

# Hálózatba kapcsolt erőforrás platformok és alkalmazásaik

Simon Csaba

TMIT

2018

# Motivation

- Process lots of data
  - Google processed about **24 petabytes** of data per day in 2009.
- **A single machine** cannot serve all the data
  - You need a distributed system to store and process **in parallel**
- Parallel programming?
  - **Threading** is hard!
  - How do you facilitate **communication** between nodes?
  - How do you **scale to more machines**?
  - How do you handle machine **failures**?

BIGDATA?

# What is BigData?

“Big Data” exceeds the capacity of traditional analytics and information management paradigms across what is known as the 4 V’s: Volume, Variety, Velocity, and Veracity

## Veracity



### Uncertainty of Data

With exponential increases of data from unfiltered and constantly flowing data sources, data quality often suffers and new methods must find ways to “sift” through junk to find meaning

PwC

## Velocity



### Analysis of Streaming Data

The speed at which data is generated and used. New data is being created every second and in some cases it may need to be analyzed just as quickly

## Variety



### Different Forms of Data

Represents the diversity of the data. Data sets will vary by type (e.g. social networking, media, text) and they will vary how well they are structured

## Volume



### Scale of Data

Reflects the size of a data set. New information is generated daily and in some cases hourly, creating data sets that are measured in terabytes and petabytes

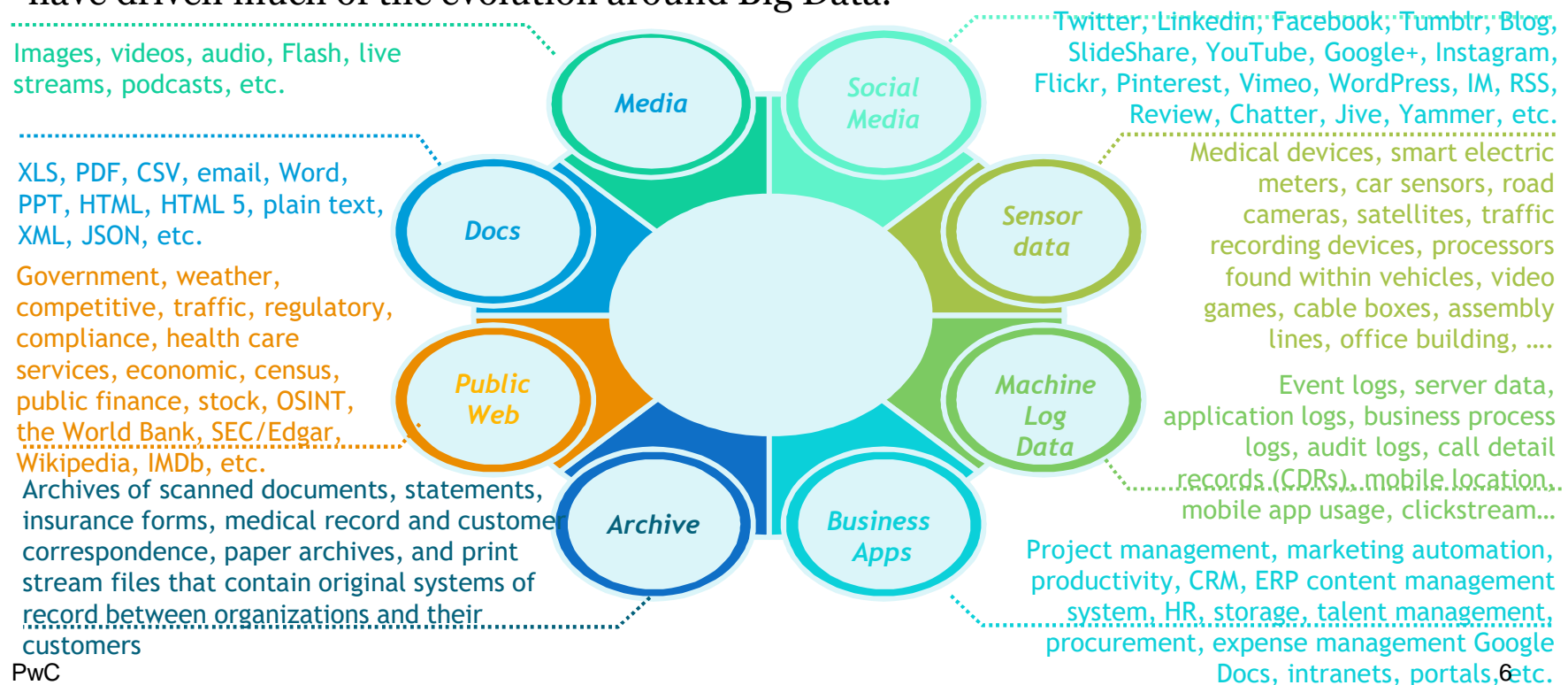
# The Promise of Big Data

Even more important than its definition is what Big Data promises to achieve:  
*intelligence in the moment.*

	<i>Traditional Techniques &amp; Issues</i>	<i>Big Data Differentiators</i>
<b>Vera city</b>	"Does not account for biases, noise and abnormality in data"	"Data is stored, and mined meaningful to the problem being analyzed" "Keeps data clean and processes to keep 'dirty data' from accumulating in your systems"
<b>Velo city</b>	"No real time analysis"	<b>In real-time:</b> "Dynamically analyze data" "Consistently integrate new information"
<b>Var iety</b>	"Compatibility issues" "Advanced analytics struggle with non-numerical data"	"Auto deletes unwanted to ensure optimal storage" "Frameworks accommodate varying data types and data models" "Insightful analysis with very few parameters"
<b>Vol.</b>	"Analysis is limited to small data sets" "Analyzing large data sets = High Costs & High Memory"	"Scalable for huge amounts of multi-sourced data" "Facilitation of massively parallel processing" "Low-cost data storage"

# Types of Big Data

Variety is the most unique aspect of Big Data. New technologies and new types of data have driven much of the evolution around Big Data.



