

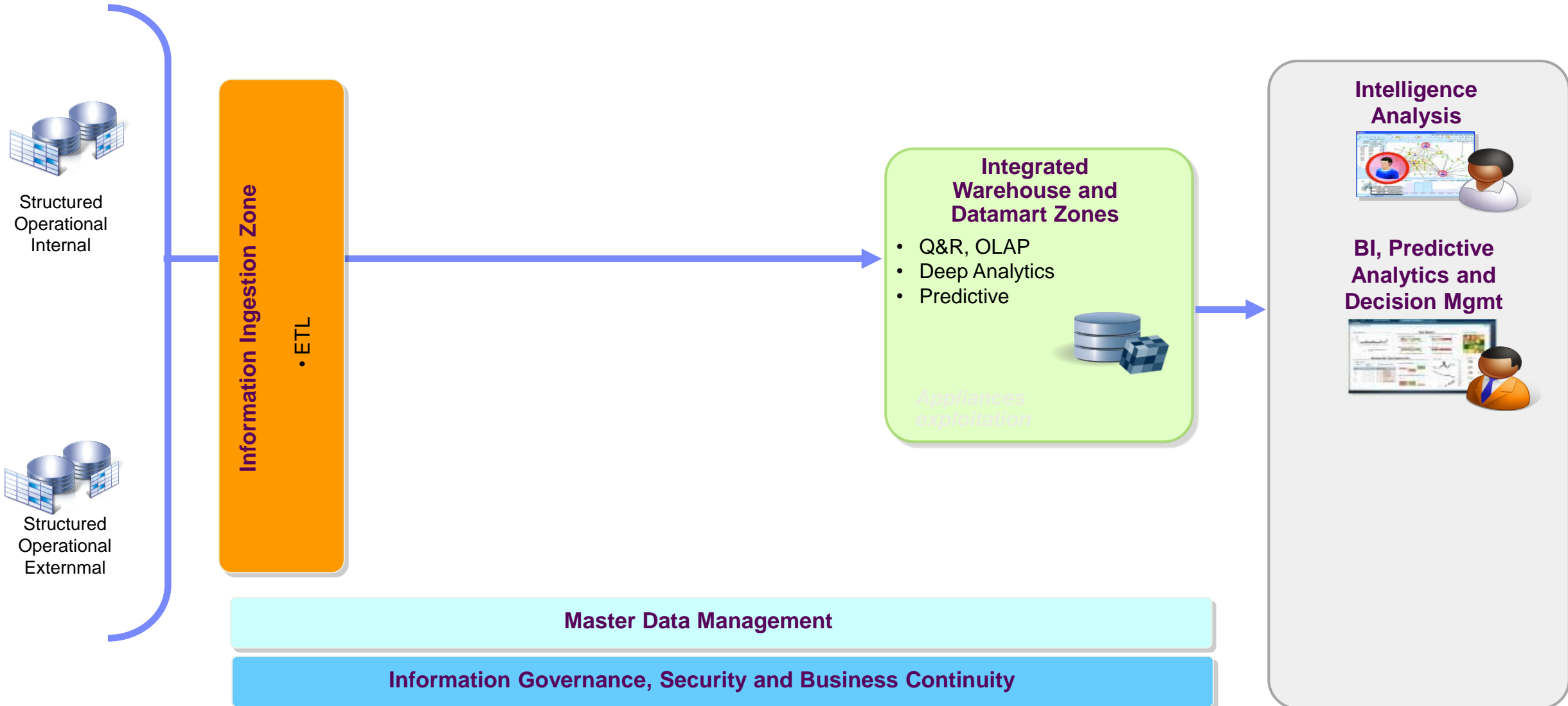
# Infosphere BigInsights

Carlo Patrini – IBM Analytics ([carlo.patrini@it.ibm.com](mailto:carlo.patrini@it.ibm.com))

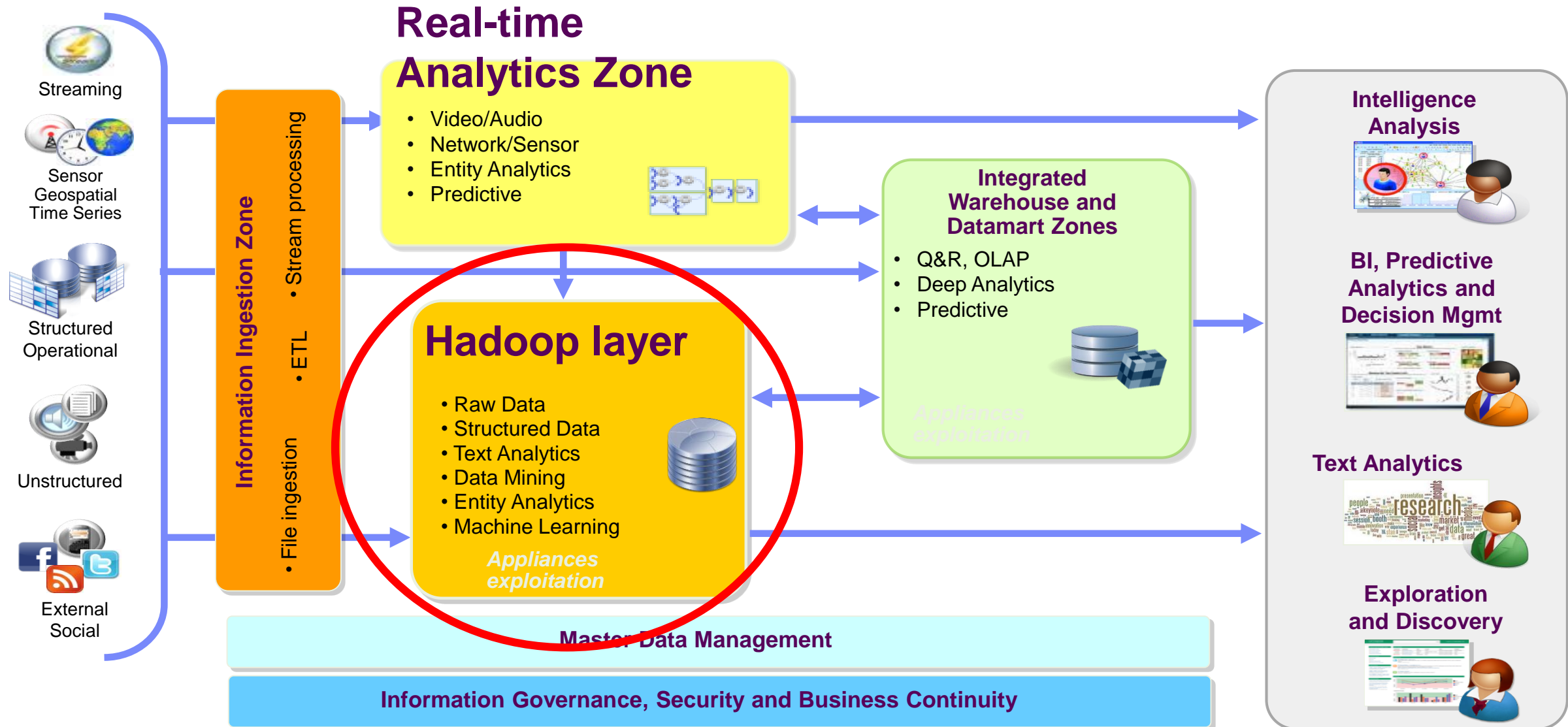
Francesco Airoidi – IBM Sales & Distribution ([airfranz@it.ibm.com](mailto:airfranz@it.ibm.com))

