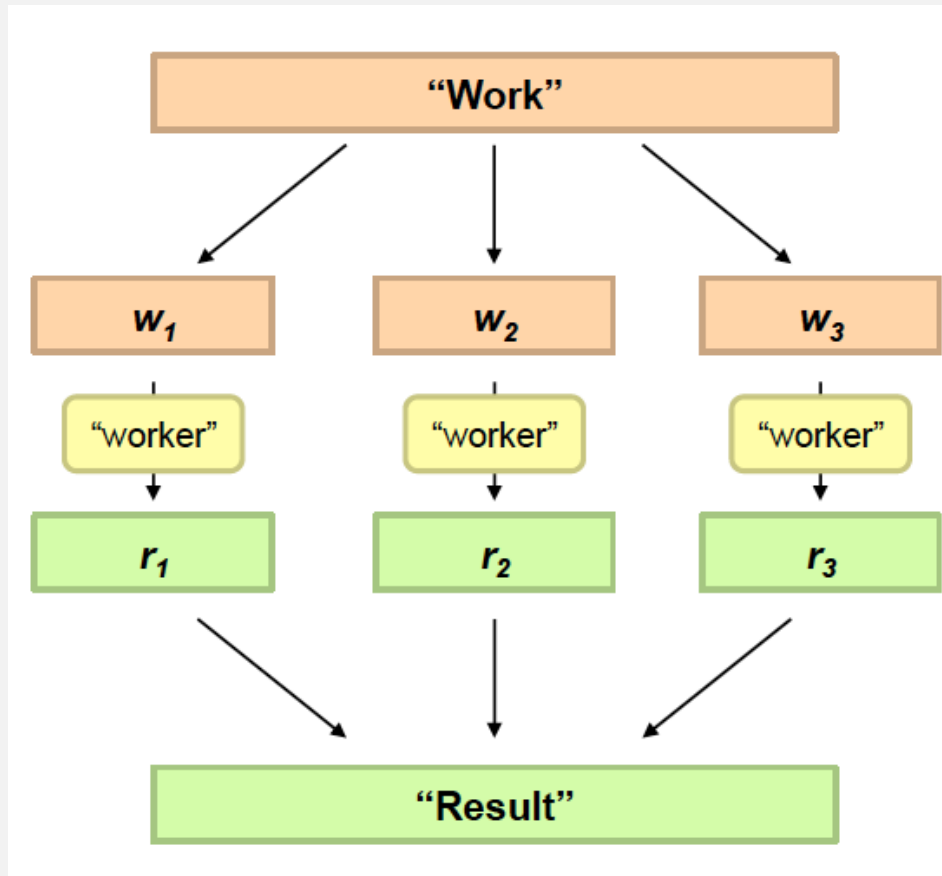# DDAM
## INTRODUCTION TO HADOOP

Docente: Patrizio Dazzi

Mail: patrizio.dazzi@isti.cnr.it

- Philosophy to Scale for Big Data?

# WORKLOAD DECOMPOSITION



**Divide Work**

**Combine Results**

# DISTRIBUTED PROCESSING IS NON-TRIVIAL

- How to assign tasks to different workers in an efficient way?

- What happens if tasks fail?

- How do workers exchange results?

- How to synchronize distributed tasks allocated to different workers?



Image courtesy of Master isolated images at FreeDigitalPhotos.net

# (PERFORMANT) BIG DATA STORAGE IS CHALLENGING

- Data Volumes are massive

- Reliability of Storing PBs of data is challenging

- All kinds of failures: Disk/Hardware/Network Failures

- Probability of failures simply increase with the number of machines …



- Performance, performance, performance

# ONE POPULAR SOLUTION: HADOOP*

Hadoop Cluster at Yahoo! (Credit: Yahoo)



* = but is not the only one

# HADOOP* OFFERS

- Redundant, Fault-tolerant data storage

- Parallel computation framework

- Job coordination



* = as well as analogous tools

# HADOOP* OFFERS

**Programmers**

*No longer need to worry about*

- Q: Where file is located?
- Q: How to handle failures & data lost?
- Q: How to divide computation?

# **There** ain't **no** such thing as a **free lunch**

> ## ~~HADOOP IS THE SOLUTION~~

- HADOOP is NOT magic

- Heuristics work… often, not always

- High Performances are not for free

- We wil see in the future how to deal with and implement optimizations

# A REAL WORLD EXAMPLE OF NEW YORK TIMES

- **Goal:** Make entire archive of articles available online: 11 million, from 1851

- **Task:** Translate 4 TB TIFF images to PDF files

- **Solution:** Used Amazon Elastic Compute Cloud (EC2) and Simple Storage System (S3)

- **Time:** ?

- **Costs:** ?