“Single sources of data are no longer sufficient to cope with the increasingly complicated problems in many policy arenas.”<sup>1</sup>

Big data “is not notable because of its size, but because of its relationality to other data. Due to efforts to mine and aggregate data, Big Data is fundamentally networked.”<sup>2</sup>

(1) M. Milakovich, “Anticipatory Government: Integrating big data for Smaller Government”, in Oxford Internet Institute “Internet, Politics, Policy 2012” Conference, Oxford, 2012

(2) D. Boyd and K. Crawford, “Six Provocations for big data,” in A Decade in Internet Time: Symposium on the Dynamics of the Internet and Society, 2011

## Why is BigData valuable?

We have identified 5 key areas where Big Data is uniquely valuable:

### Accessibility to Data

Enhanced visibility of relevant information and better transparency to massive amounts of data. **Improved reporting** to stakeholders.

### Decision Making

Next generation analytics can enable **automated decision making** (inventory management, financial risk assessment, sensor data management, machinery tuning).

### Marketing Trends

Segmentation of population to **customize** offerings and **marketing** campaigns (consumer goods, retail, social, clinical data, etc).

### Performance Improvement

Exploration for, and discovery of, new needs, can drive organizations to fine tune for **optimal performance** and efficiency (employee data).

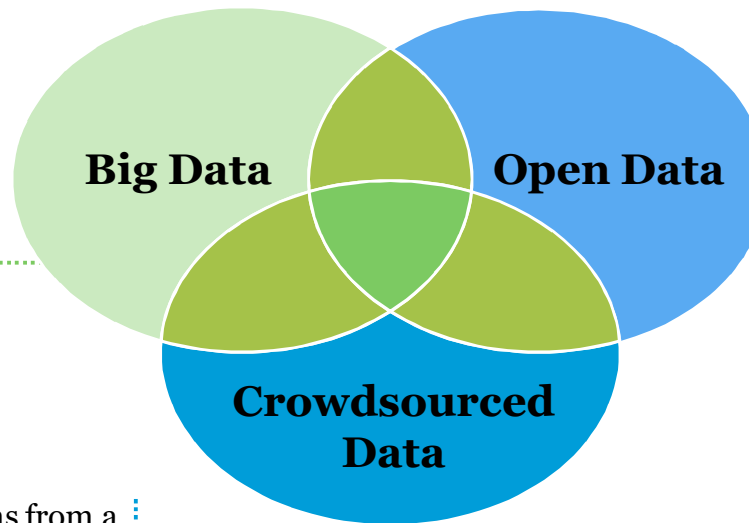
### New Business Models/Services

**Discovery of trends** will lead organizations to form new business models to adapt by creating new service offerings for their customers. Intermediary companies with big data expertise will provide analytics to 3<sup>rd</sup> parties.



## Not to be confused with...

Structured, semi-structured or unstructured information distinguished by one or more of the four “V”s: Veracity, Velocity, Variety, Volume.



Public, freely available data

Data collected through contributions from a large number of individuals.

# It's not just about the data...

It is important to understand the distinction between Big Data sets (large, unstructured, fast, and uncertain data) and 'Big Data Analytics'.

