Collect all pairs with same key

**Provided by the programmer**

**Reduce:** Collect all values belonging to the key and output

The crew of the space shuttle Endeavor recently returned to Earth as ambassadors, harbingers of a new era of space exploration. Scientists at NASA are saying that the recent assembly of the Dextre bot is the first step in a long-term space-based man/mache partnership. '"The work we're doing now -- the robotics we're doing - - is what we're going to need ...................

**Big document**

(The, 1)
(crew, 1)
(of, 1)
(the, 1)
(space, 1)
(shuttle, 1)
(Endeavor, 1)
(recently, 1)
....

**(key, value)**

(crew, 1)
(crew, 1)
(space, 1)
(the, 1)
(the, 1)
(the, 1)
(shuttle, 1)
(recently, 1)
...

**(key, value)**

(crew, 2)
(space, 1)
(the, 3)
(shuttle, 1)
(recently, 1)
...

**(key, value)**

Only sequential reads

# Example: Word Count

```python
@abstractmethod
def map(k, v):
    pass




@abstractmethod
def reduce(k, vs):
    pass
```

# Example: Word Count (v1)

```python
def map(k, v):
    for w in tokenize(v):
        yield (w,1)




def reduce(k, vs):
    return len(vs)
```

# Example: Word Count (v1)

```
def map(k, v):
    for w in tokenize(v):
        yield (w,1)
```

```
def tokenize(s):
    #simple version
    return s.split(' ')
```

```
def reduce(k, vs):
    return len(vs)
```

# Example: Word Count (v2)

```
def map(k, v):
    counts = dict()
    for w in tokenize(v):
```

counts each word within the chunk
(try/except is faster than
"if w in counts")

# Example: Word Count (v2)

```python
def map(k, v):
    counts = dict()
    for w in tokenize(v):
        try:
            counts[w] += 1
        except KeyError:
            counts[w] = 1
    for item in counts.iteritems():
        yield item
```

counts each word within the chunk (try/except is faster than "if w in counts")

# Example: Word Count (v2)

```python
def map(k, v):
    counts = dict()
    for w in tokenize(v):
        try:
            counts[w] += 1
        except KeyError:
            counts[w] = 1
    for item in counts.iteritems():
        yield item
```

counts each word within the chunk (try/except is faster than "if w in counts")

```python
def reduce(k, vs):
    return (k, sum(vs) )
```

sum of counts from different chunks

# Distributed Architecture (Cluster)

Challenges for IO Cluster Computing

1. Nodes fail

   1 in 1000 nodes fail a day

   Duplicate Data **(Distributed FS)** ✅

2. Network is a bottleneck

   Typically 1-10 Gb/s throughput

   Bring computation to nodes, rather than data to nodes.

3. Traditional distributed programming is often ad-hoc and complicated

   Stipulate a programming system that can easily be distributed

# Distributed Architecture (Cluster)

Challenges for IO Cluster Computing

1. Nodes fail

   1 in 1000 nodes fail a day

   Duplicate Data **(Distributed FS)**

2. Network is a bottleneck

   Typically 1-10 Gb/s throughput

   Bring computation to nodes, rather than data to nodes. **(Sort and Shuffle)**

3. Traditional distributed programming is often ad-hoc and complicated

   Stipulate a programming system that can easily be distributed

# Distributed Architecture (Cluster)

Challenges for IO Cluster Computing

1. Nodes fail

   1 in 1000 nodes fail a day

   Duplicate Data **(Distributed FS)**

2. Network is a bottleneck

   Typically 1-10 Gb/s throughput

   Bring computation to nodes, rather than data to nodes. **(Sort and Shuffle)**

3. Traditional distributed programming is often ad-hoc and complicated **(Simply define a map and reduce)**

   Stipulate a programming system that can easily be distributed

# Example: Relational Algebra

Select

Project

Union, Intersection, Difference

Natural Join

Grouping

# Example: Relational Algebra

**Select**

Project

Union, Intersection, Difference

**Natural Join**

Grouping

# Example: Relational Algebra

**Select**

$R(A_1, A_2, A_3, ...)$, Relation  $R$, Attributes $A_*$

return only those attribute tuples where condition $C$ is true

# Example: Relational Algebra

**Select**

$R(A_1, A_2, A_3, ...)$, Relation $R$, Attributes $A_*$

return only those attribute tuples where condition $C$ is true

```
def map(k, v): #v is list of attribute tuples: [(...,), (...,), ...]
    r = []
    for t in v:
        if t satisfies C:
            r += [(t, t)]
    return r
```