# A REAL WORLD EXAMPLE OF NEW YORK TIMES

- **Goal:** Make entire archive of articles available online: 11 million, from 1851
- **Task:** Translate 4 TB TIFF images to PDF files
- **Solution:** Used Amazon Elastic Compute Cloud (EC2) and Simple Storage System (S3)
- **Time:** **< 24 hours**
- **Costs:** **$240**

# A LITTLE HISTORY ON HADOOP

- Hadoop is an open-source implementation based on **Google File System** (GFS) and **MapReduce** from Google

- Hadoop was created by **Doug Cutting** and **Mike Cafarella** in 2005

- Hadoop was donated to **Apache** in 2006

# WHO ARE USING HADOOP?

### Social
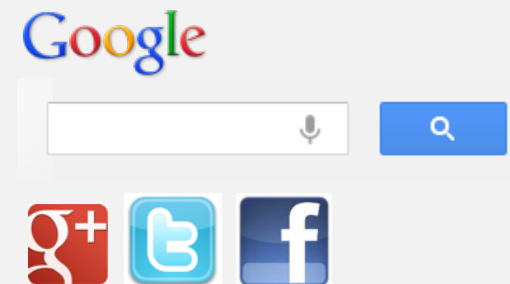


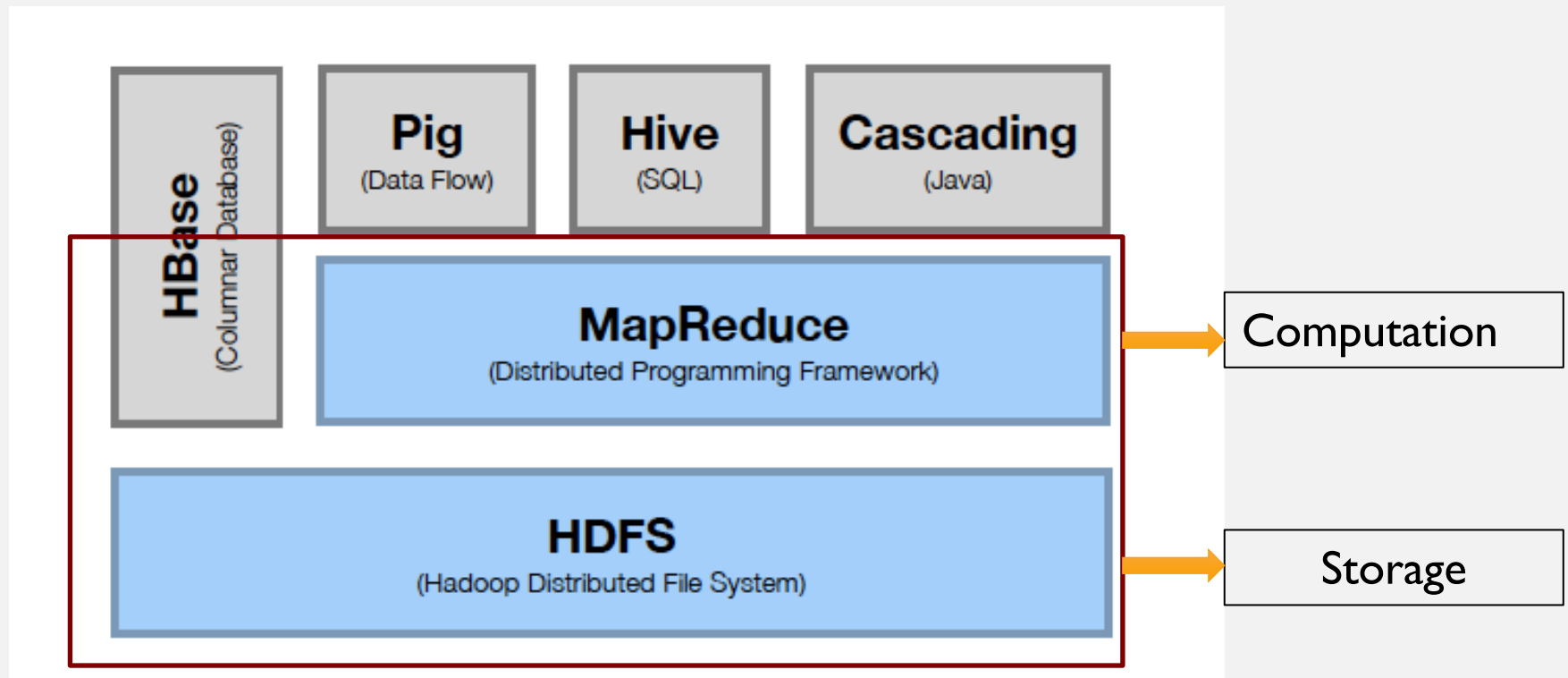### User Tracking & Engagement



### Homeland Security



### eCommerce



### Financial Services



### Real Time Search

# HADOOP STACK

| HBase (Columnar Database) | Pig (Data Flow) | Hive (SQL) | Cascading (Java) |

**MapReduce** (Distributed Programming Framework) → Computation

**HDFS** (Hadoop Distributed File System) → Storage

# HADOOP RESOURCES

- Hadoop at ND:

http://ccl.cse.nd.edu/operations/hadoop/

- Apache Hadoop Documentation:

http://hadoop.apache.org/docs/current/

- Data Intensive Text Processing with Map-Reduce

http://lintool.github.io/MapReduceAlgorithms/

- Hadoop Definitive Guide:

http://www.amazon.com/Hadoop-Definitive-Guide-Tom-White/dp/1449311520

# HDFS

HADOOP DISTRIBUTED FILE SYSTEM

# MOTIVATION QUESTIONS

- **Problem 1:** Data is too big to store on one machine.