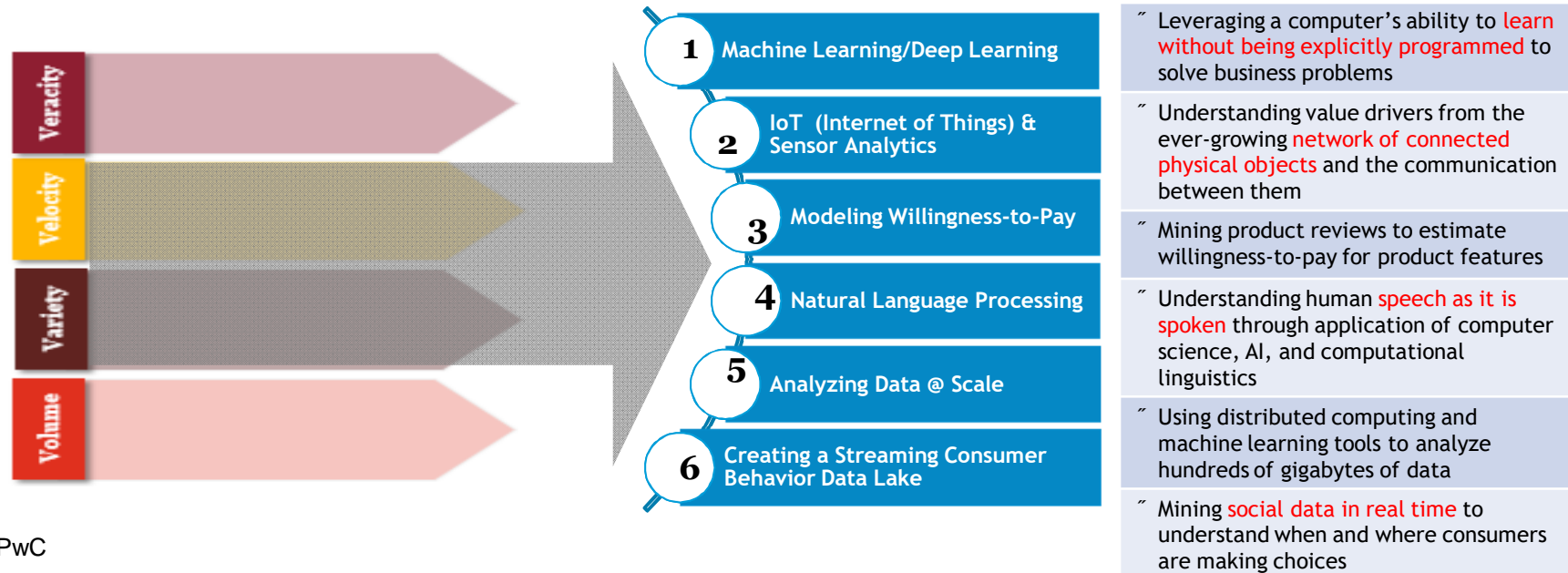
**BigData**

Refers to the DATA only

+


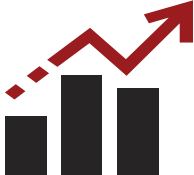



**BigData Analytics**

**Methods** of using Big Data to generate insight



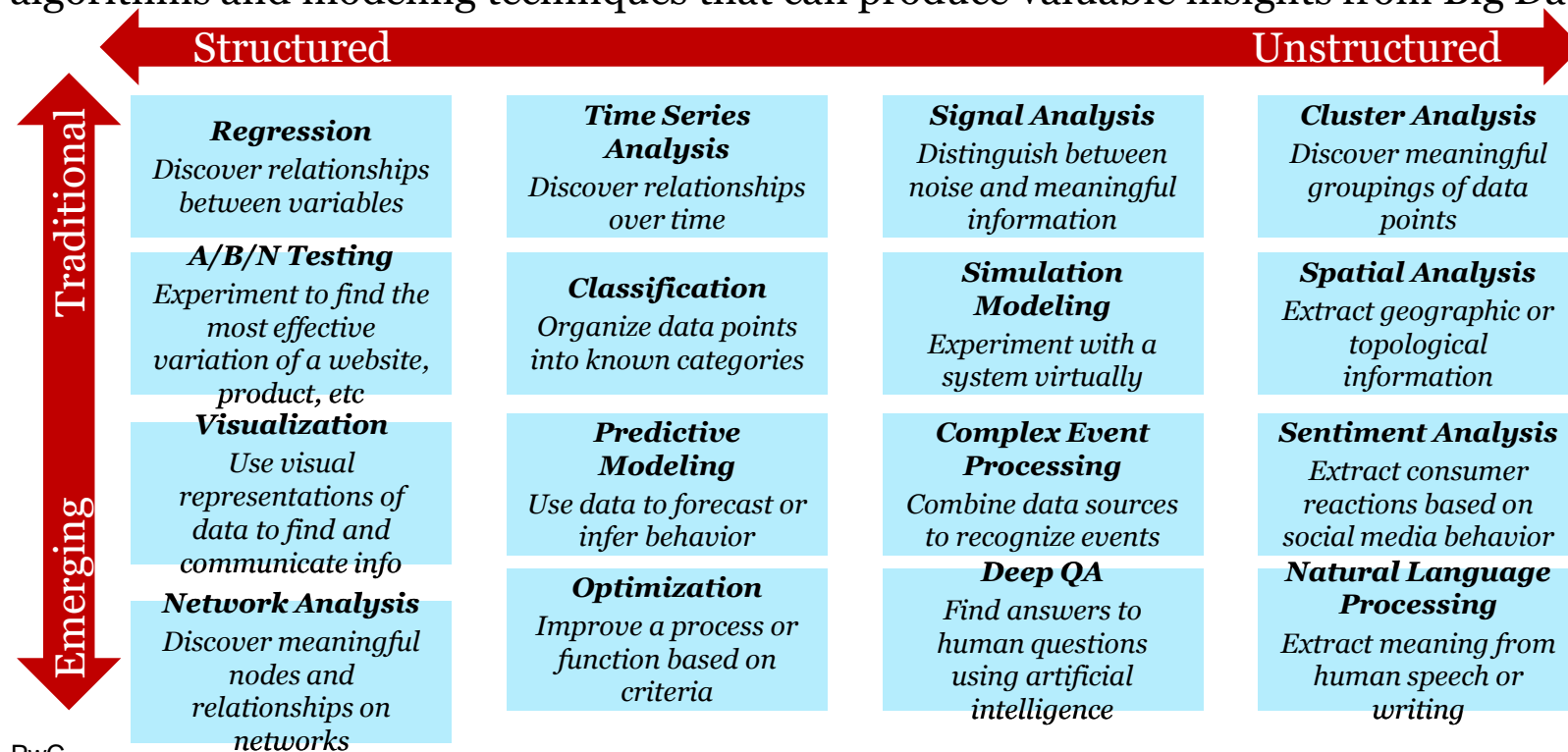
## ... It's also about what, how, and why you use it

BigData Analytics – the process of harnessing Big Data to yield actionable insights – is a combination of five key elements:

<i>Decisions</i>	<i>Analytics</i>	<i>Data</i>	<i>Technology</i>	<i>Mindset &amp; Skills</i>
<p>The value of Big Data Analytics is driven by the unique <b>decisions facing leaders</b>, companies, and countries today. In turn, the type, frequency, speed, and complexity of decisions drive how Big Data Analytics is deployed.</p>	<p>To <b>leverage the variety and volume</b> of Big Data while managing its volatility, advanced analytical approaches are necessary, such as natural language processing, network analysis, simulative modeling, artificial intelligence, etc.</p>	<p>Big Data Analytics is about operationalizing new and more data, but it is also about <b>data quality</b>, data interoperability, data disaggregation, and the ability to modularize data structures to <b>quickly absorb new data</b> and new types of data.</p>	<p>To store, manage, and use Big Data often requires investments in <b>new technologies</b> and data processing methods, such as distributed processing (e.g., Hadoop), NoSQL storage, and Cloud computing.</p>	<p>Big Data Analytics requires firm <b>commitment</b> to using analytics in decision-making; a decisive mentality capable of employing in-the-moment intelligence; and <b>investment</b> in analytical technology, resources, and skills.</p>
				