Giugno 2015

# La tradizionale filiera BI



# La filiera delle IBM Big Data & Analytics



## ■ CPU istruzioni al secondo – miglioramenti significativi

1990 44 Mips at 40 Mhz

2000 3.562 Mips at 1.2 Ghz

2010 147.600 Mips at 3.3 Ghz

## ■ RAM Memory - miglioramenti significativi

1990 640 K

2000 64 Mb

2010 8-32 GB

## ■ Disk capacity - miglioramenti significativi

1990 20 MB

2000 10 GB

2010 1 TB

## ■ Disk latency (velocità di leggere e scrivere su disco ) - miglioramenti poco significativi

Negli ultimi 7-10 anni non ci sono state enormi migliorie  
correntemente la velocità è di circa 70 – 80 MB / sec

# Quanto tempo ci vuole per scandire 1 TB ?

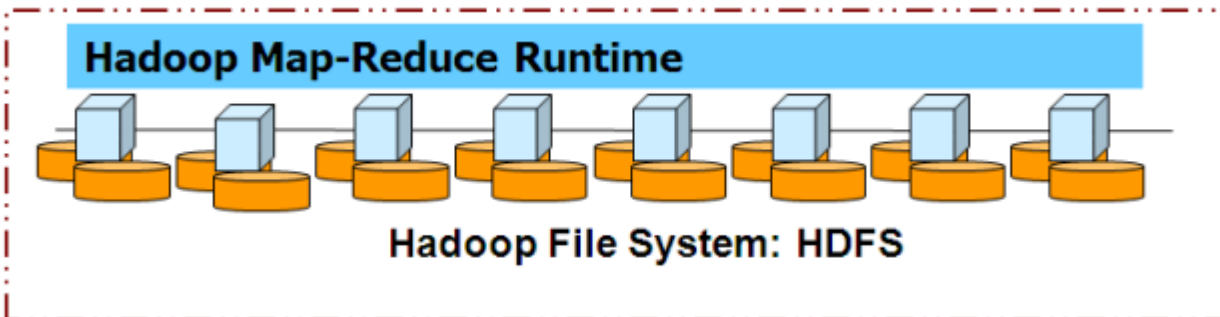
## ❑ 1 TB (at 80 MB / sec)

1 disk	3.4 hours
10 disks	20 min
100 disks	2 min
1000 disks	12 sec

## ❑ Per ovviare alla Disc Latency la risposta è la ..elaborazione parallela

## ❑ Hadoop : un nuovo modo per memorizzare ed elaborare i dati

- Scritto in Java
- Progettato per lavorare su hardware non specializzato
- Gira in ambiente Linux
- Scalabile, Flessibile, Robusto



# What is Hadoop?

- **Apache Hadoop = free, open source framework** for data-intensive applications
  - Inspired by **Google technologies** (MapReduce, GFS)
  - **Yahoo has been** the largest contributor to the project (Doug Cutting),
  - Well-suited to batch-oriented, read-intensive applications
  - Originally built to address scalability problems of Nutch, an open source Web search technology
- Enables applications to work with thousands of nodes and petabytes of data in a highly parallel, cost effective manner
  - CPU + disks of commodity box = Hadoop “node”
  - Boxes can be combined into clusters
  - New nodes can be added as needed without changing
    - Data formats
    - How data is loaded
    - How jobs are written



## ■ **MapReduce framework**

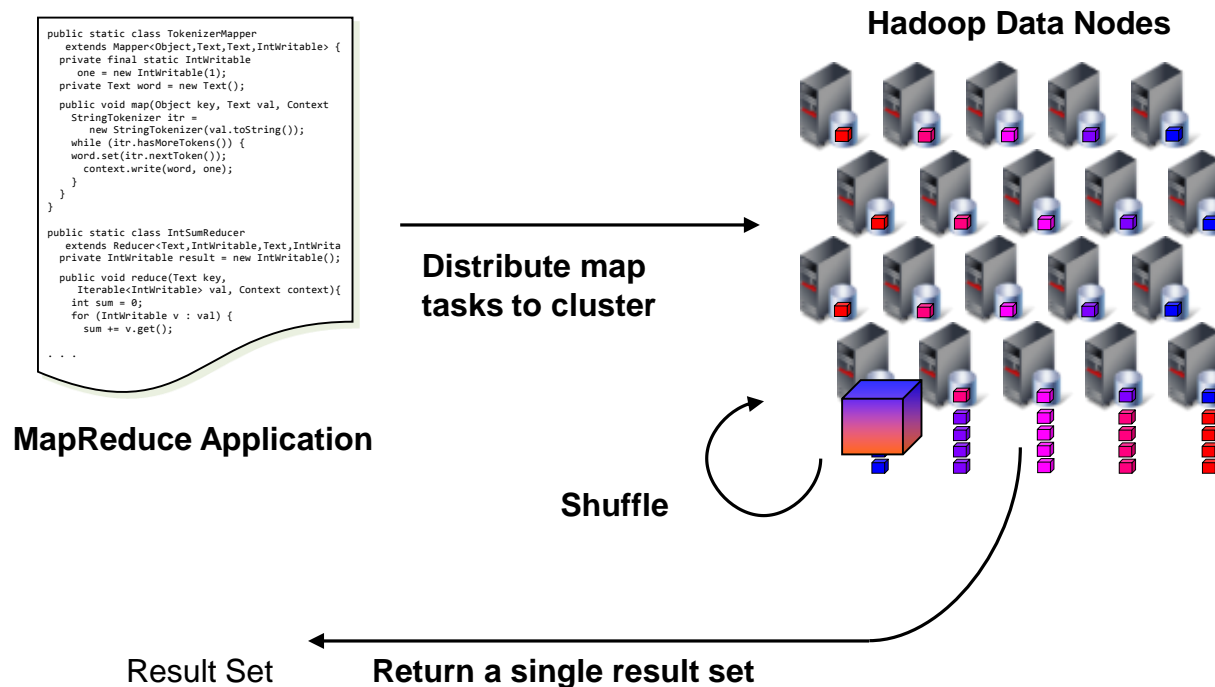
- MapReduce is a software framework introduced by Google to support distributed computing on large data sets of clusters of computers.
- How Hadoop understands and assigns work to the nodes (machines)

## ■ **Hadoop Distributed File System = HDFS**

- Where Hadoop stores data
- A file system that spans all the nodes in a Hadoop cluster
- It links together the file systems on many local nodes to make them into one big file system

# Hadoop ed il paradigma Map Reduce

- I dati sono memorizzati su un sistema distribuito di server
- Le funzioni elaborative vengono inviate dove ci sono I dati
- Ogni server elabora I dati di propria competenza e condivide i risultati
- Il sistema può scalare raggiungendo migliaia di nodi e PB di dati

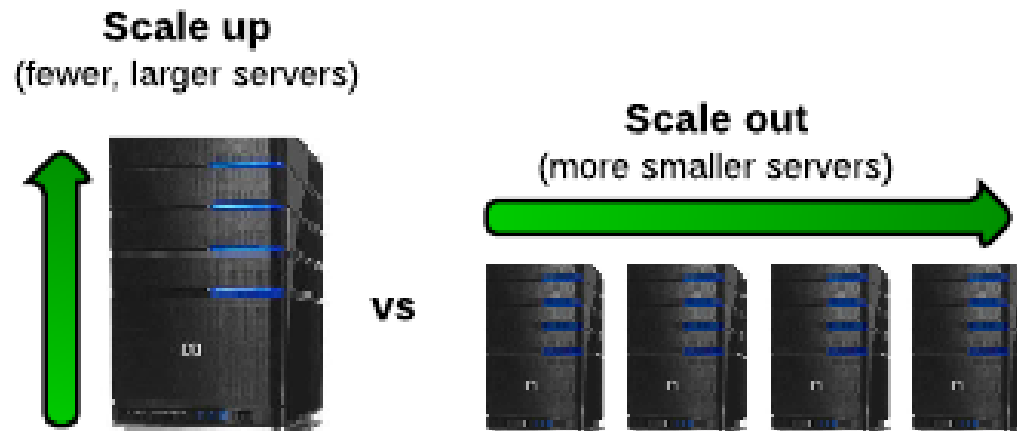


1. Map Phase  
(spezza il job in piccole parti)
2. Shuffle  
(riordina I risultati parziali per le elaborazione finale)
3. Reduce Phase  
(rielabora il tutto per ottenere un singolo risultato)