# Example: Relational Algebra

**Select**

*R(A$_1$,A$_2$,A$_3$,...)*, Relation  *R,* Attributes *A$_*$*
return only those attribute tuples where condition *C* is true

```python
def map(k, v): #v is list of attribute tuples: [(...,), (...,), ...]
    r = []
    for t in v:
        if t satisfies C:
            r += [(t, t)]
    return r
                def reduce(k, vs):
                  r = []
                  for each v in vs:
                    r += [(k, v)]
                  return r
```

# Example: Relational Algebra

**Select**

$R(A_1, A_2, A_3, ...)$, Relation $R$, Attributes $A_*$

return only those attribute tuples where condition $C$ is true

```
def map(k, v): #v is list of attribute tuples
    for t in v:
        if t satisfies C:
            yield (t, t)


def reduce(k, vs):
    For each v in vs:
        yield  (k, v)
```

## Natural Join

Given $R_1$ and $R_2$ return $R_{join}$
         -- union of all pairs of tuples that match given attributes.

```
def map(k, v): #k \in {R1, R2}, v is (A, B) for R1, (B, C) for R2
               #B are matched attributes
```

# Example: Relational Algebra

## Natural Join

Given $R_1$ and $R_2$ return $R_{join}$
                       -- union of all pairs of tuples that match given attributes.

```
def map(k, v): #k \in {R1, R2}, v is (A, B) for R1, (B, C) for R2
               #B are matched attributes
    if k=='R1':
        (a, b) = v
        return (b,('R₁',a))
    if k=='R2':
        (b,c) = v
        return (b,('R₂',c))
```

# Example: Relational Algebra

## Natural Join

Given $R_1$ and $R_2$ return $R_{join}$

                 -- union of all pairs of tuples that match given attributes.

```
def map(k, v): #k \in {R1, R2}, v is (A, B) for R1, (B, C) for R2
               #B are matched attributes

    if k=='R1':
        (a, b) = v
        return (b,('R₁',a))
    if k=='R2':
        (b,c) = v
        return (b,('R₂',c))
```

```
def reduce(k, vs):
    r1, r2, rjn = [], [], []
    for (s, x) in vs: #separate rs
        if s == 'R1': r1.append(x)
        else: r2.append(x)
    for a in r1: #join as tuple
        for each c in r2:
            rjn += ('R_join', (a, k, c)) #k is b
    return rjn
```

# Data Flow



**MAP:**
Read input and produces a set of key-value pairs

**Group by key:**
Collect all pairs with same key
(Hash merge, Shuffle, Sort, Partition)

**Reduce:**
Collect all values belonging to the key and output

Input | Big document
Intermediate | k1:v k1:v k2:v | k1:v | k3:v k4:v | k4:v k5:v | k4:v | k1:v k3:v

Group by Key

Grouped | k1:v,v,v,v | k2:v | k3:v,v | k4:v,v,v | k5:v

Output

J. Leskovec, A. Rajaraman, J. Ullman: Mining of Massive Datasets, http://www.mmds.org

# Data Flow



(Leskovec at al., 2014; http://www.mmds.org/)

# Data Flow



(Leskovec at al., 2014; http://www.mmds.org/)

# Data Flow

DFS ➡ Map ➡ Map's Local FS ➡ Reduce ➡ DFS

# Data Flow

MapReduce system handles:

- Partitioning

- Scheduling map / reducer execution

- Group by key


- Restarts from node failures

- Inter-machine communication

# Data Flow

DFS ⟹ MapReduce ⟹ DFS

- Schedule map tasks near physical storage of chunk
- Intermediate results stored locally
- Master / Name Node coordinates

# Data Flow

DFS $\Rightarrow$ MapReduce $\Rightarrow$ DFS

- Schedule map tasks near physical storage of chunk
- Intermediate results stored locally
- Master / Name Node coordinates
  - Task status: idle, in-progress, complete
  - Receives location of intermediate results and schedules with reducer
  - Checks nodes for failures and restarts when necessary
    - All map tasks on nodes must be completely restarted
    - Reduce tasks can pickup with reduce task failed