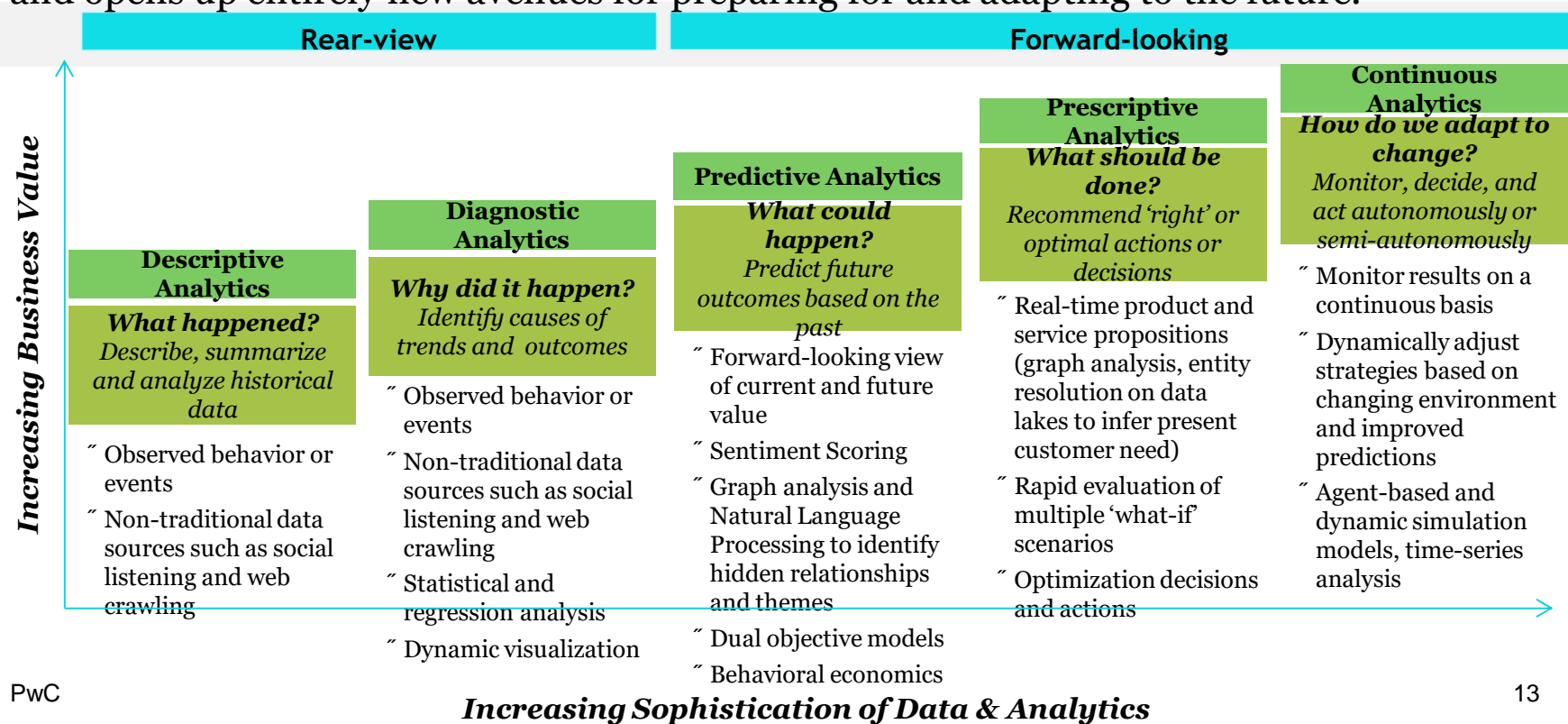
# BigData Analytical Capabilities

Continuing increases in processing capacity have opened the door to a range of advanced algorithms and modeling techniques that can produce valuable insights from Big Data.



# Forward-Looking vs. Rear-View Analytics

Big Data Analytics improves the speed and efficiency with which we understand the past, and opens up entirely new avenues for preparing for and adapting to the future.



# MAPREDUCE

THE BIRTH OF BIGDATA TECHNOLOGY

# MapReduce

- MapReduce [OSDI'04] provides
  - Automatic parallelization, distribution
  - I/O scheduling
    - Load balancing
    - Network and data transfer optimization
  - Fault tolerance
    - Handling of machine failures
- **Need more power: Scale out, not up!**
  - Large number of **commodity servers** as opposed to some high end specialized servers

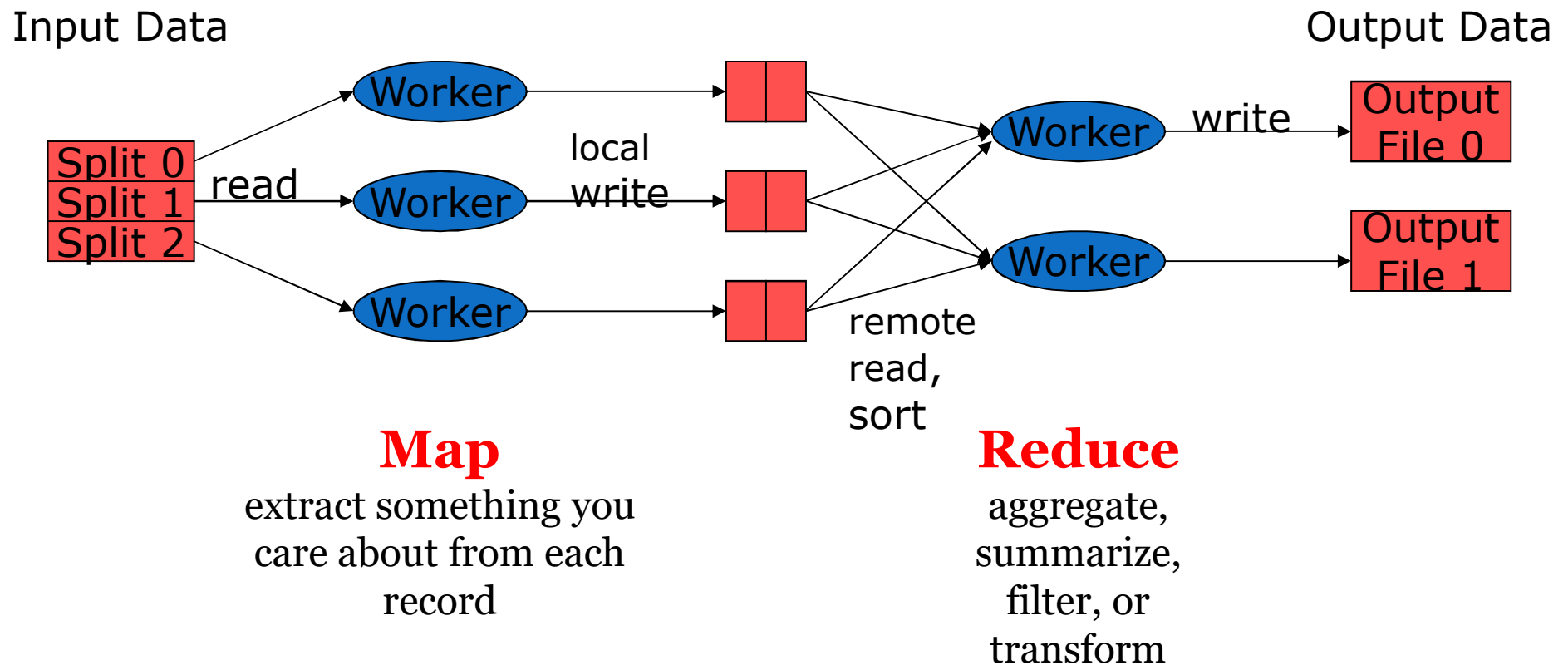
**Apache  
Hadoop:**  
Open source  
implementation of  
MapReduce