- **HDFS:** Store the data on multiple machines!

# MOTIVATION QUESTIONS

- **Problem 2:** Very high end machines are too expensive

- **HDFS:** Run on commodity hardware!

# MOTIVATION QUESTIONS

- **Problem 3:** Commodity hardware will fail!

- **HDFS:** Software is intelligent enough to handle hardware failure!

# MOTIVATION QUESTIONS

- **Problem 4:** What happens to the data if the machine stores the data fails?

- **HDFS:** Replicate the data!

- **Problem 5:** How can distributed machines organize the data in a coordinated way?

- **HDFS:** Master-Slave Architecture!

**Master**



**Name Node (NN)**

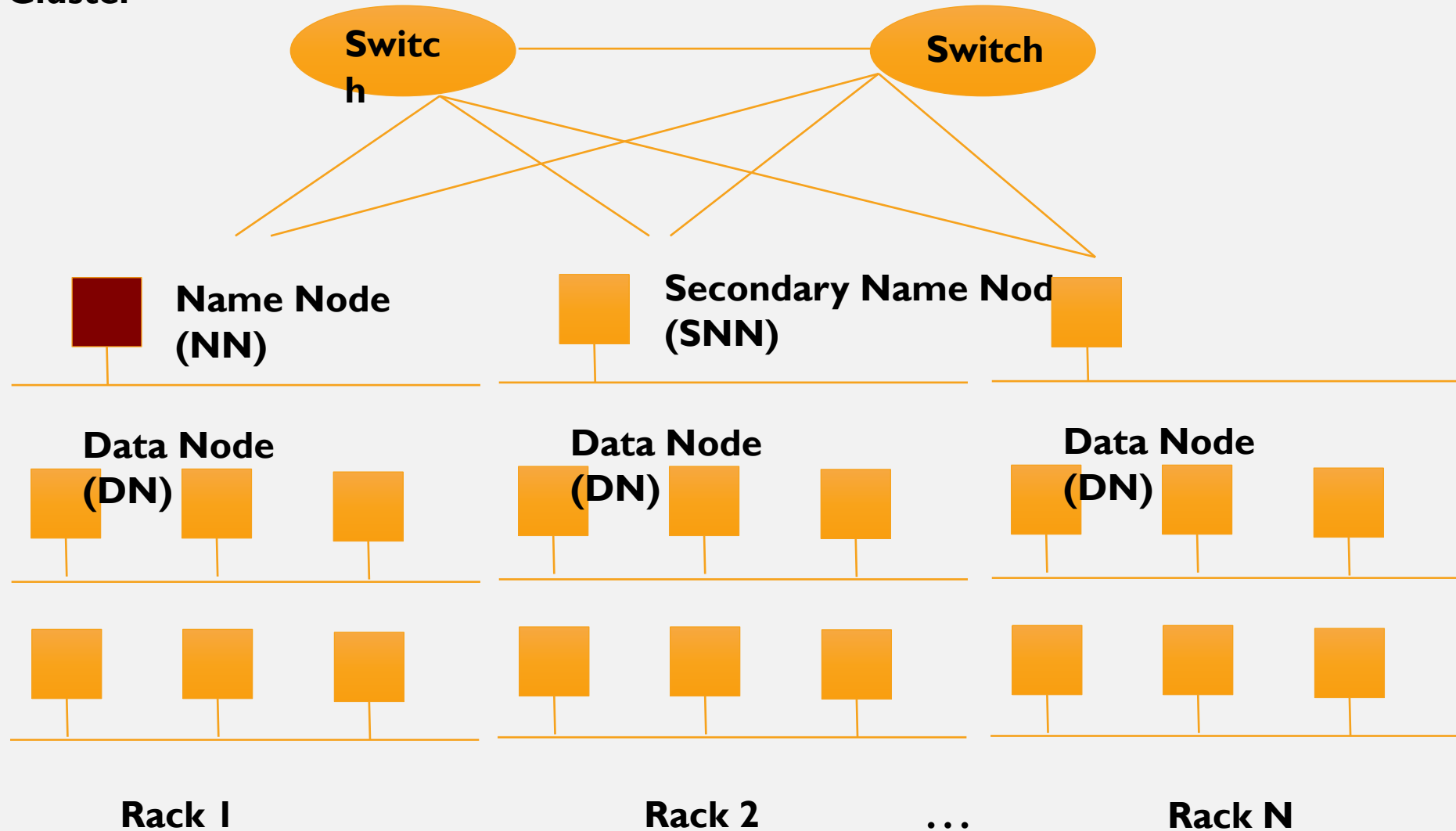**Secondary Name Node (SNN)**

**Data Node (DN)**

**Slaves**    **Single Rack Cluster**

- Name Node: Controller
  - File System Name Space Management
  - Block Mappings

- Data Node: Work Horses
  - Block Operations
  - Replication
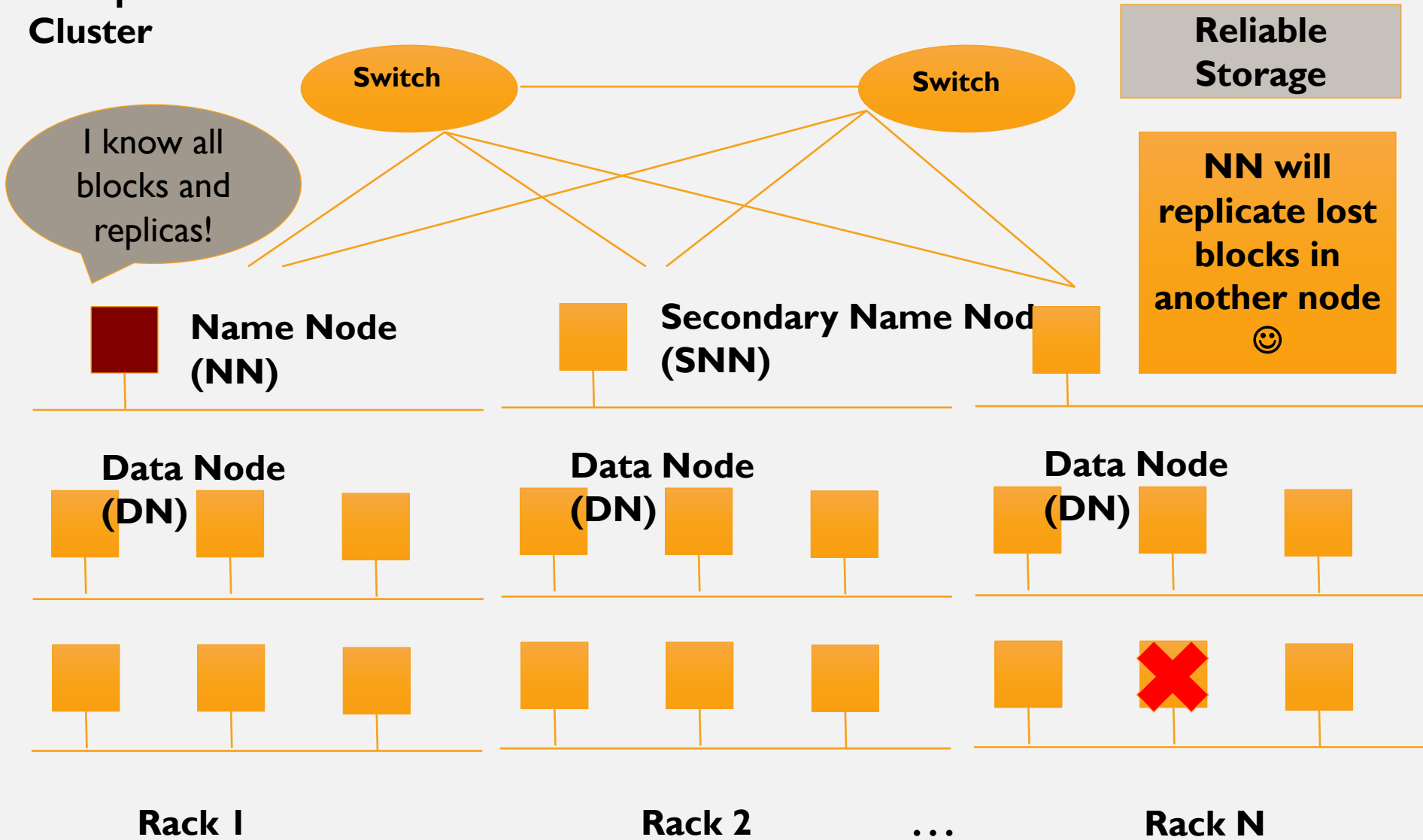
- Secondary Name Node:
  - Checkpoint node