# Scale-out vs scale-up

- ❑ **Just reading 100 terabytes is slow**
  - Standard computer (100 MBPS) ~11 days
  - Across 10Gbit link (high end storage) 1 day
  - 1000 standard computers 15 minutes!
- ❑ **Seek times for random disk access is a problem**
  - 1 TB data set with  $10^{10}$  100-byte records
    - Updates to 1% would require 1 month*
    - Reading and rewriting the whole data set would take 1 day\**
- ❑ **One node is not enough!**
- ❑ **Need to scale out not up!**



- **Bad news: nodes fail, especially if you have many**
  - Mean time between failures for 1 node = 3 years, 1000 nodes = 1 day
  - Super-fancy hardware still fails and commodity machines give better performance per dollar
  
- **Bad news II: distributed programming is hard**
  - Communication, synchronization, and deadlocks
  - Recovering from machine failure
  - Debugging
  - Optimization
  
- **Bad news III: repeat for every problem**

In 2003/4 Google publishes seminal whitepapers on a new programming paradigm to handle data at Internet Scale.

Original assumptions:

- scale-out architecture
- build a super computer on commodity hardware (commodity != low-end - not tied to expensive, proprietary offerings from a single vendor)
- hides system-level details from the developers (the datacenter is the computer)
- move processing to the data (because cluster have limited bandwidth)
- process data sequentially (seeks are expensive)
- expect failure
- shared-nothing architecture (processing tasks have no dependency on one other)

## **The Google File System**

Sanjay Ghemawat, Howard Gobioff, and Shun-Tak Leung  
Lake George, NY, October, 2003.

## **MapReduce: Simplified Data Processing on Large Clusters**

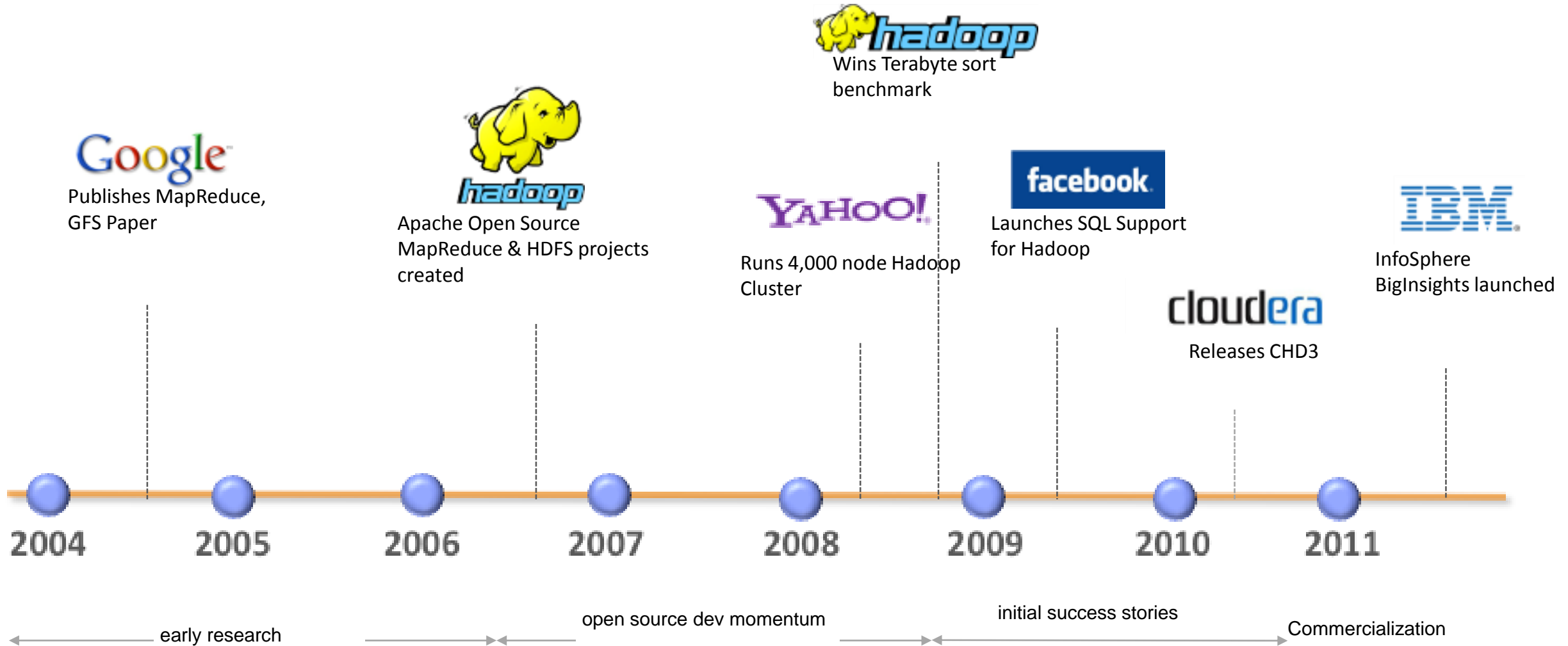
Jeffrey Dean and Sanjay Ghemawat  
San Francisco, CA, December, 2004.

**“A simple and powerful interface that enables automatic parallelization and distribution of large-scale computations, combined with an implementation of this interface that achieves high performance on large clusters of commodity PCs”**

## **Bigtable: A Distributed Storage System for Structured Data**

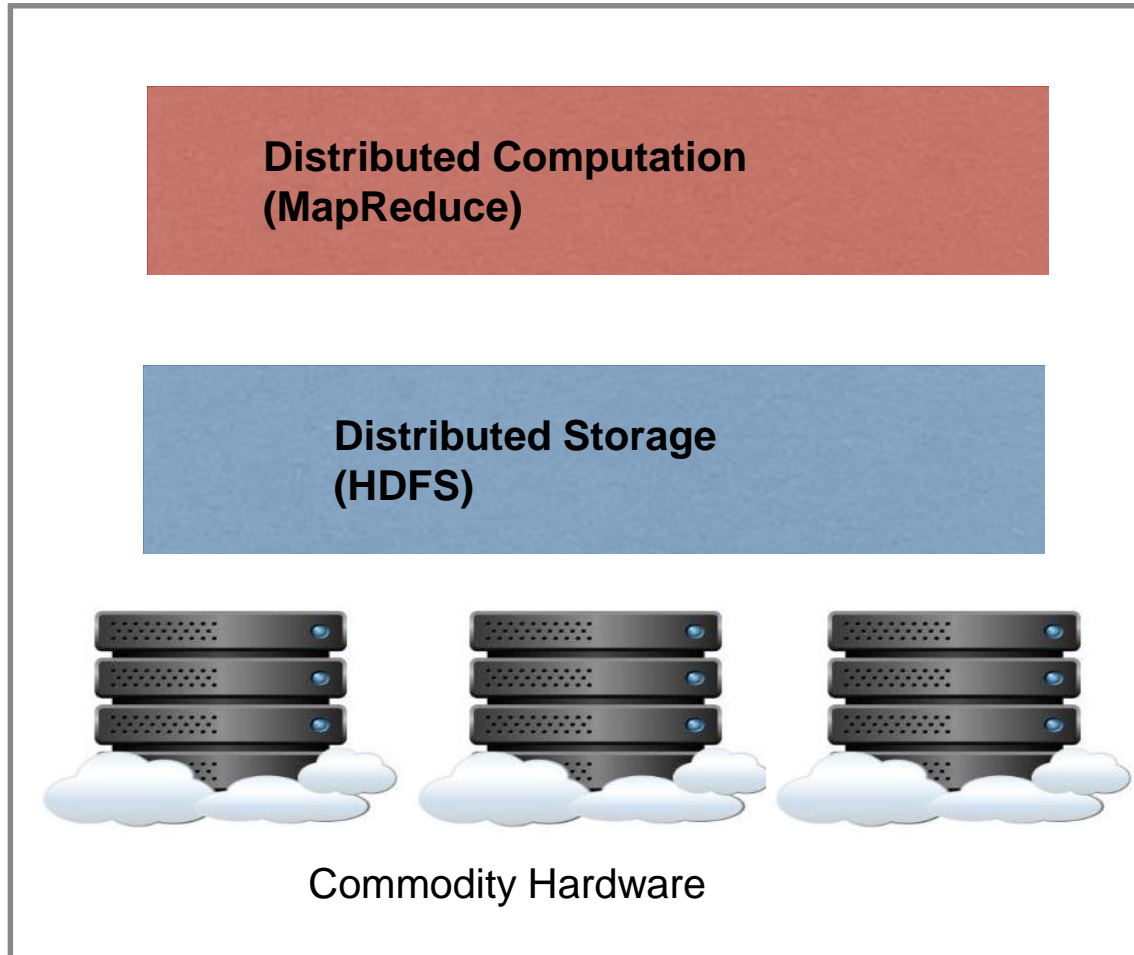
Fay Chang, Jeffrey Dean, Sanjay Ghemawat, Wilson C. Hsieh, Deborah A. Wallach, Mike Burrows, Tushar Chandra, Andrew Fikes, and Robert E. Gruber  
Seattle, WA, November, 2006.