## Typical problem solved by MapReduce

- Read a lot of data
- **Map**: extract something you care about from each record
- Shuffle and Sort
- **Reduce**: aggregate, summarize, filter, or transform
- Write the results



# MapReduce workflow



# Mappers and Reducers

- Need to handle **more data**? Just add **more Mappers/Reducers**!
- No need to handle **multithreaded code** 😊
  - Mappers and Reducers are typically single threaded and **deterministic**
    - **Determinism** allows for restarting of failed jobs
  - Mappers/Reducers run **entirely independent** of each other
    - In Hadoop, they run in separate JVMs

# Example: Word Count

## Input Files

Apple Orange Mango  
Orange Grapes Plum

Apple Plum Mango  
Apple Apple Plum

<http://kickstarthadoop.blogspot.ca/2011/04/word-count-hadoop-map-reduce-example.html>

# 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

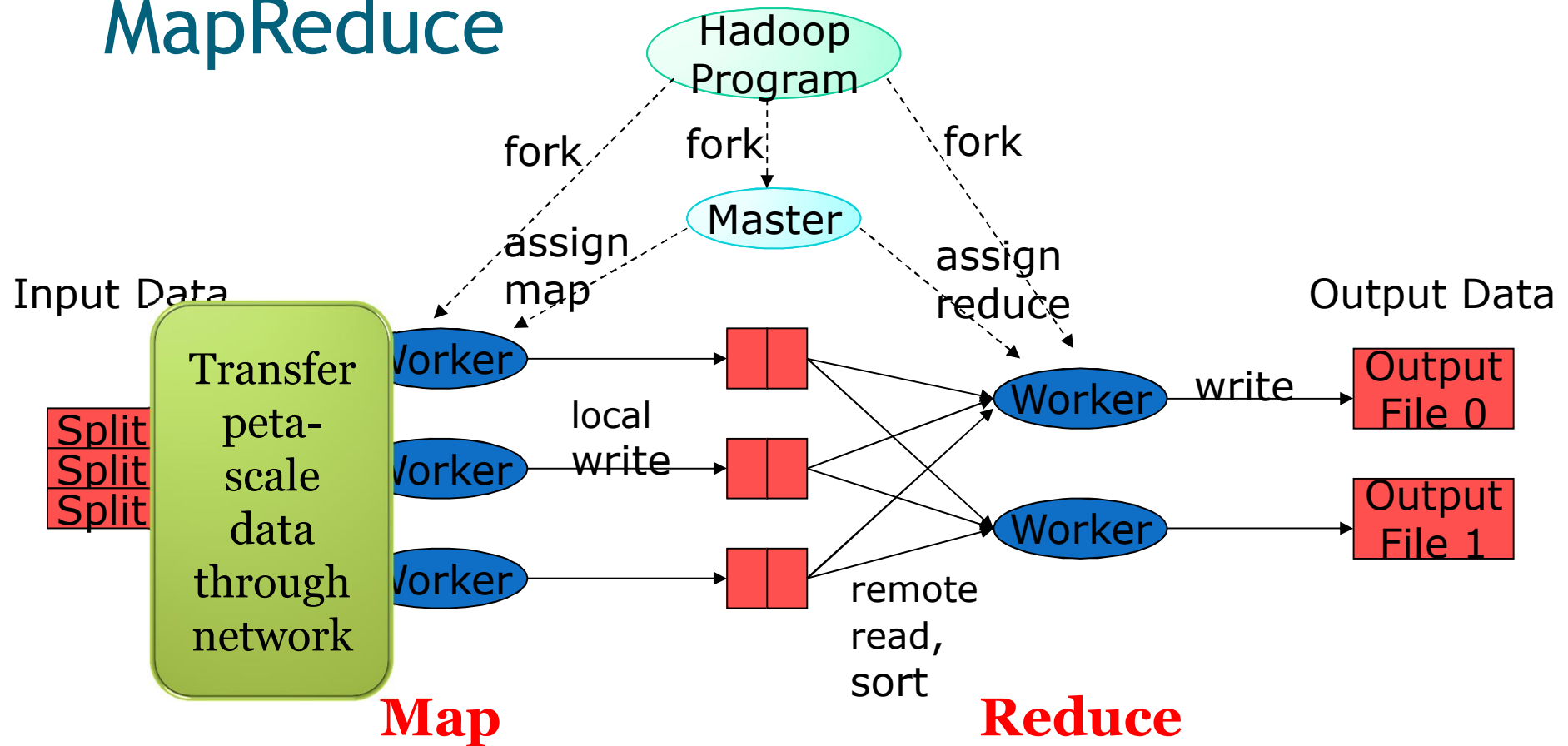
- For our example, the reducer input would be:

<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>

- The output would be:

<The, 6> <teacher, 1> <went, 1> <to, 1> **<store, 3>** <was, 1>  
<closed, 1> <opens, 1> <morning, 1> <at, 1> <9am, 1>