DFS ⟹ MapReduce ⟹ DFS

- Schedule map tasks near physical storage of chunk
- Intermediate results stored locally
- Master / Name Node coordinates
  - Task status: idle, in-progress, complete
  - Receives location of intermediate results and schedules with reducer
  - Checks nodes for failures and restarts when necessary
    - All map tasks on nodes must be completely restarted
    - Reduce tasks can pickup with reduce task failed

DFS ⟹ MapReduce ⟹ DFS ⟹ MapReduce ⟹ DFS

Skew: The degree to which certain tasks end up taking much longer than others.

Handled with:

- More reducers (i.e. partitions) than reduce tasks
- More reduce tasks than nodes

**Key Question:** *How many Map and Reduce jobs?*

*M:* map tasks, *R:* reducer tasks

# Data Flow

**Key Question:** *How many Map and Reduce jobs?*

*M:* map tasks, *R:* reducer tasks

**Answer: 1)** If possible, one chunk per map task

*(maximizes flexibility for scheduling)*

       **2)** $M >> |nodes| \approx\approx |cores|$

*(better handling of node failures, better load balancing)*

       **3)** $R <= M$

*(reduces number of parts stored in DFS)*

| CPU | CPU | | CPU |
|---|---|---|---|
| Mem | Mem | . | Mem |
| | | . | |
| Disk | Disk | . | Disk |

# Data Flow

□ Tasks (Map Task or Reduce Task)

version 1: few reduce tasks
(same number of reduce tasks as nodes)



node1

node2

node3

node4

node5

→ time

tasks represented by
**time to complete task**
(some tasks take much longer)

# Data Flow

Tasks (Map Task or Reduce Task)

version 1: few reduce tasks
(same number of reduce tasks as nodes)

version 2: more reduce tasks
(more reduce tasks than nodes)



| node1 |
| node2 |
| node3 |
| node4 |
| node5 |

time

tasks represented by
**time to complete task**
(some tasks take much longer)

tasks represented by
**time to complete task**
(some tasks take much longer)

# Data Flow

Tasks (Map Task or Reduce Task)
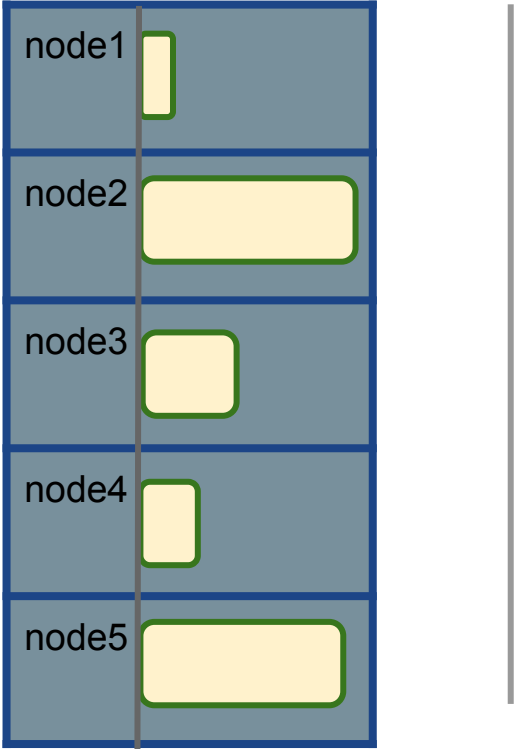
version 1: few reduce tasks
(same number of reduce tasks as nodes)

version 2: more reduce tasks
(more reduce tasks than nodes)



Last task completed

Can redistribute these tasks to other nodes

node1
node2
node3
node4
node5

time

tasks represented by
**time to complete task**
(some tasks take much longer)

tasks represented by
**time to complete task**
(some tasks take much longer)

(the last task now completes
much earlier )

# Communication Cost Model

How to assess performance?

(1) Computation: Map + Reduce + System Tasks

(2) Communication: Moving (key, value) pairs

# Communication Cost Model

How to assess performance?

(1) Computation: Map + Reduce + System Tasks

(2) Communication: Moving (key, value) pairs

Ultimate Goal: wall-clock Time.

# Communication Cost Model

How to assess performance?

