

# CSC 555 and DSC 333

## Mining Big Data

### Lecture 5

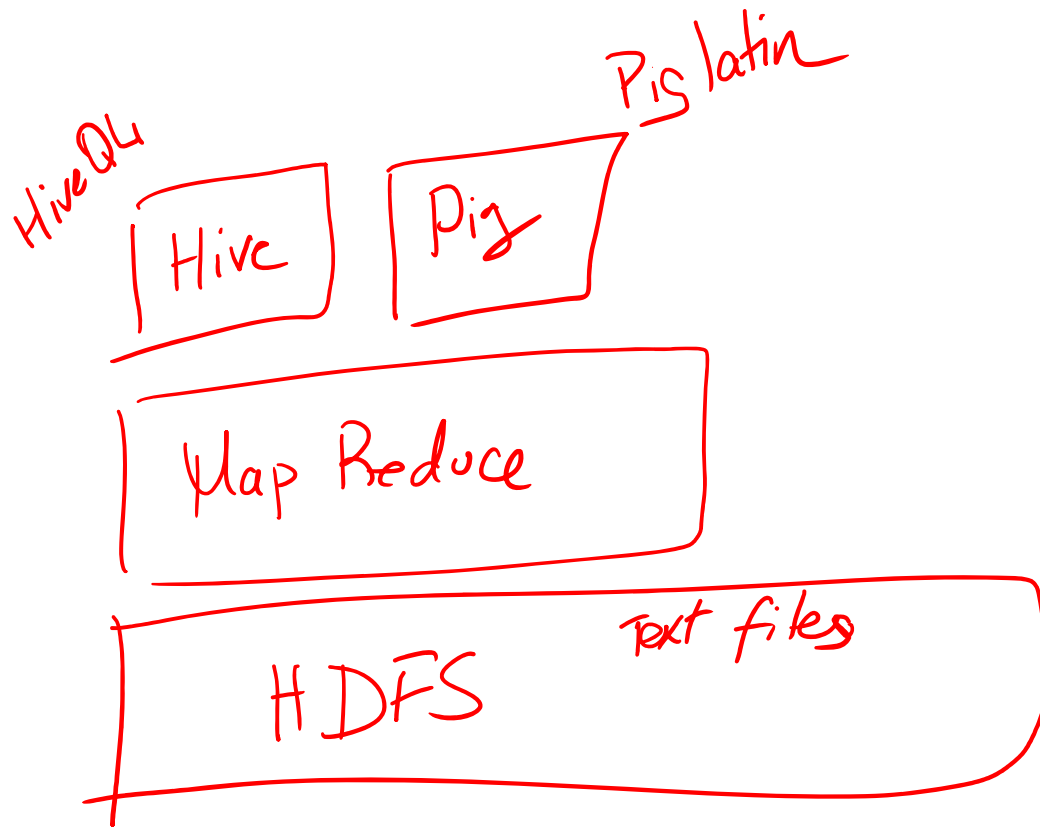
Alexander Rasin

College of CDM, DePaul University

October 12<sup>th</sup>, 2021

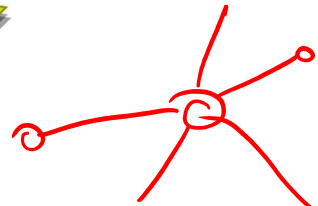
# Tonight

- Star schema / SSBM
- Hive and Pig
- Hadoop Streaming
- HBase





# Star-Schema



"Fact" table

**Transactions**

tid	cid	pid	date	amount
1	3	4	01-01-10	200.00
2	1	10	01-01-10	45.95
3	3	9	01-02-10	110.00

"Dimension" table

**Customers**

cid	name	address
1	...	
2	...	
3	Alex	

"Dimension" table

**Products**

pid	name	desc
1		
2		
3		
4	LCD Monitor	...

"Dimension" table

**Date**

date	day	holiday

# SSBM Schema

- 4 Dimension tables
  - Part: manufacturer, name, size, ...
  - Supplier: name, address, phone, ...
  - Customer: name, address, ...
  - Date: month, year, dayofweek, ...
- 1 Fact table: lineorder
  - Composite key (orderkey+linenumber)
  - 4 foreign keys (one per dimension)
  - Orderpriority, quantity, discount, tax, ...