**Multiple-Rack Cluster**

**How about network performance?**
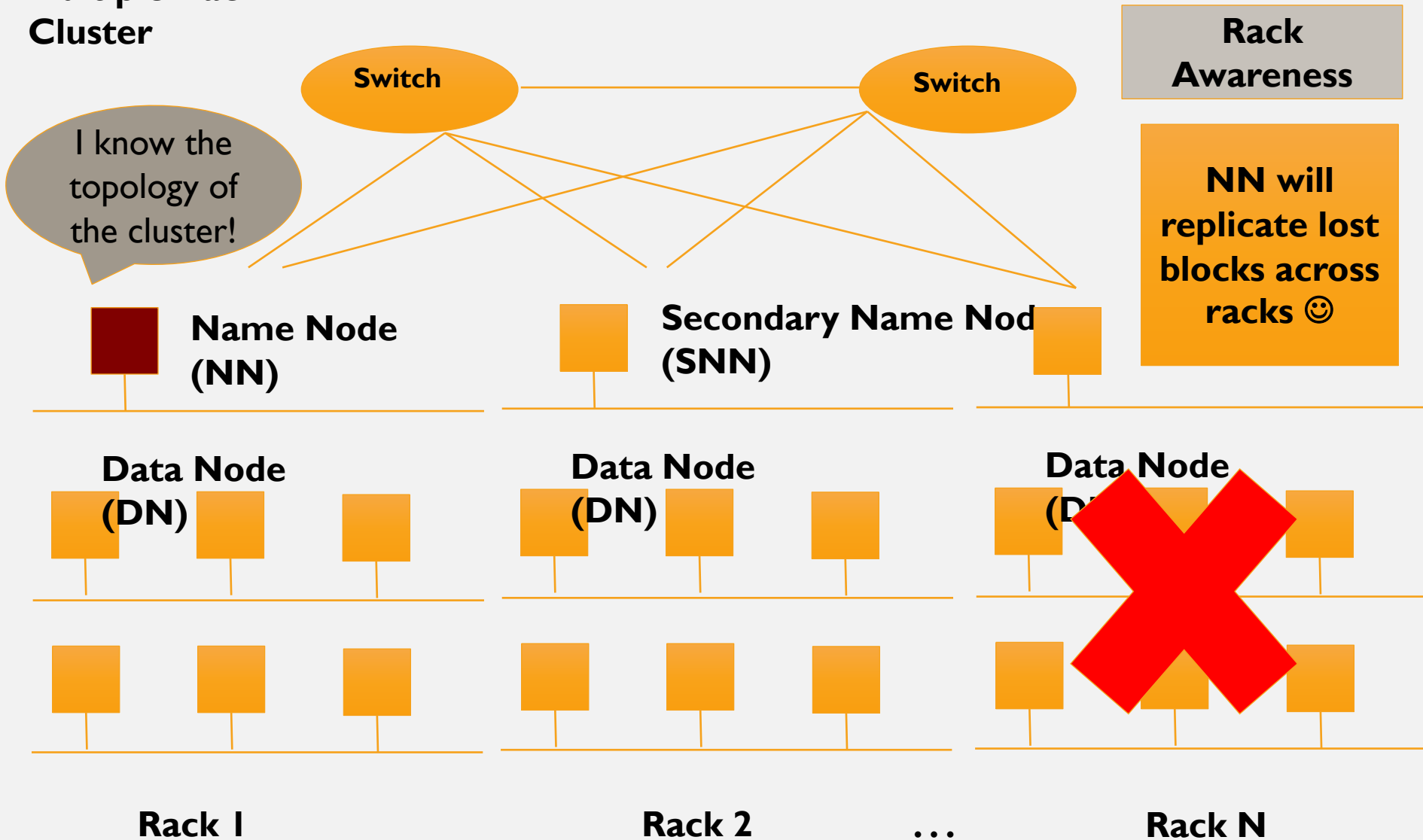
Switch

Switch

**Keep bulky communication within a rack!**

**Name Node (NN)**

**Secondary Name Node (SNN)**

**Data Node (DN)**

**Data Node (DN)**

**Data Node (DN)**

**Rack 1**

**Rack 2**

. . .

**Rack N**

# HDFS INSIDE: NAME NODE

**Name Node**

Snapshot of FS

Edit log: record changes to FS

| Filename | Replication factor | Block ID |
|----------|-------------------|----------|
| File 1 | 3 | [1, 2, 3] |
| File 2 | 2 | [4, 5, 6] |
| File 3 | 1 | [7,8] |