# A brief history of Hadoop



# So what exactly is Hadoop?

A framework for running applications (aka jobs) on large clusters built on **commodity hardware** capable of processing **petabytes of data**



- ✦ it implements a computational paradigm named **Map/Reduce**, where the application is divided into self contained units of work, each of which may be executed or re-executed on any node in the cluster
- ✦ it provides a distributed file system (**HDFS**) that stores data on the compute nodes, providing very high aggregate bandwidth across the cluster.
- ✦ **node failures are automatically handled** by the framework.

## The largest known production environments:

- By 2011: Facebook, which consisted of 4,000 nodes (including 8- and 16-core CPUs) and supported 21PB of storage
- By (October) 2013: Yahoo!, spanning more than 35,000 nodes

## Facebook Datawarehousing Hadoop cluster (2011 data):

- 12 TB of compressed data added per day
- 800 TB of compressed data scanned per day
- 25,000 map-reduce jobs per day
- 65 millions files in HDFS
- 30,000 simultaneous clients to the HDFS NameNode

Nowadays Facebook Datawarehouse is 300 PB and add about 600 TB of compressed data per day

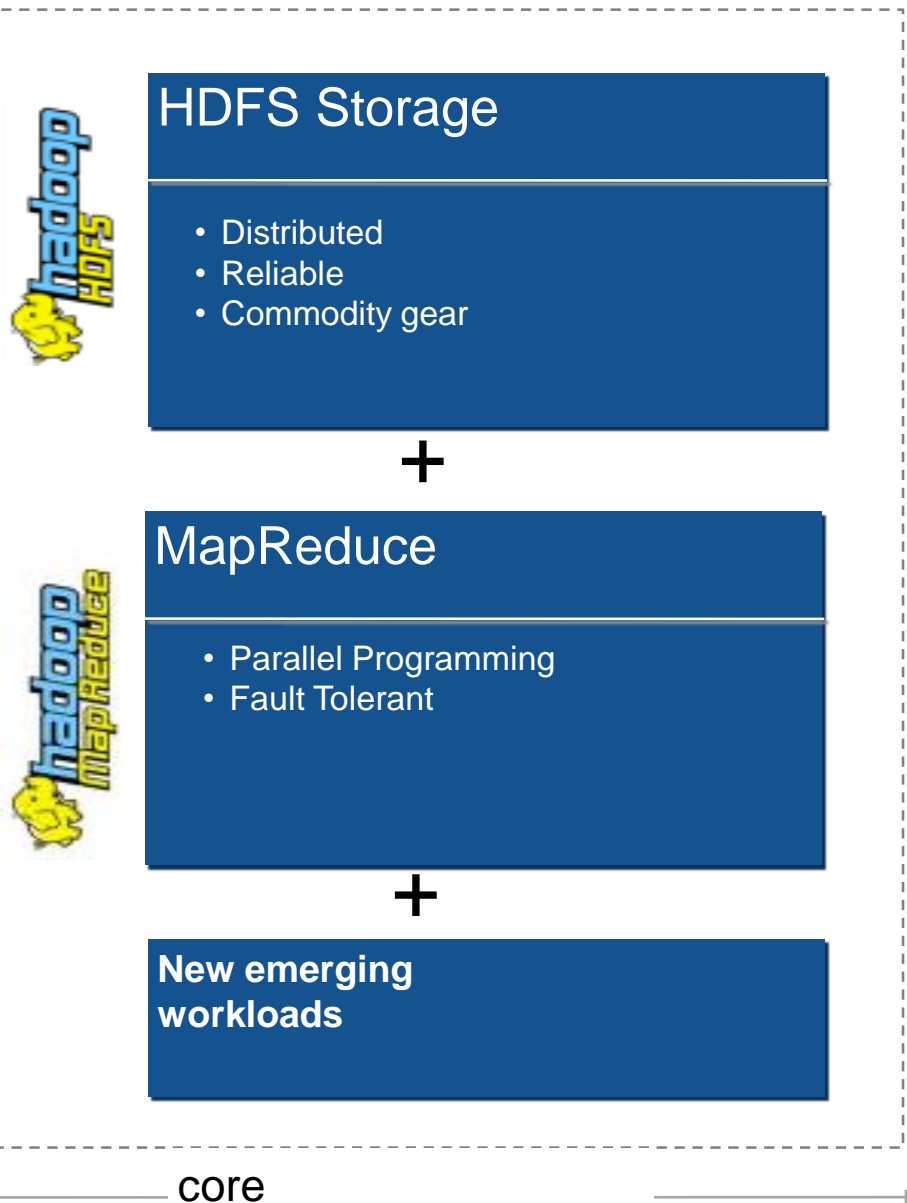
## Not good for

- Not to process transactions (random access)
- Not good when work cannot be parallelized
- Not good for low latency data access
- Not good for processing lots of small files
- Not good for intensive calculations with little data

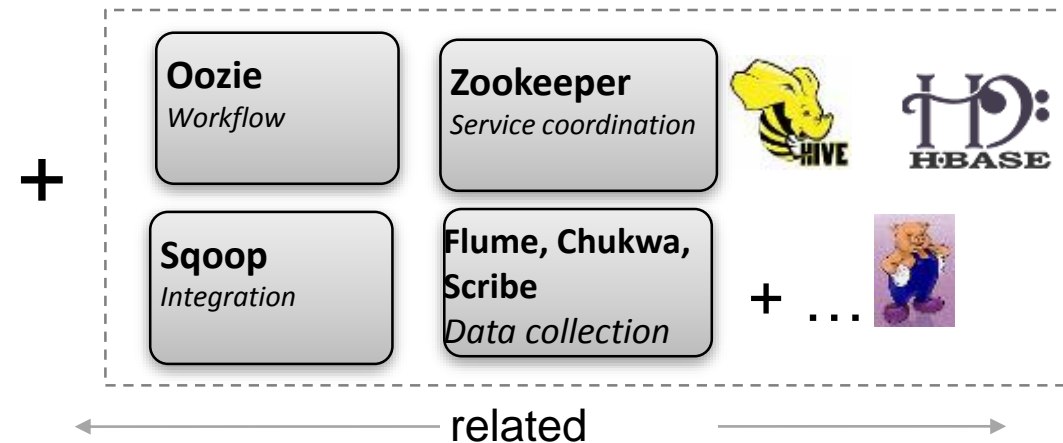
## Good for

- ◆ Data intensive, bulk-processing oriented, typically “embarrassingly parallel”
- ◆ Examples
  - ☞ Index building at Google and Yahoo!
  - ☞ Article clustering
  - ☞ Statistical machine translation
  - ☞ Spam detection
  - ☞ Ad optimization
  - ☞ Natural Language Processing
  - ☞ Image analysis
  - ☞ OCR
  - ☞ IBM’s Watson