# Join in Hive

*st Map1-a id Address*  
*BT Map1-b id -BT*

- Hive StarSchema optimization

SELECT Address, Count(\*)

*Map2 Address -*

FROM SmallTable st JOIN

BigTable bt ON (st.id = bt.id)

GROUP BY Address

- 2 MapReduce jobs
- Network transfer

hadoop hadoop-streaming -Mapper 1.py -reducer 2.py  
Map-Join -file 1.py -file 2.py  
-file SmallTable.txt

SELECT /\*+ MAPJOIN(st) \*/

Address, Count(\*)

Map1 Address -

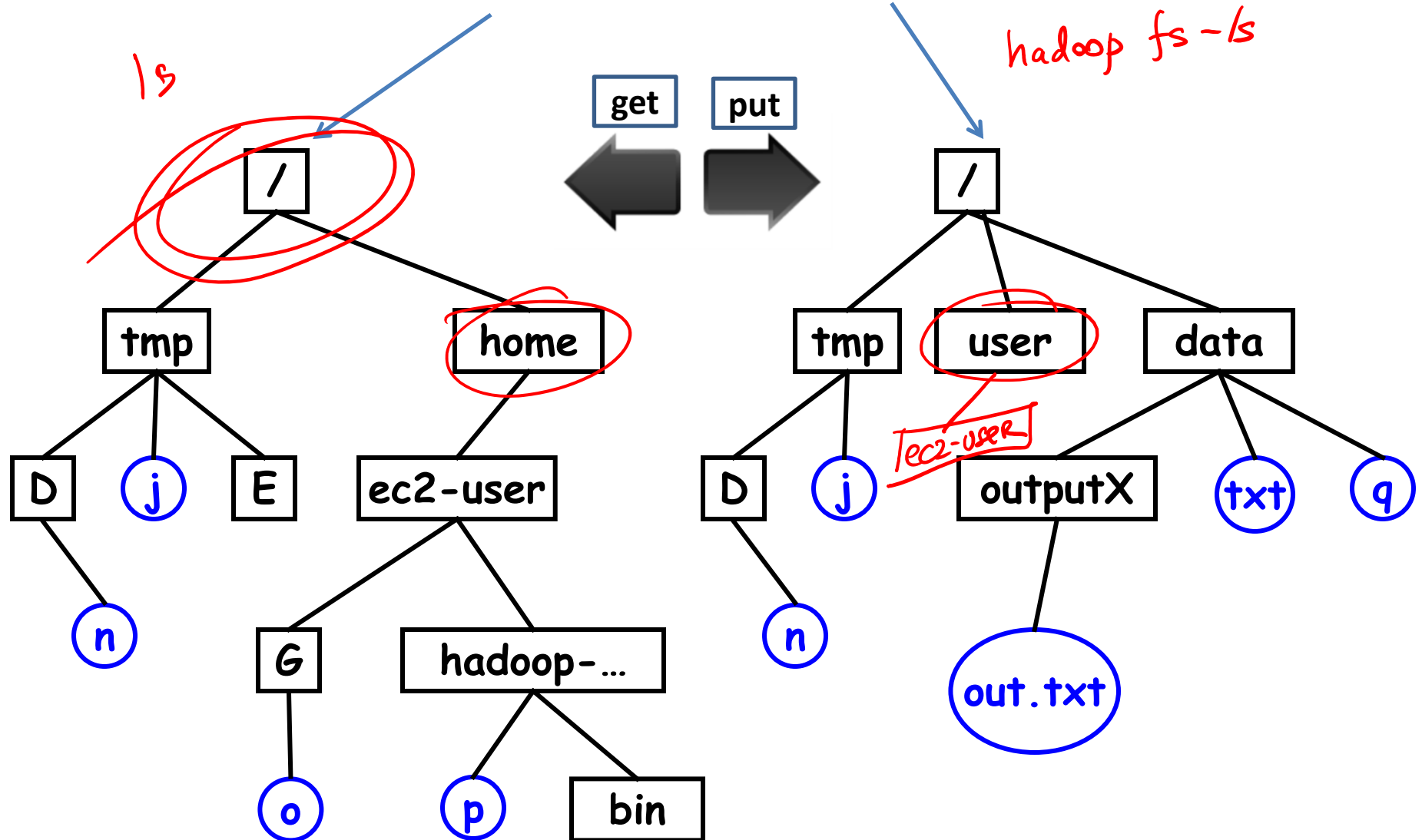
FROM SmallTable st JOIN

BigTable bt ON (st.id = bt.id)

GROUP BY Address

- Replicate the SmallTable
- set hive.auto.convert.join = true;
  - hive.mapjoin.smalltable.filesize (25MB default)