**Data Nodes**

1, 2, 5, 7, 4, 3
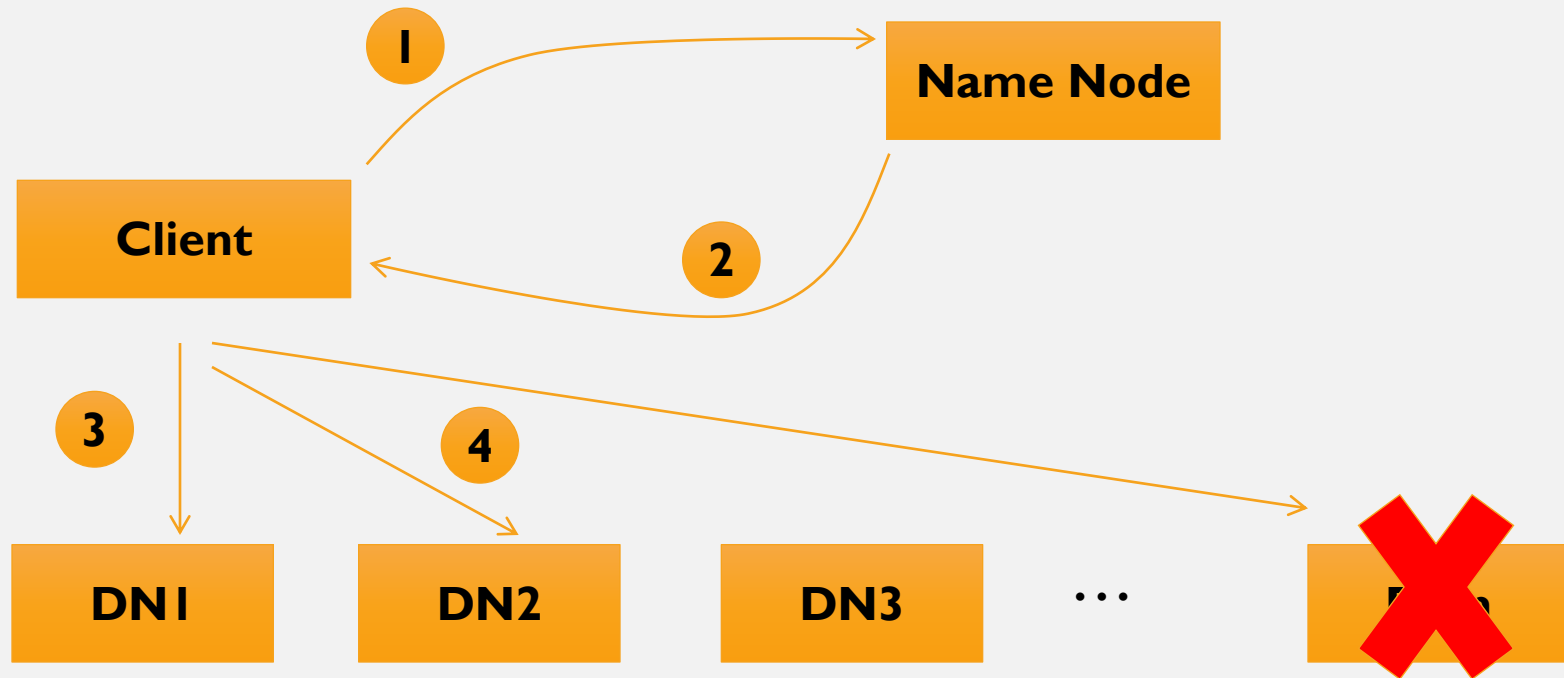
1, 5, 3, 2, 8, 6

1, 4, 3, 2, 6

# HDFS INSIDE: BLOCKS

- Q: Why do we need the abstraction "Blocks" in addition to "Files"?

- Reasons:
  - File can be larger than a single disk
  - Block is of fixed size, easy to manage and manipulate
  - Easy to replicate and do more fine grained load balancing

# HDFS INSIDE: BLOCKS

- HDFS Block size is by default **64 MB**, why it is much larger than regular file system block?

- Reasons:

  - Minimize overhead: disk seek time is almost constant

# HDFS INSIDE: READ



1. Client connects to NN to read data
2. NN tells client where to find the data blocks
3. Client reads blocks directly from data nodes (without going through NN)
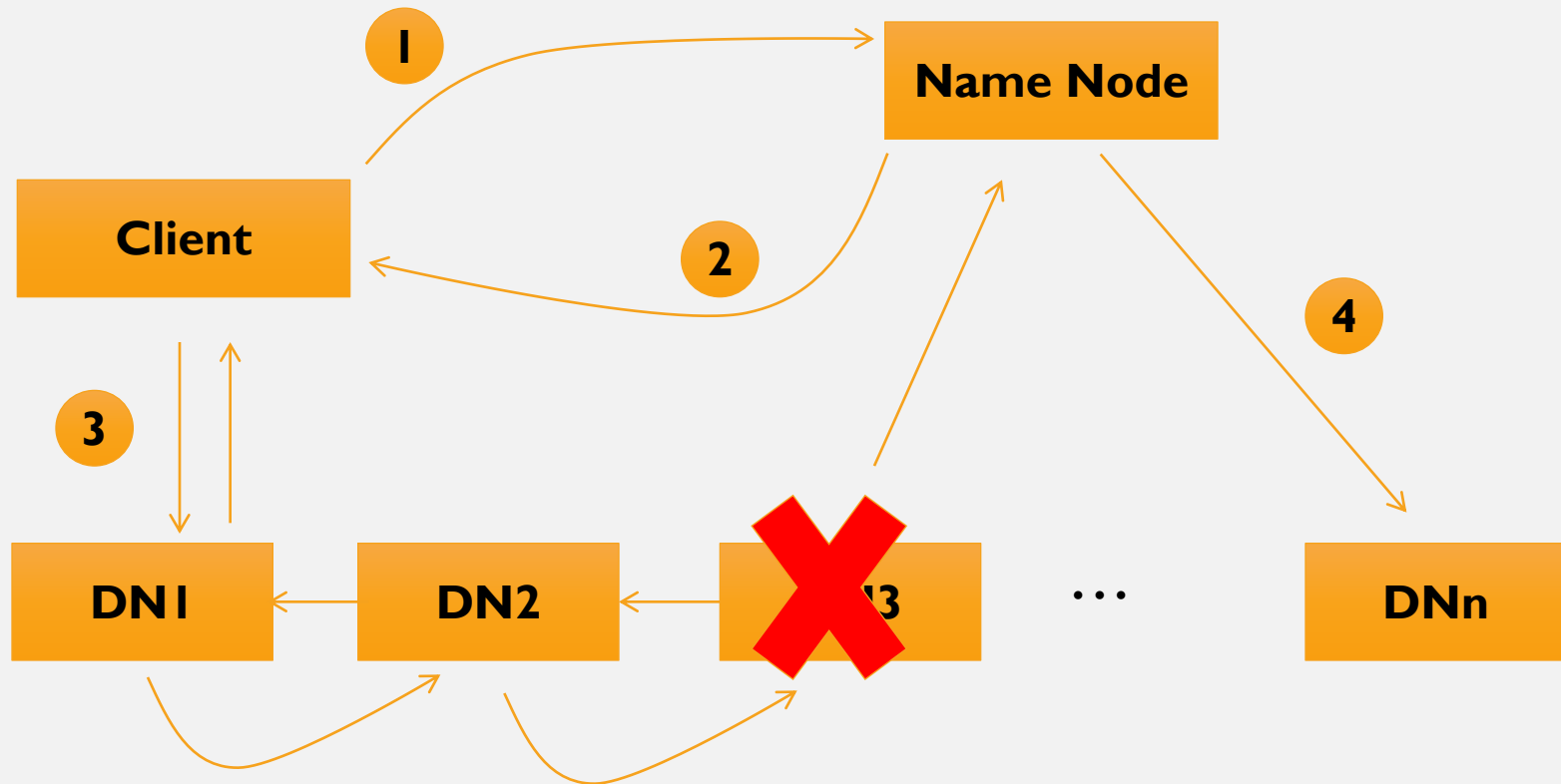4. In case of node failures, client connects to another node that serves the missing block