Hadoop term is also used for a family of related projects that fall under the umbrella of infrastructure for distributed computing and large-scale data processing





## IBM Open Platform with Apache Hadoop

Ambari\*

Avro

Flume

Hadoop

HDFS/MapReduce/YARN\*

Hive

Knox

Open JD

Oozie

Slider

Snappy

Solr

Sqoop

# HDFS - Terminology review



Node 1

# HDFS - Terminology review

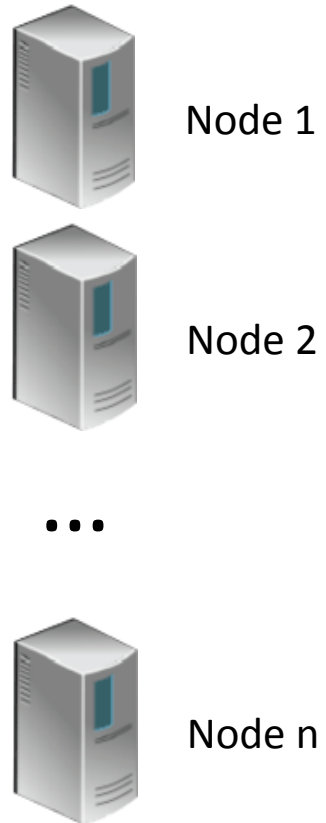


Node 1

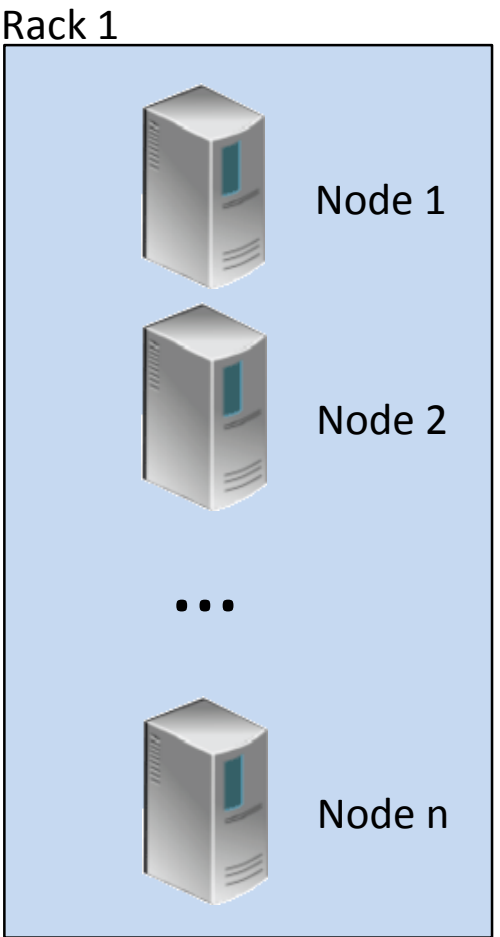


Node 2

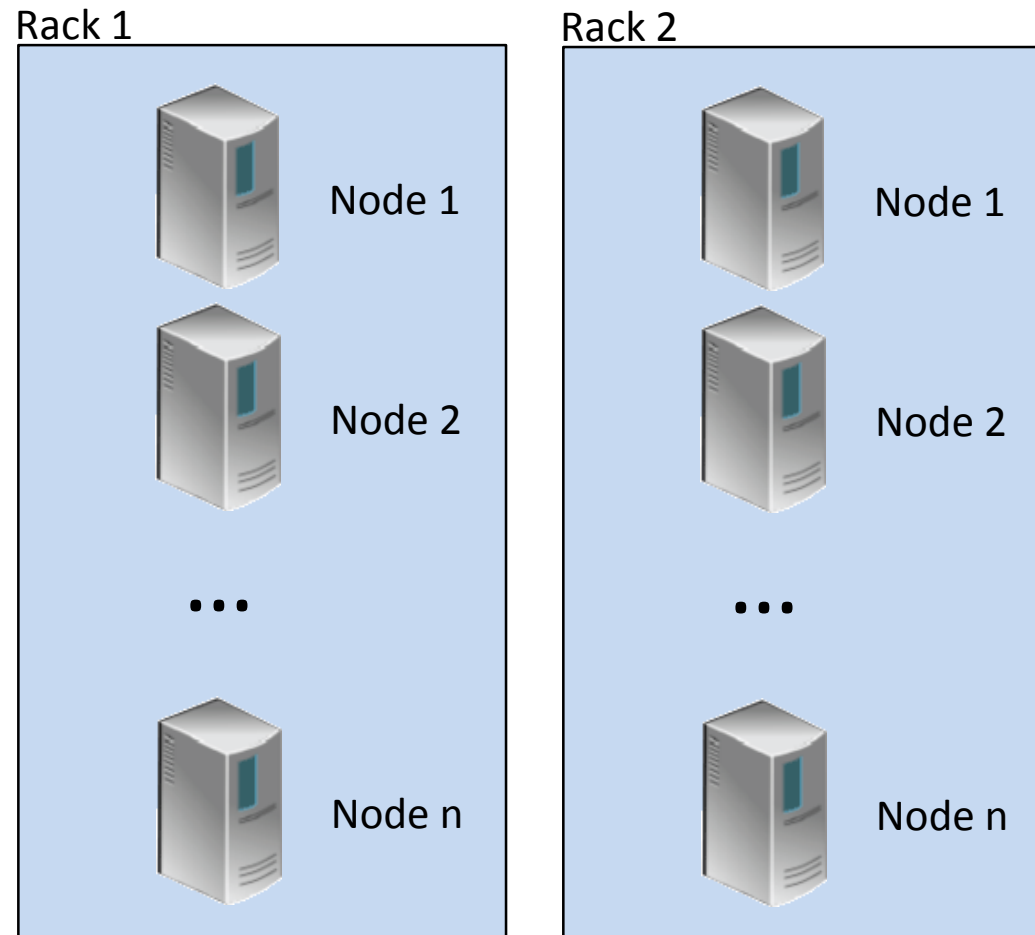
# HDFS - Terminology review



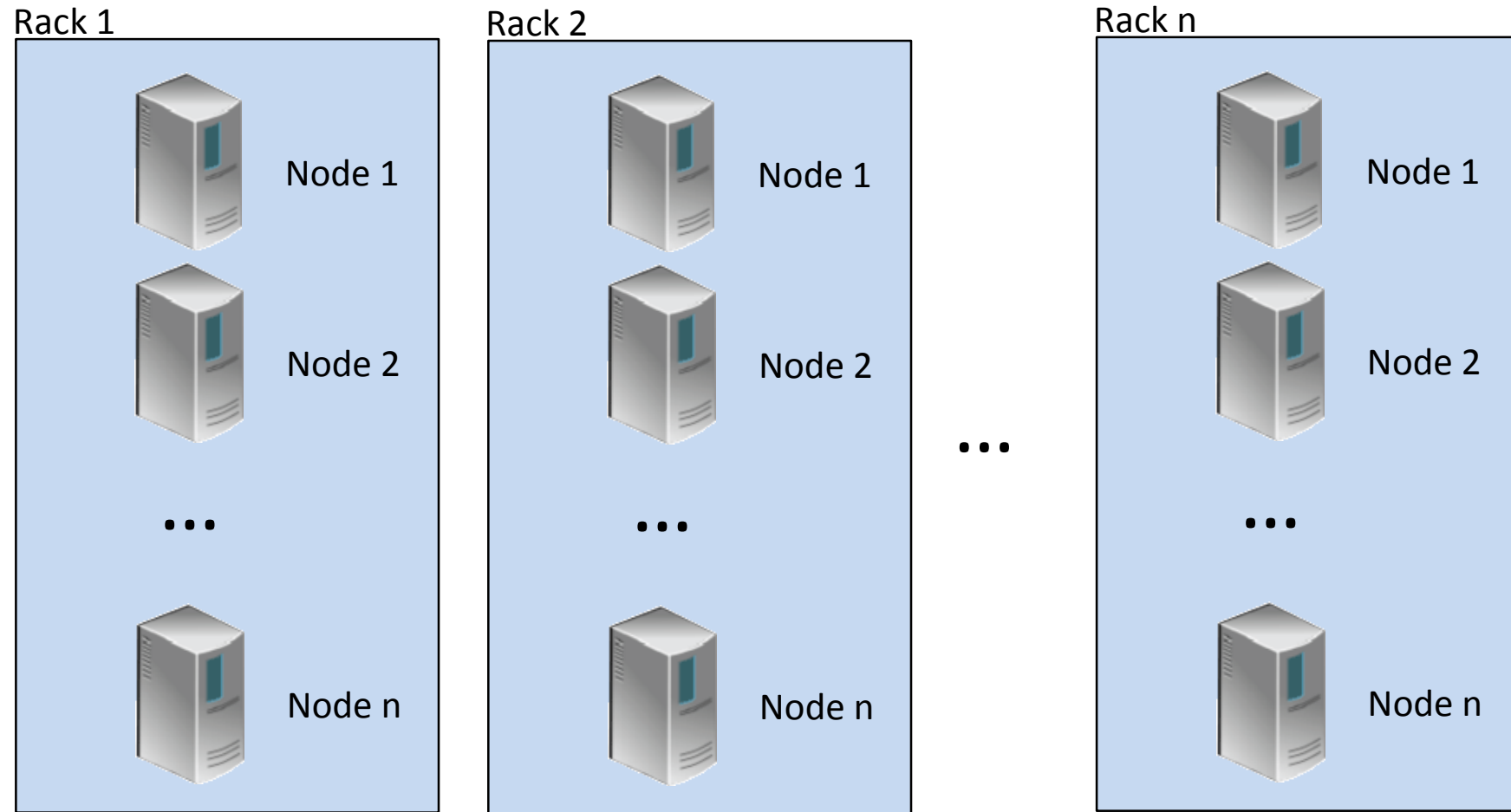
# HDFS - Terminology review



# HDFS - Terminology review

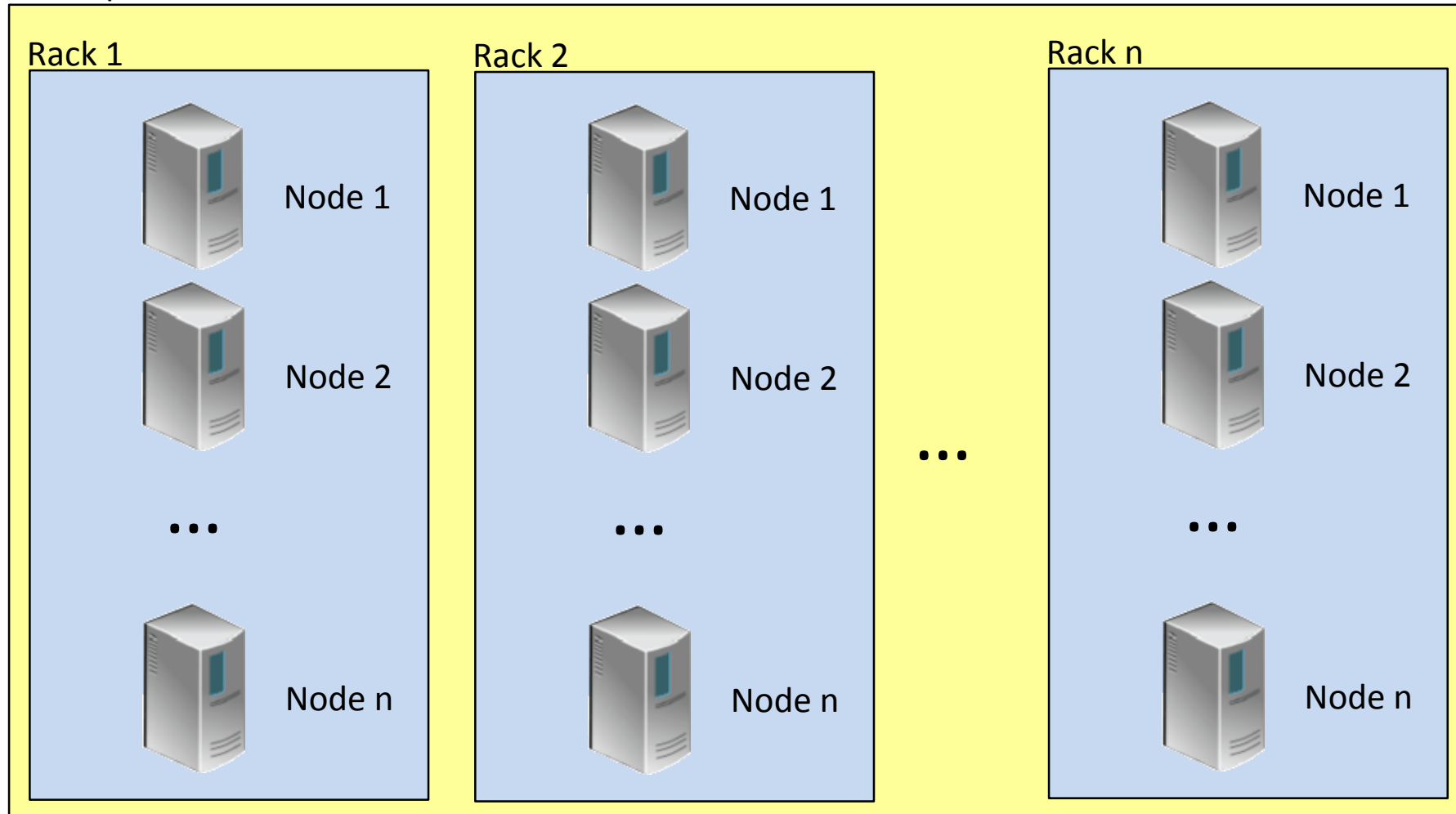


# HDFS - Terminology review

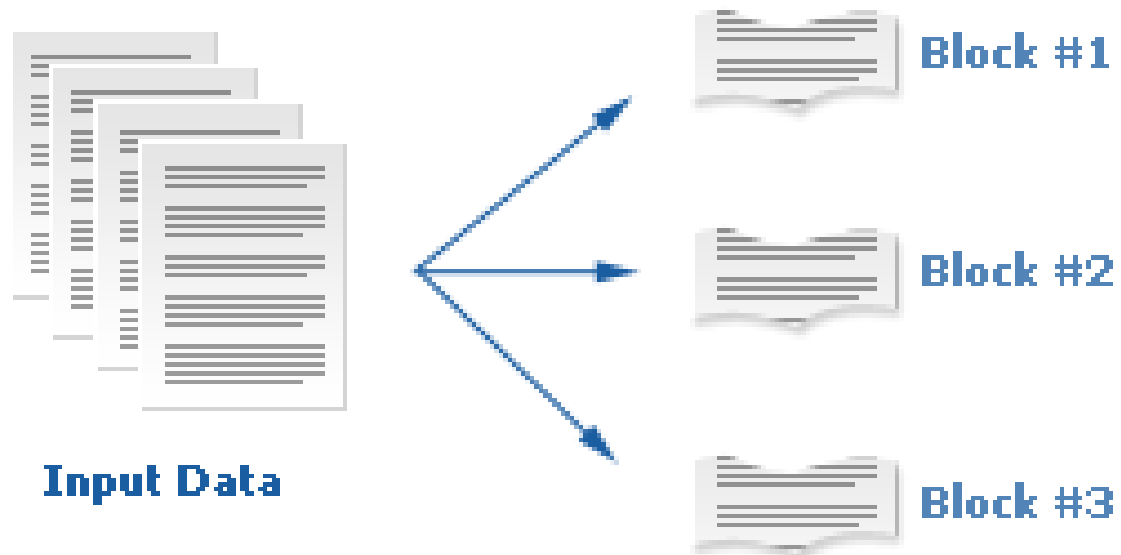




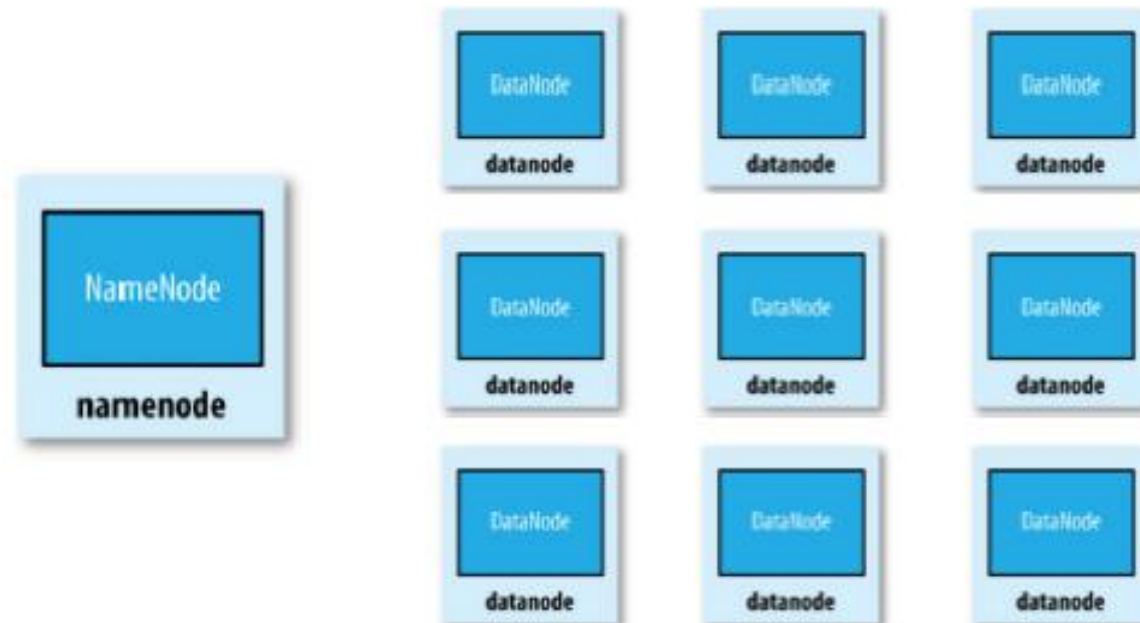
## Hadoop cluster



- Hadoop file system that runs on top of existing file system
- Designed to handle very large files with streaming data access patterns
- Uses blocks to store a file or parts of a file
- Can create, delete, copy, but NOT update



- Files are divided into chunks (blocks)
- Chunks are replicated at different compute nodes (usually 3+)
- Nodes holding copies of one chunk are located on different racks
- Chunk size and the degree of replication can be decided by the user