# MapReduce



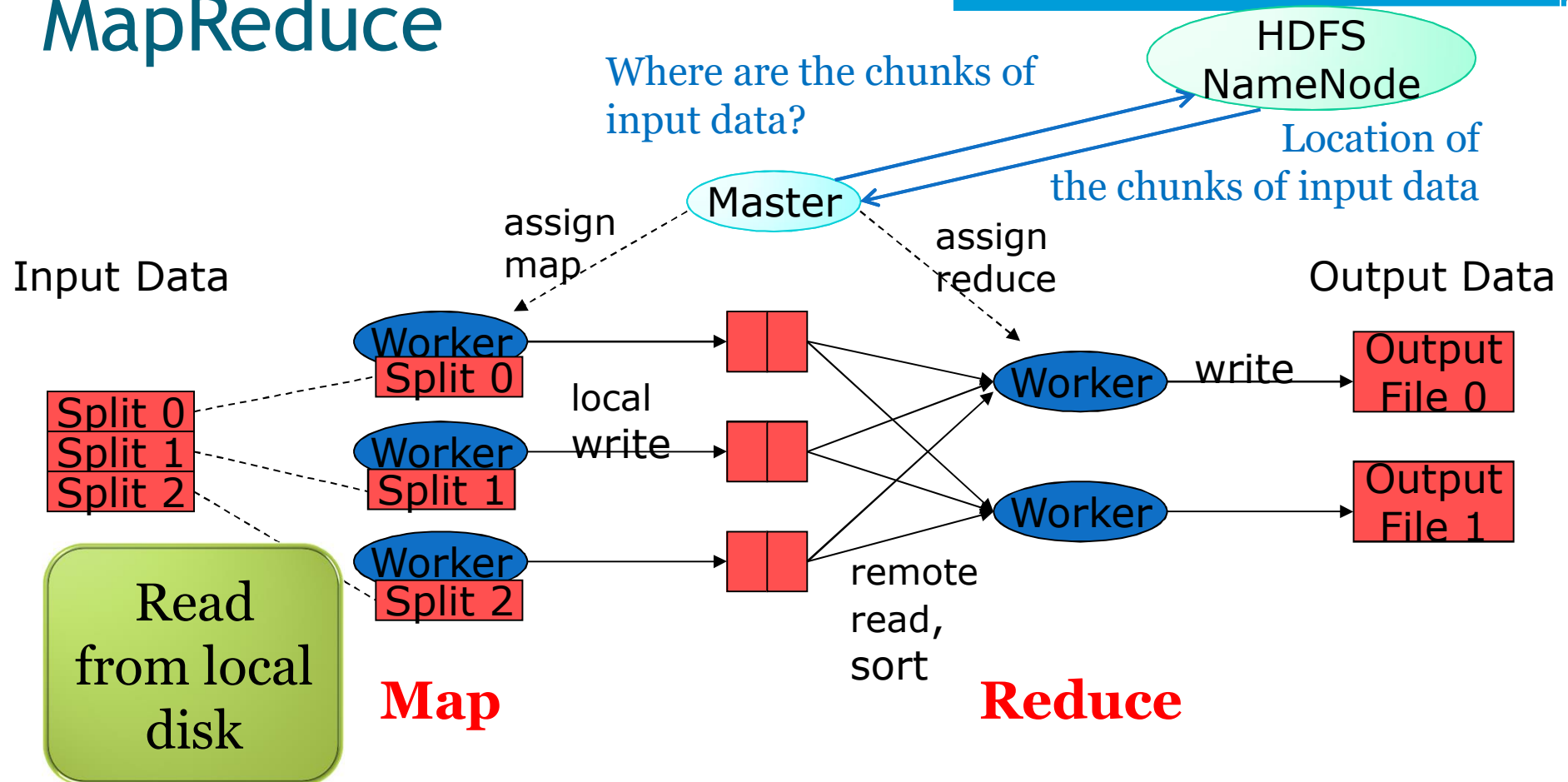
# Google File System (GFS)

## Hadoop Distributed File System (HDFS)

- Split data and store 3 replica on commodity servers



# MapReduce





# Locality Optimization

- **Master scheduling policy:**
  - Asks GFS for locations of replicas of input file blocks
  - Map tasks scheduled so GFS input block replica are on same machine or same rack
- Effect: Thousands of machines read input at local disk speed
  - Eliminate network bottleneck!

# Failure in MapReduce

- **Failures** are **norm** in commodity hardware
- **Worker** failure
  - Detect failure via periodic heartbeats
  - Re-execute in-progress map/reduce tasks
- **Master** failure
  - Single point of failure; Resume from Execution Log
- **Robust**
  - Google's experience: lost 1600 of 1800 machines once!, but finished fine

# Fault tolerance: Handled via re-execution

- On worker **failure**:
  - Detect failure via periodic heartbeats
  - Re-execute completed and in-progress *map* tasks
  - Task completion committed through master
- Robust: [Google's experience] lost 1600 of 1800 machines, but finished **fine**

# Refinement: Redundant Execution

- **Slow workers** significantly lengthen completion time
  - Other jobs consuming resources on machine
  - Bad disks with soft errors transfer data very slowly
  - Weird things: processor caches disabled (!!)
- **Solution**: spawn backup copies of tasks
  - Whichever one finishes first "wins"

## Refinement: Skipping Bad Records

Map/Reduce functions sometimes fail for particular inputs

- Best solution is to debug & fix, but not always possible
- If master sees **two failures** for the **same record**:
  - Next worker is told to **skip the record**

# Summary

- MapReduce
  - Programming paradigm for data-intensive computing
  - Distributed & parallel execution model
  - Simple to program
    - The framework automates many tedious tasks (machine selection, failure handling, etc.)

# Details of GFS

# Motivation: Large Scale Data Storage

- Manipulate large (**Peta Scale**) sets of data
- Large number of machines with **commodity hardware**
- Component failure is the norm
- Goal: **Scalable, high performance, fault tolerant** distributed file system

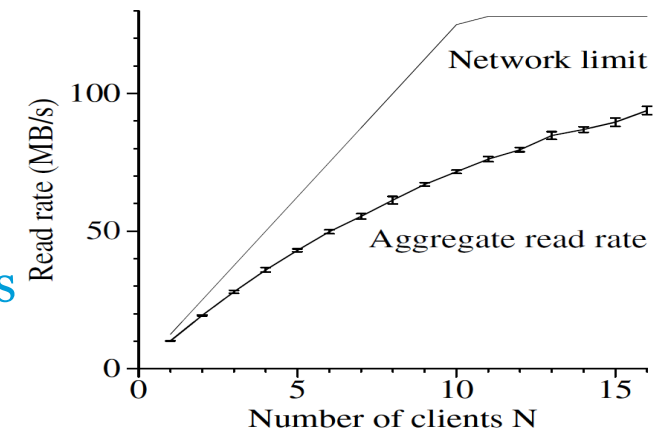


## Why a new file system?

- None designed for their failure model
- Few scale as highly or dynamically and easily
- Lack of special primitives for large distributed computation