**(1) Computation: Map + Reduce + System Tasks**

- Mappers and reducers often single pass O(n) within node
- System: sort the keys is usually most expensive
- Even if map executes on same node, disk read usually dominates
- In any case, can add more nodes

(2) Communication: Moving key, value pairs

Ultimate Goal: Wall-clock Time.

# Communication Cost Model

## How to assess performance?

(1) Computation: Map + Reduce + System Tasks

**(2) Communication: Moving key, value pairs**

Often dominates computation.
- Connection speeds: 1-10 giga**bits** per sec;
  HD read: 50-150 giga**bytes** per sec
- Even reading from disk to memory typically takes longer than operating on the data.

# Communication Cost Model

How to assess performance?

**Communication Cost =** input size +
(sum of size of all map-to-reducer files)

**(2)  Communication: Moving key, value pairs**

Often dominates computation.
- Connection speeds: 1-10 giga**bits** per sec;
  HD read: 50-150 giga**bytes** per sec
- Even reading from disk to memory typically takes longer than operating on the data.

# Communication Cost Model

How to assess performance?

**Communication Cost =** input size +
(sum of size of all map-to-reducer files)

**(2) Communication: Moving key, value pairs**

Often dominates computation.
- Connection speeds: 1-10 giga**bits** per sec;
  HD read: 50-150 giga**bytes** per sec
- Even reading from disk to memory typically takes longer than operating on the data.
- Output from reducer ignored because it's either small (finished summarizing data) or being passed to another mapreduce job.

# Communication Cost: Natural Join

R, S: Relations (Tables)     *R(A, B) ⋈ S(B, C)*

**Communication Cost =**     input size +
                             (sum of size of all map-to-reducer files)

DFS ⇨ Map ⇨ LocalFS ⇨ Network ⇨ Reduce ⇨ DFS ⇨ ?

# Communication Cost: Natural Join

R, S: Relations (Tables)    $R(A, B) \bowtie S(B, C)$

**Communication Cost =** input size +
(sum of size of all map-to-reducer files)

```
def reduce(k, vs):
    r1, r2 = [], []
    for (rel, x) in vs: #separate rs
        if rel == 'R': r1.append(x)
        else: r2.append(x)
    for a in r1: #join as tuple
        for each c in r2:
            yield (R_join', (a, k, c)) #k is
```

```
def map(k, v):
    if k=="R1":
        (a, b) = v
        yield (b,(R_1,a))
    if k=="R2":
        (b,c) = v
        yield (b,(R_2,c))
```

$b$

# Communication Cost: Natural Join

R, S: Relations (Tables)     $R(A, B) \bowtie S(B, C)$

**Communication Cost =** input size +
(sum of size of all map-to-reducer files)

= |R1| + |R2| + (|R1| + |R2|)
= **$O(|R1| + |R2|)$**

```python
def map(k, v):
    if k=="R1":
        (a, b) = v
        yield (b,(R1,a))
    if k=="R2":
        (b,c) = v
        yield (b,(R2,c))
```

```python
def reduce(k, vs):
    r1, r2 = [], []
    for (rel, x) in vs: #separate rs
        if rel == 'R': r1.append(x)
        else: r2.append(x)
    for a in r1: #join as tuple
        for each c in r2:
            yield (Rjoin,, (a, k, c)) #k is
```

b

# MapReduce: Final Considerations

- Performance Refinements:
  - Combiners (like word count version 2 but done via reduce)




  - Backup tasks (aka speculative tasks)




  - Override partition hash function to organize data

# MapReduce: Final Considerations

- Performance Refinements:
  - Combiners (like word count version 2 but done via reduce)
    - Run reduce right after map from same node before passing to reduce (MapTask can execute)
    - Reduces communication cost but requires commutative reduce steps
  - Backup tasks (aka speculative tasks)



  - Override partition hash function to organize data

# MapReduce: Final Considerations

- Performance Refinements:
  - Combiners (like word count version 2 but done via reduce)
    - Run reduce right after map from same node before passing to reduce (MapTask can execute)
    - Reduces communication cost but requires commutative reduce steps
  - Backup tasks (aka speculative tasks)
    - Schedule multiple copies of tasks when close to the end to mitigate certain nodes running slow.

  - Override partition hash function to organize data
    E.g. instead of `hash(url)` use `hash(hostname(url))`