- A special node (the **NameNode**) stores, for each file, the positions of its chunks



- **Entire metadata is kept in RAM**
  - Ensure enough RAM in NameNode
  - If run out of RAM, NameNode will crash
- **NameNode mainly consists of:**
  - fsimage: Contains the metadata on disk (not exact copy of what is in RAM, but a checkpoint copy)
  - edit logs: Records all write operations, synchronizes with metadata in RAM after each write
- **In case of 'power failure' on NameNode**
  - Can recover using fsimage + edit logs
- **Need to format NameNode to use it:**
  - `hadoop namenode -format`

- Many per Hadoop cluster
- Manages blocks with data and serves them to clients
- Periodically reports to NameNode the list of blocks it stores
- Use inexpensive commodity hardware for this node

- 1) Some number of **Map tasks** each are given one or more chunks of data.
- 2) These Map tasks turn the chunk into a sequence of key-value pairs. The way key-value pairs are produced is determined by the code written by the user for the Map function.
- 3) The key-value pairs from each Map task are collected by a master controller and sorted and grouped by key (**Shuffle and sort**).
- 4) The keys are divided among all the **Reduce tasks**, so all key-value pairs with the same key wind up at the same Reduce task.
- 5) The Reduce tasks work on one key at a time, and combine all the values associated with that key in some way. The way values are combined is determined by the code written by the user for the Reduce function.
- 6) Output key-value pairs from each reducer are written persistently back onto the distributed file system
- 7) The output ends up in  $r$  files, where  $r$  is the number of reducers. The  $r$  files often serve as input to yet another MapReduce job

**Problem:** counting the number of occurrences for each word in a collection of documents.

**Input:** a repository of documents, each document is an element

**Map:** reads a document and emits a sequence of key-value pairs where keys are words of the documents and values are equal to 1:

$$(w_1, 1), (w_2, 1), \dots, (w_n, 1)$$

**Grouping:** groups by key and generates pairs of the form

$$(w_1, [1, 1, \dots, 1]), \dots, (w_n, [1, 1, \dots, 1])$$

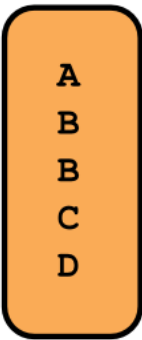
**Reduce:** adds up all the values and emits:

$$(w_1, k), \dots, (w_n, l)$$

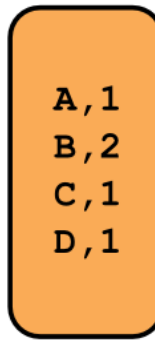
**Output:**  $(w, m)$  pairs, where  $w$  is a word that appears at least once among all the input documents and  $m$  is the total number of occurrences of  $w$  among all those documents.