# LocalFS and HDFS





X Processor 

2cpu  
2G RAM

 2.4c/hour

✓ Disk 

8G x

 10c/G/Month

# Pig Example

- `cd hdfs:///`
- `ls`
- `mkdir testpig`
- `copyFromLocal /etc/passwd passwd`
- `passwd = LOAD 'hdfs:///testpig/passwd'  
USING PigStorage(':') AS (user:chararray,  
passwd:chararray, uid:int, gid:int,  
userinfo:chararray, home:chararray,  
shell:chararray);`

# Pig Use

- DUMP passwd;
- groupShell = GROUP passwd BY shell;
- DUMP groupShell;
- groupCount = FOREACH groupShell GENERATE group, COUNT(passwd);

# Hadoop Streaming

```
hadoop jar hadoop-streaming.jar
```

```
-input myInputDirs
```

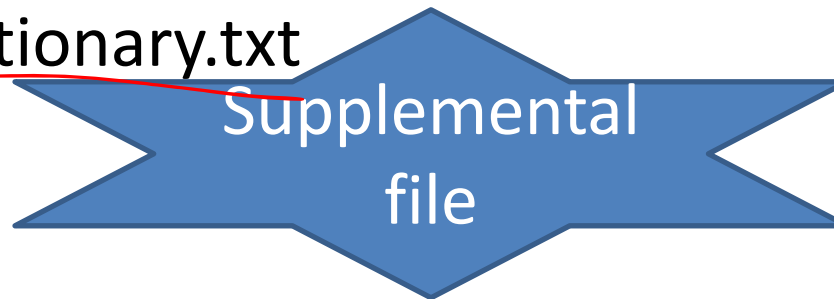
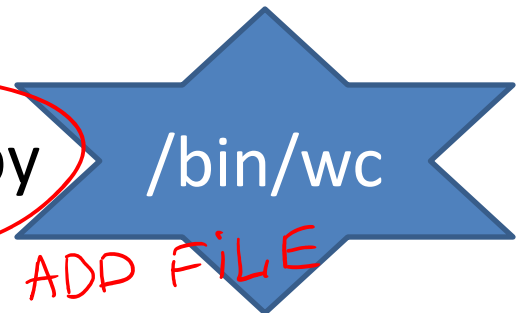
```
-output myOutputDir
```

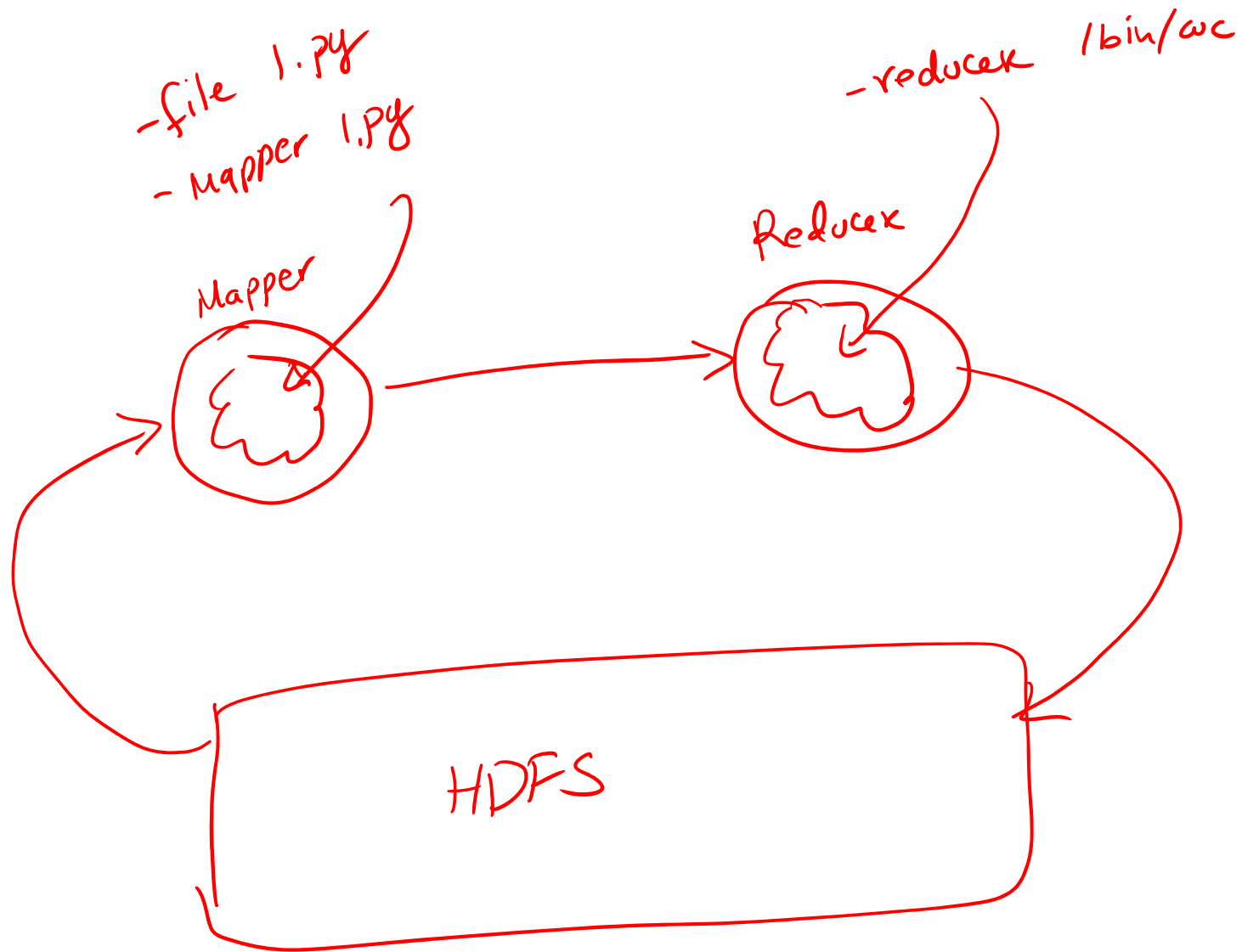
```
-mapper myPythonScript.py
```

```
-reducer myOtherPythonScript.py
```

```
-file myPythonScript.py
```

```
-file myDictionary.txt
```