# HDFS INSIDE: READ

- Q: Why does HDFS choose such a design for read? Why not ask client to read blocks through NN?

- Reasons:
  - Prevent NN from being the bottleneck of the cluster
  - Allow HDFS to scale to large number of concurrent clients
  - Spread the data traffic across the cluster

# HDFS INSIDE: READ

- Q: Given multiple replicas of the same block, how does NN decide which replica the client should read?

- HDFS Solution:
  - Rack awareness based on network topology

# HDFS INSIDE: WRITE



1. Client connects to NN to write data
2. NN tells client write these data nodes
3. Client writes blocks directly to data nodes with desired replication factor
4. In case of node failures, NN will figure it out and replicate the missing blocks

# HDFS INSIDE: WRITE

- Q: Where should HDFS put the three replicas of a block? What tradeoffs we need to consider?

- Tradeoffs:

  - Reliability

  - Write Bandwidth

  - Read Bandwidth

Q: What are some possible strategies?

- Replication Strategy vs Tradeoffs

| | Reliability | Write Bandwidth | Read Bandwidth |
|---|---|---|---|
| Put all replicas on one node | 🙁 | 🙂 | 🙁 |
| Put all replicas on different racks | 🙂 | 🙁 | 🙁 |
| | | | |

# HDFS INSIDE: WRITE

- Replication Strategy vs Tradeoffs

| | Reliability | Write Bandwidth | Read Bandwidth |
|---|---|---|---|
| Put all replicas on one node | 🙁 | 🙂 | 🙁 |
| Put all replicas on different racks | 🙂 | 🙁 | 🙁 |
| HDFS: 1-> same node as client 2-> a node on different rack 3-> a different node on the same rack as 2 | 🙂 | OK | OK |

# MAPREDUCE



How MapReduce Works?

Map()      Shuffle      Reduce()     http://blog.sqlauthority.com

**Map**
extract something you care about from each record

**Reduce**
aggregate, summarize, filter, or transform

**Input Files**

Apple Orange Mango
Orange Grapes Plum

Apple Plum Mango
Apple Apple Plum

**Each line passed to inividual mapper instances**

Apple Orange Mango

Orange Grapes Plum

Apple Plum Mango

Apple Apple Plum

**Map Key Value Splitting**

Apple,1
Orange,1
Mango,1

Orange,1
Grapes,1
Plum,1

Apple,1
Plum,1
Mango,1

Apple,1
Apple,1
Plum,1

**Sort and Shuffle**

Apple,1
Apple,1
Apple,1
Apple,1

Grapes,1

Mango,1
Mango,1

Orange,1
Orange,1

Plum,1
Plum,1
Plum,1

**Reduce Key Value Pairs**

Apple,4

Grapes,1

Mango,2

Orange,2

Plum,3

**Final Output**

Apple,4
Grapes,1
Mango,2
Orange,2
Plum,3

# MAPPER

- Reads in input pair <Key,Value>

- Outputs a pair <K', V'>

  - Let's count number of each word in user queries (or Tweets/Blogs)