# What should expect from GFS

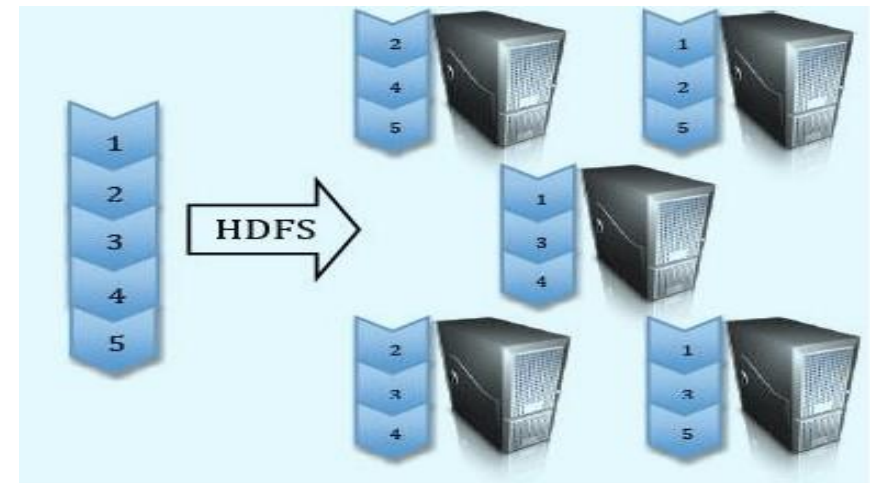
- Designed for Google's application
  - Control of both file system and application
  - Applications use a few specific access patterns
    - Append to large files
    - Large streaming reads
  - **Not** a good fit for
    - low-latency data access
    - lots of small files, multiple writers, arbitrary file modifications
- Not POSIX, although mostly traditional
  - Specific operations: RecordAppend



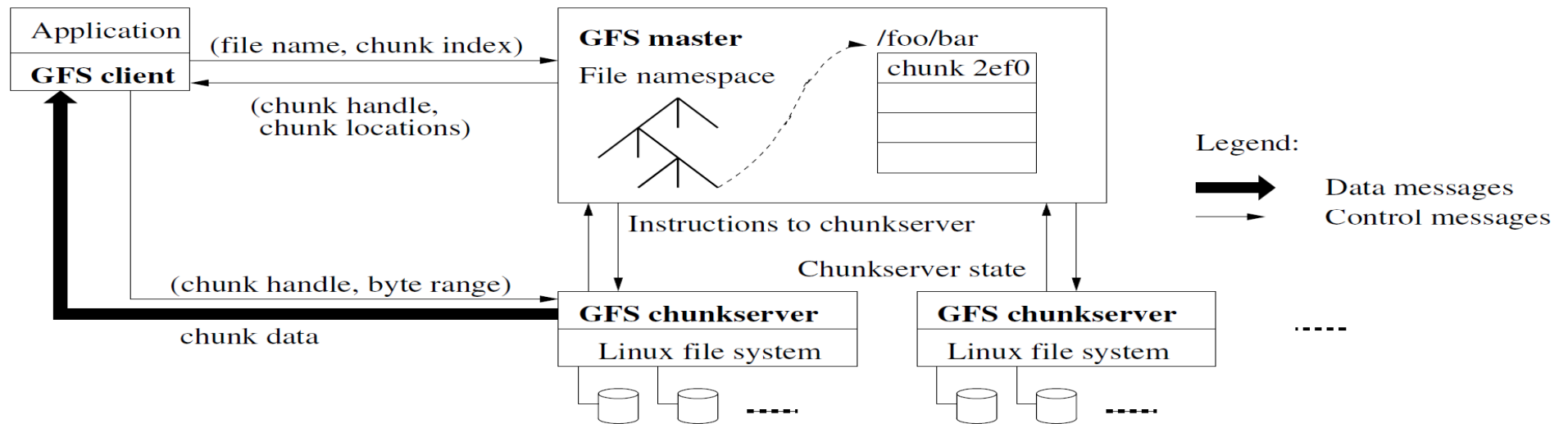
- **Different** characteristic than **transactional** or the “customer order” data : “**write once read many (WORM)**”
  - e.g. web logs, web crawler’s data, or healthcare and patient information
  - WORM inspired MapReduce programming model
- Google exploited this characteristics in its Google file system [SOSP’03]
  - **Apache Hadoop: Open source project**

# Components

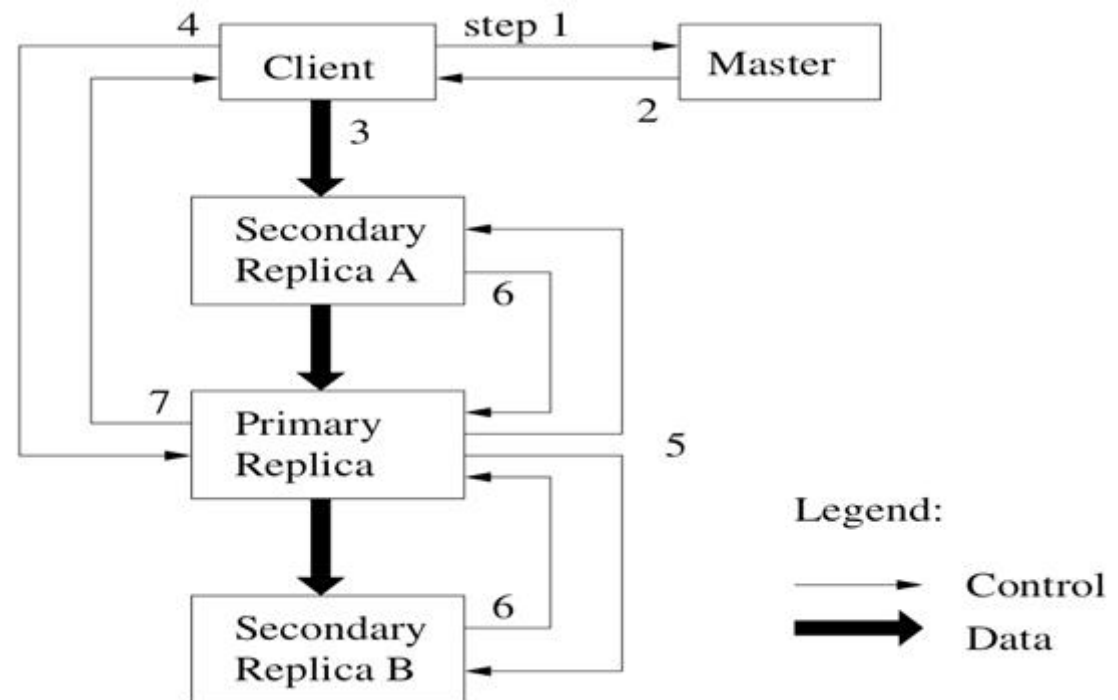
- **Master (NameNode)**
  - Manages metadata (namespace)
  - Not involved in data transfer
  - Controls allocation, placement, replication
- **Chunkserver (DataNode)**
  - Stores chunks of data
  - No knowledge of GFS file system structure
  - Built on local Linux file system