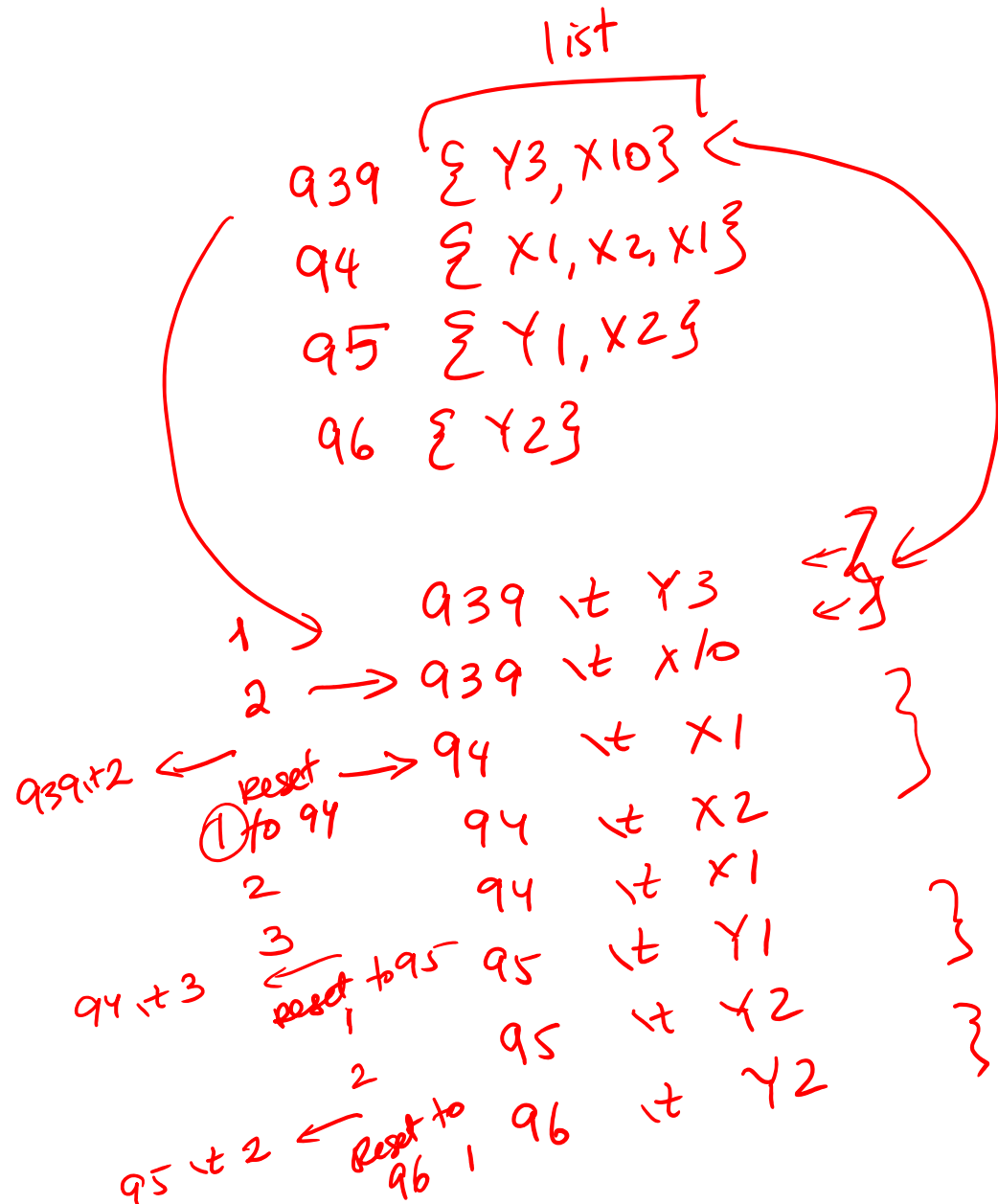




# Hadoop Streaming

- Map: cat binary
- Reduce: value count in python
- `hadoop jar hadoop-streaming-2.6.4.jar -input /home/ec2-user/mlens -output /data/output1 -mapper /bin/cat -reducer myReducer.py -file myReducer.py`

q4 it x1  
 q5 it x2  
 q39 it x10  
 q4 it x1  
 q4 it x2  
 q5 it x1  
 q5 it x2  
 q6 it x3  
 q39 it x3



# Hadoop Streaming

- Map: Timestamp=>Weekday in python
- Reduce: value count in python
- `hadoop jar hadoop-streaming-2.6.4.jar -input /user/ec2-user/mlens -output /data/output4 -mapper myMapper.py -reducer myReducer.py -file myReducer.py -file myMapper.py`





# Hadoop Streaming Options

- `-D stream.map.output.field.separator=|`
- `-mapper`  
`org.apache.hadoop.mapred.lib.IdentityMapper`
- `-partitioner`  
`org.apache.hadoop.mapred.lib.KeyFieldBasedPartitioner`
- `-D mapred.reduce.tasks=2`
- `-D mapred.text.key.comparator.options=-nr`

*# of reducers*

*numeric*

*reverse*

# Hadoop Streaming

- Map: cat binary
- Reduce: value count in python
- `hadoop jar hadoop-streaming-2.6.4.jar -D  
mapred.reduce.tasks=3 -D  
mapred.output.key.comparator.class=org.apache.hadoop  
.mapred.lib.KeyFieldBasedComparator -D  
mapred.text.key.comparator.options=-n -input  
/user/ec2-user/mlens -output /data/output3 -mapper  
/bin/cat -reducer myReducer.py -file myReducer.py`

# MapReduce Debugging

- `#!/usr/bin/python`
- Map code:
  - `cat test.txt | python myMapper.py`
- Shuffle code:
  - `cat test.txt | python myMapper.py | sort`
- Reduce code:
  - `cat test.txt | python myMapper.py | sort | python myReducer.py`

# Hadoop DistributedCache

- Lookup tables
- Data transferred before JVM starts
  - `-file myMapper.py` *Mapper myMapper.py*  
*copy* *run*
- Hive
  - `add FILE weekday_mapper.py;` *copy*