  - The input to the mapper will be <queryID, QueryText>:

    ```
    <Q1, "The teacher went to the store. The store was closed; the store
    opens in the morning. The store opens at 9am." >
    ```

  - The output would be:

    ```
    <The, 1> <teacher, 1> <went, 1> <to, 1>
    <the, 1> <store,1> <the, 1> <store, 1>
    <was, 1> <closed, 1> <the, 1> <store,1>
    <opens, 1> <in, 1> <the, 1> <morning,
    1> <the 1> <store, 1> <opens, 1> <at,
    1> <9am, 1>
    ```

# REDUCER

- Accepts the Mapper output, and aggregates values on the key

    - For our example, the reducer input would be:

        <store, 1> <store, 1> <store, 1><store, 1>

    - The output would be:

        **<store, 4>**

# JAVA MAP-REDUCE

```java
36 public static class MapForWordCount extends Mapper<LongWritable, Text, Text, IntWritable>{
37 public void map(LongWritable key, Text value, Context con) throws IOException, InterruptedException
38 {
39 String line = value.toString();
40 String[] words=line.split(",");
41 for(String word: words )
42 {
43     Text outputKey = new Text(word.toUpperCase().trim());
44   IntWritable outputValue = new IntWritable(1);
45   con.write(outputKey, outputValue);
46 }
47 }
48 }
49
50 public static class ReduceForWordCount extends Reducer<Text, IntWritable, Text, IntWritable>
51 {
52 public void reduce(Text word, Iterable<IntWritable> values, Context con) throws IOException, InterruptedE
53 {
54 int sum = 0;
55     for(IntWritable value : values)
56     {
57     sum += value.get();
58     }
59     con.write(word, new IntWritable(sum));
60 }
61
62 }
```

# PYTHON MAP-REDUCE

mapper.py

```python
#!/usr/bin/env python

import sys

# input comes from STDIN (standard input)
for line in sys.stdin:
    # remove leading and trailing whitespace
    line = line.strip()
    # split the line into words
    words = line.split()
    # increase counters
    for word in words:
        # write the results to STDOUT (standard output);
        # what we output here will be the input for the
        # Reduce step, i.e. the input for reducer.py
        #
        # tab-delimited; the trivial word count is 1
        print '%s\t%s' % (word, 1)
```

reducer.py

```python
#!/usr/bin/env python

from operator import itemgetter
import sys

current_word = None
current_count = 0
word = None

# input comes from STDIN
for line in sys.stdin:
    # remove leading and trailing whitespace
    line = line.strip()

    # parse the input we got from mapper.py
    word, count = line.split('\t', 1)

    # convert count (currently a string) to int
    try:
        count = int(count)
    except ValueError:
        # count was not a number, so silently
        # ignore/discard this line
        continue

    # this IF-switch only works because Hadoop sorts map output
    # by key (here: word) before it is passed to the reducer
    if current_word == word:
        current_count += count
    else:
        if current_word:
            # write result to STDOUT
            print '%s\t%s' % (current_word, current_count)
        current_count = count
        current_word = word

# do not forget to output the last word if needed!
if current_word == word:
    print '%s\t%s' % (current_word, current_count)
```

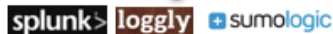# The Big Data Landscape

dave@vcdave.com

bigdatalandcsape.com

# Big Data Landscape (Version 2.0)

© Matt Turck (@mattturck) and ShivonZilis (@shivonz) Bloomberg Ventures

# BIG DATA LANDSCAPE, VERSION 3.0

Exited: Acquisition or IPO

## Infrastructure

- NoSQL Databases
- Hadoop On Prem
- Big Data
- NewSQL Databases
- Cluster Service
- MPP
- Management / Monitoring
- Graph Databases
- Data Transformation
- Security
- Storage
- Crowd-sourcing
- App Dev

## Analytics

- Analytics Platforms
- For Business Analysts
- Data Science Platforms / Tools
- BI Platforms
- Unstructured Data
- Data Visualization
- Machine Learning
- AI
- Social Analytics
- Location/People/Events
- Analytic Services
- Big Data Search
- Statistical Computing
- Log Analysis
- Crowd-Sourced
- Real-Time
- SMB

## Applications

- Ad Optimization
- Publisher Tools
- Marketing
- Finance
- Human
- Legal
- Government / Regulation
- Security
- Education n/Learning
- Health
- Industries

## Cross Infrastructure /

SAP · SAS · IBM · Google · Microsoft · vmware · amazon · 1010data · talend · Teradata

## Open Source

- Framework
- Query / Data Flow
- Data Access
- Coordination / Workflow
- Real-Time
- Stat Tools
- Machine
- Cloud Deploy
- Search

## Data Sources

- Data Mkts
- Data Source
- Sensor Data
- Incubators & School

© Matt Turck (@mattturck), Sutian Dong (@sutiandong) & FirstMark Capital (@firstmarkcap)