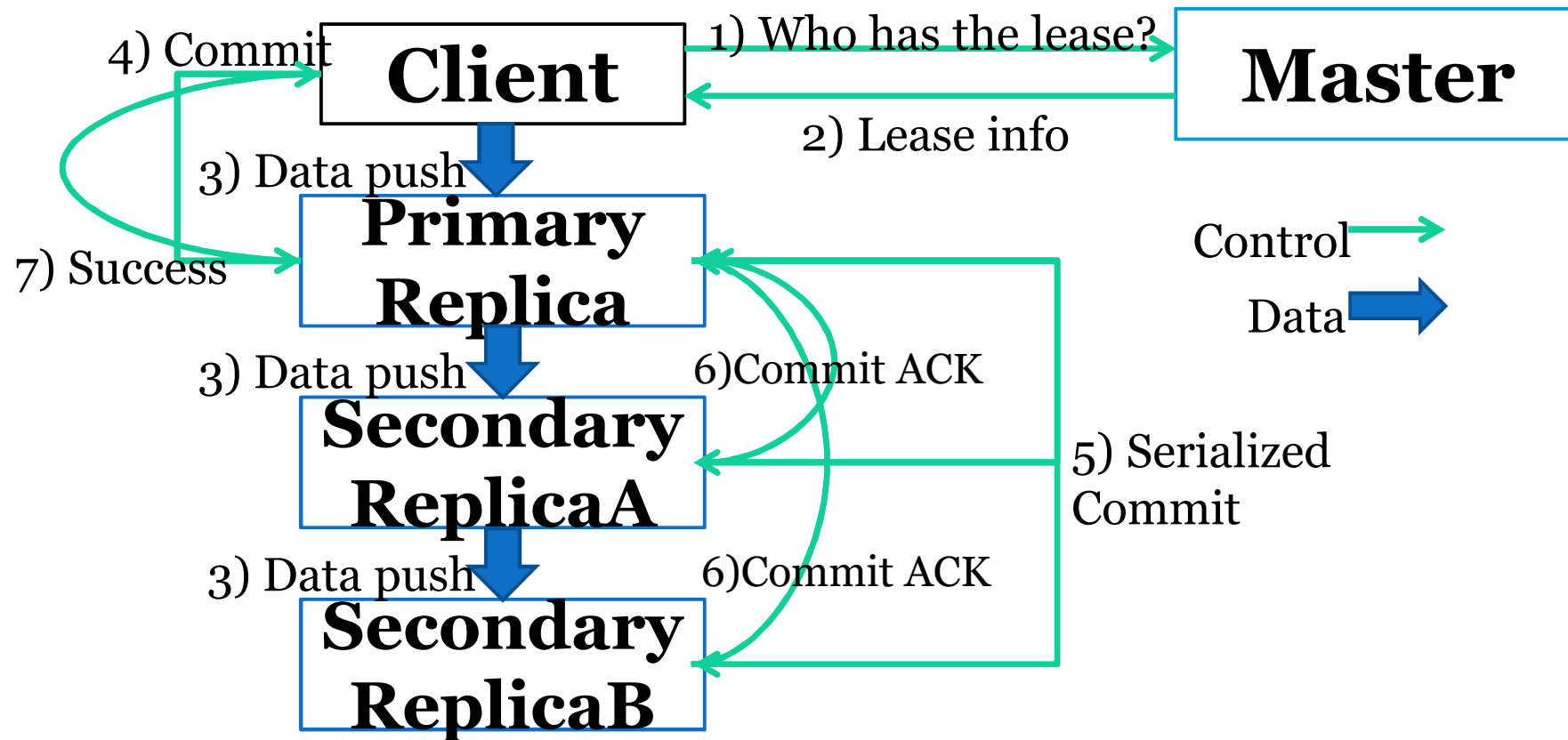
# GFS Architecture



# Write operation



# Write(filename, offset, data)



# RecordAppend(filename, data)

- Significant use in distributed apps. For example at Google production cluster:
  - 21% of bytes written
  - 28% of write operations
- **Guaranteed:** All data appended at least once as a single consecutive byte range
- Same basic structure as write
  - Client obtains information from master
  - Client sends data to data nodes (chunkservers)
  - Client sends “append-commit”
  - Lease holder serializes append
- **Advantage:** Large number of concurrent writers with minimal coordination



## RecordAppend (2)

- Record size is limited by chunk size
- When a record does not fit into available space,
  - chunk is padded to end
  - and client retries request.

# Fault tolerance

- Replication
  - High availability for reads
  - User controllable, default 3 (non-RAID)
  - Provides read/seek bandwidth
  - Master is responsible for directing re-replication if a data node dies
- Online checksumming in data nodes
  - Verified on reads

# Replica Management

- Bias towards topological spreading
  - Rack, data center
- Rebalancing
  - Move chunks around to balance disk fullness
  - Gently fixes imbalances due to:
    - Adding/removing data nodes

# Replica Management (Cloning)

- Chunk replica lost or corrupt
- **Goal:** minimize app disruption and data loss
  - Approximately in priority order
    - More replica missing-> priority boost
    - Deleted file-> priority decrease
    - Client blocking on a write-> large priority boost
  - Master directs copying of data
- Performance on a production cluster
  - Single failure, full recovery (600GB): 23.2 min
  - Double failure, restored 2x replication: 2min

# Garbage Collection

- Master does **not** need to have a **strong knowledge** of what is stored on each data node
  - Master regularly scans namespace
  - After GC interval, deleted files are removed from the namespace
  - Data node periodically polls Master about each chunk it knows of.
  - If a chunk is forgotten, the master tells data node to delete it.

# Limitations

- Master is a central point of failure
- Master can be a scalability bottleneck
- Latency when opening/stating thousands of files
- Security model is weak

## Conclusion

- Inexpensive commodity components can be the basis of a large scale reliable system
- Adjusting the API, e.g. RecordAppend, can enable large distributed apps
- Fault tolerant
- Useful for many similar apps