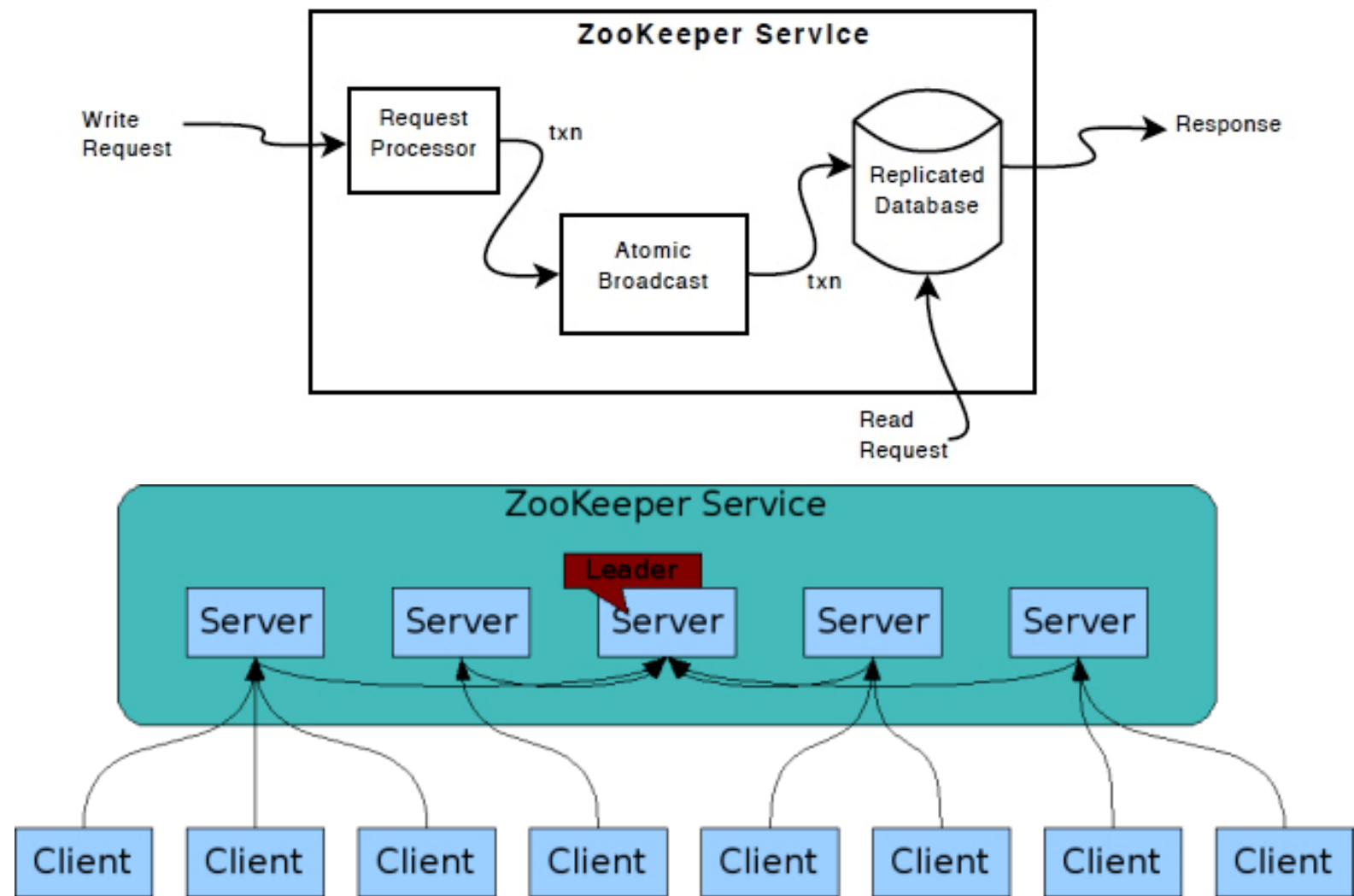


# Zookeeper

- “Because coordinating distributed systems is a zoo”
- Coordinate service for distributed applications
  - Naming
  - Configuration
  - Synchronization
  - Group management/leader election



# Zookeeper Overview



# Zookeeper Uses

- Barriers
  - Delay processing until a condition is met
- Queues
  - (priority) Distributed queue
- Locks
  - Exclusive access to data
- Leader election
  - Select one “leader” node
  - Find a replacement when needed



# Wide-Column Stores

Keys and Values

Optimize column family distribution

Add columns (no penalty)

	<u>row keys</u>	column family "color"	column family "shape"
row	"first"	"red": "#F00" "blue": "#00F" "yellow": "#FF0"	"square": "4"
row	"second"		"triangle": "3" "square": "4"

Table 1

row keys

- 1
- 2
- 3
- 4
- 5
- 6
- 7
- 8
- 9
- 10
- 11
- 12
- 14

cd-fam1				fam2	
a	b	c	d	x	y

1

7

a  
b

1 7

1 1

x

1

2

4

5

fam4				c4
c1	c2	c3		

1

5

a

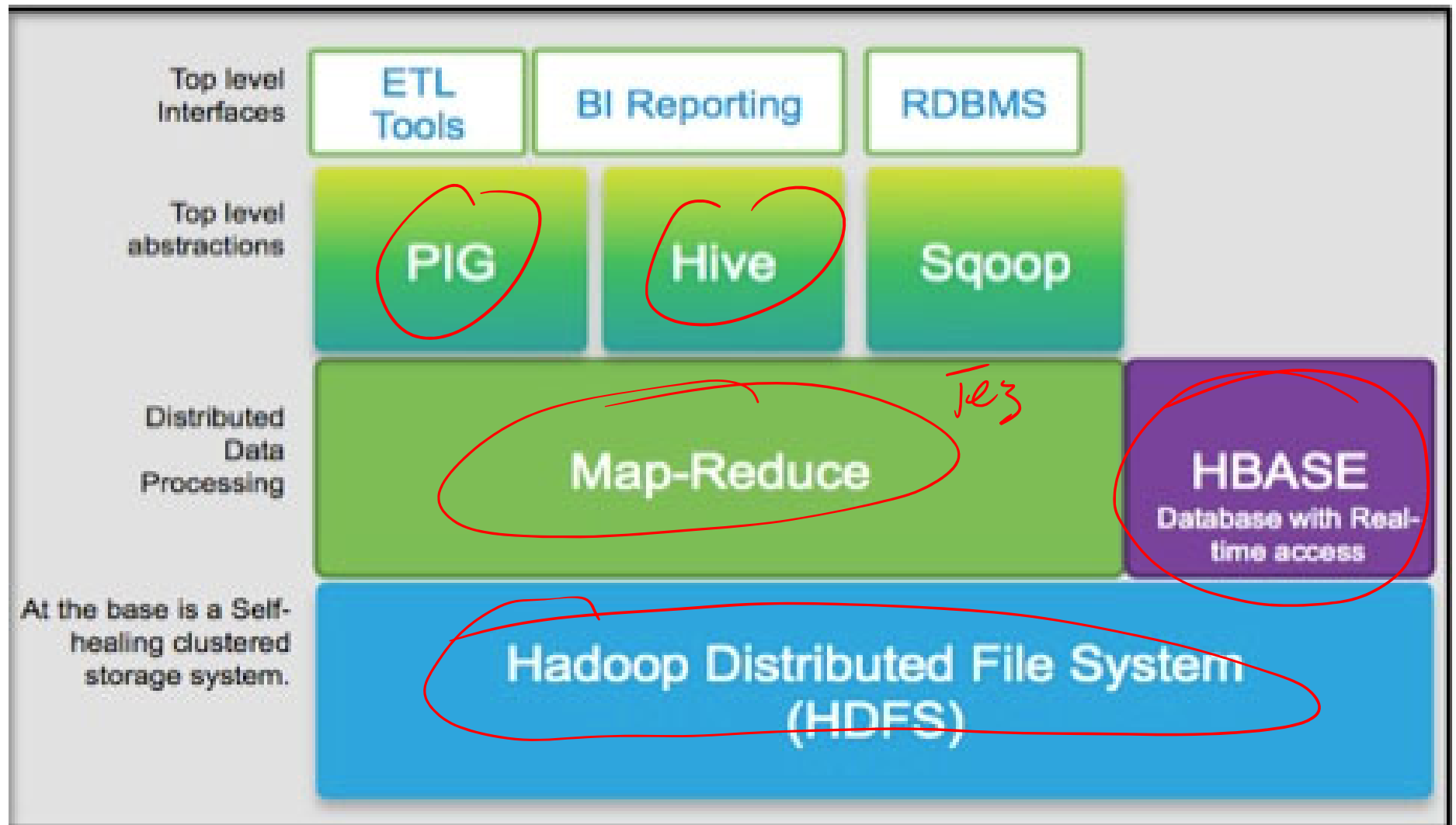
b

d

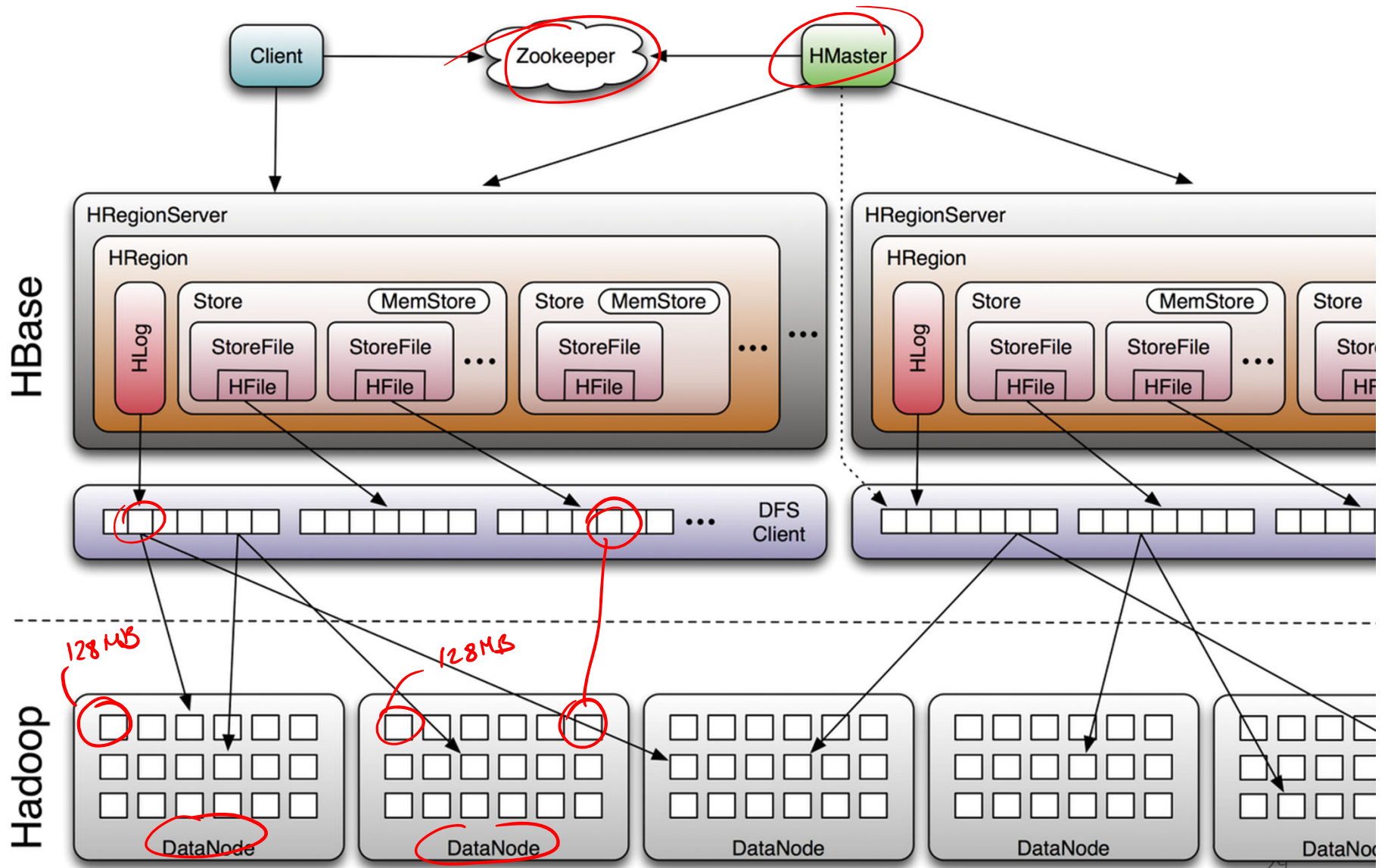
# HBase

- Column-oriented store
  - Google's Bigtable
- Random, real-time updates
  - Individual record access
- Highly available/distributed
- No SQL
  - No joins
  - (some) indexing

# HBase



# HBase Architecture



# HBase Example

- Create/populate a table:
  - create 'test', 'col\_fam'
  - list 'test'
  - put 'test', 'row1', 'col\_fam:a', 'value1'
  - put 'test', 'row2', 'col\_fam:b', 'value2'
  - put 'test', 'row3', 'col\_fam:c', 'value3'
  - put 'test', 'row5', 'col\_fam:a', 'value4'
  - put 'test', 'row3', 'col\_fam:b', 'value5'
  - scan 'test'

Test

	<u>col - fam</u>		
	<u>a</u>	<u>b</u>	<u>c</u>
row 1	Value1		
row 2		Value2	
row 3		<u>Value3</u>	Value3
row 5	Value4		

new Values -

# Using HBase

```
put 'test', 'row1', 'col_fam:z', 'value6'  
put 'test', 'row5', 'col_fam:z', 'value7'  
put 'test', 'row6', 'col_fam:z', 'value8'  
scan 'test', {COLUMNS => 'col_fam:c'}  
scan 'test', {STARTROW => 'row3'}  
put 'test', 'row3', 'col_fam:b', 'newValue5'  
  
scan 'test', {COLUMNS => 'col_fam:b'}  
scan 'test', {TIMERANGE=>[a, b]}
```



# We Stopped Here

# Hive SELECT TRANSFORM

```
CREATE TABLE u_data_new ( userid INT,  
movieid INT, rating INT, weekday String)  
ROW FORMAT DELIMITED FIELDS  
TERMINATED BY '\t';  
add FILE weekday_mapper.py;  
INSERT OVERWRITE TABLE u_data_new  
SELECT TRANSFORM (userid, movieid, rating, unixtime)  
USING 'python weekday_mapper.py'  
AS (userid, movieid, rating, weekday) FROM u_data;
```

# SQL to MapReduce

```
SELECT AVG(Grade) as GPA, FirstName,  
       COUNT(DISTINCT LastName), MAX(s.Year)  
FROM Student s, Grades g  
WHERE s.ID = g.SID  
GROUP BY FirstName  
ORDER BY GPA
```

# Next Time:

- Performance and Compression
- NoSQL
- DNS and Firewalls
- Multi-node cluster setup