# Simple data flow example

input



output

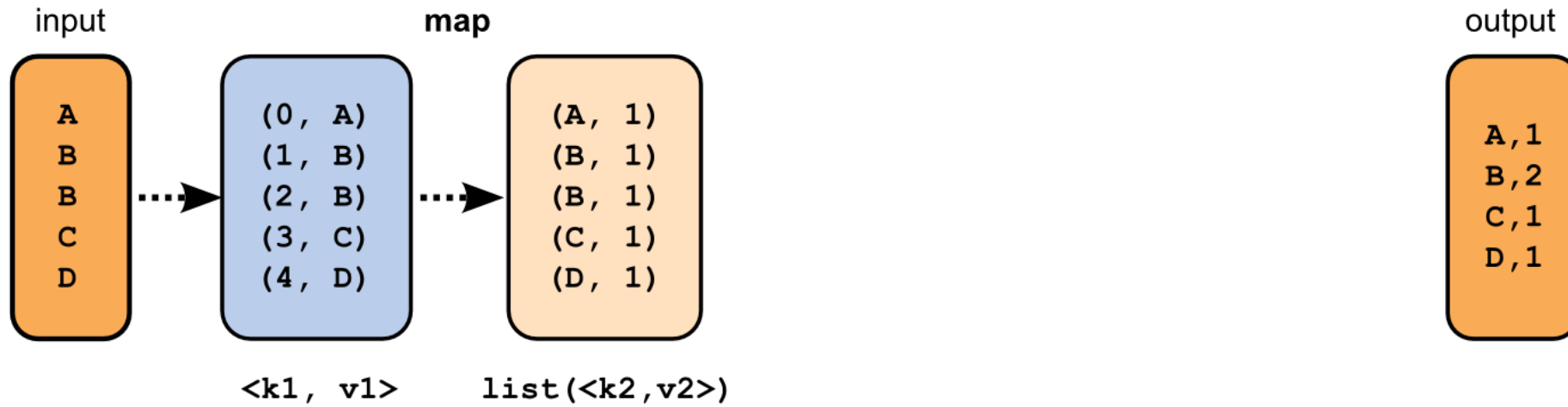




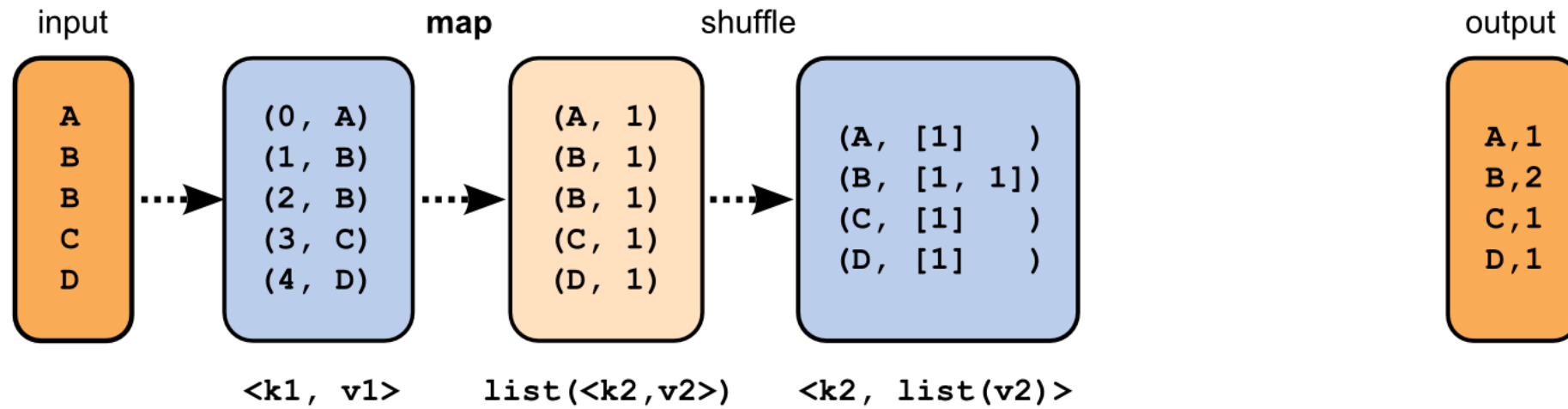
# Simple data flow example



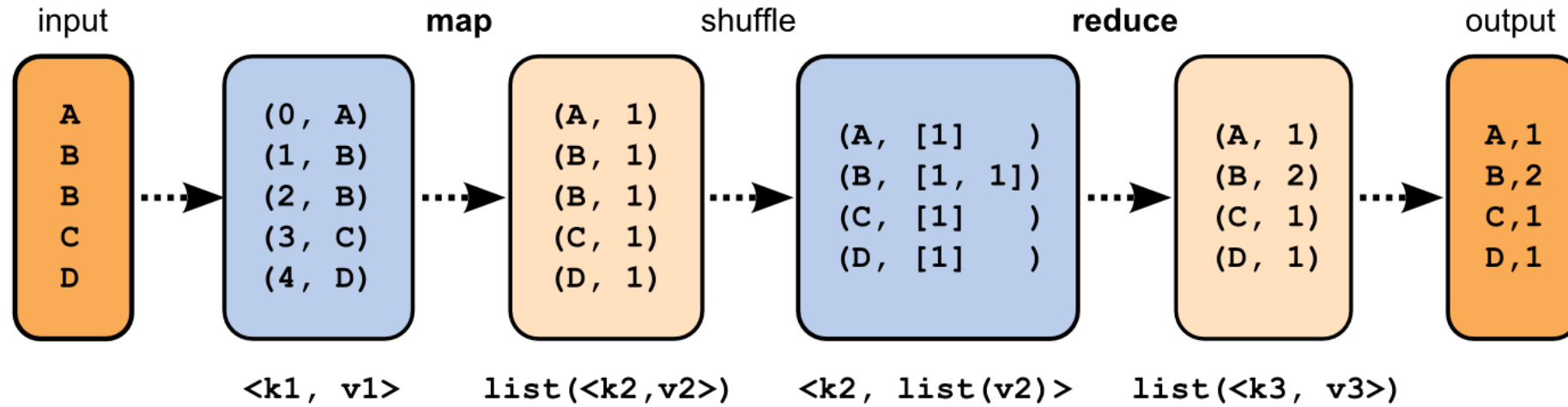
# Simple data flow example



# Simple data flow example



# Simple data flow example



# BIGINSIGHTS

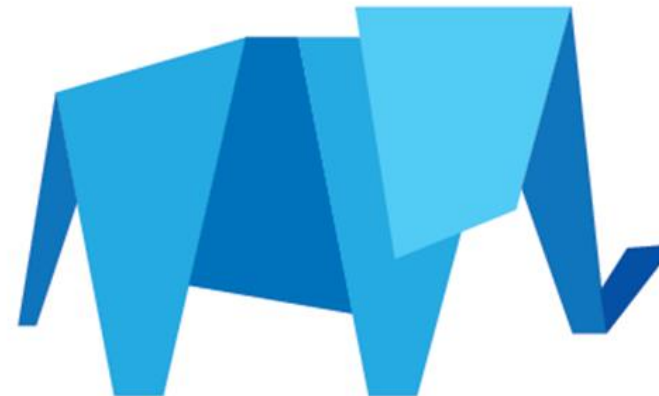
IBM InfoSphere BigInsights offre una serie completa di funzionalità di analytics avanzate che consentono alle aziende di analizzare volumi elevati di dati strutturati e non strutturati nel loro formato nativo.

Il software unisce la tecnologia Apache Hadoop open source con le innovazioni IBM, quali analytics dei testi avanzata, IBM BigSheets e Big SQL per la consultazione dei dati e una serie di funzioni amministrative, di sicurezza e prestazioni.

Il risultato è una soluzione economica e facile da utilizzare per l'analytics di big data complessi

# InfoSphere BigInsights is 100% standard, open source Hadoop

**Avoid proprietary lock-in.** BigInsights is based on Open Data Platform and includes the rich tools that Hadoop users expect. Where IBM does provide value-added features, they are carefully implemented so that customers have a choice whether to use IBM enhancements or standard Hadoop functionality.





## IBM BigInsights for Apache Hadoop

<b>SQL on Hadoop</b> Big SQL – optimized ANSI compliant SQL	<b>Data Visualization</b> BigSheets spreadsheet interface
<b>Predictive Modeling</b> Big R, Machine Learning	<b>Text Analytics</b> Advanced text processing with AQL, Text extraction web interface
<b>Storage Integration</b> GPFS - POSIX Distributed Filesystem	<b>Enterprise Manageability</b> Adaptive MapReduce, Multi-tenant scheduling
<b>Application Tooling</b> Toolkits and accelerators	<b>Search &amp; Exploration</b> Watson Explorer
<b>Real-time Analytics</b> InfoSphere Streams	<b>Data Governance and Security</b> DataClick, LDAP, Secure cluster



## IBM Open Platform with Apache Hadoop

Ambari*	Avro	Flume	Hadoop
Hive	Knox	Open JD	HDFS/MapReduce/YARN*
Slider	Snappy	Solr	Oozie
			Sqoop

# VESTAS ha ottimizzato una analisi predittiva su 2.5 Peta

**IBM**



**Vestas optimizes capital investments based on 2.5 Petabytes of information**

**Need**

- Model the weather to optimize placement of turbines, maximizing power generation and longevity

**Benefits**

- Reduce time required to identify placement of turbine from weeks to hours
- Reduces IT footprint and costs, and decreases energy consumption by 40 % – while increasing computational power
- Incorporate 2.5 PB of structured and semi-structured information flows. Data volume expected to grow to 6 PB

**Vestas**

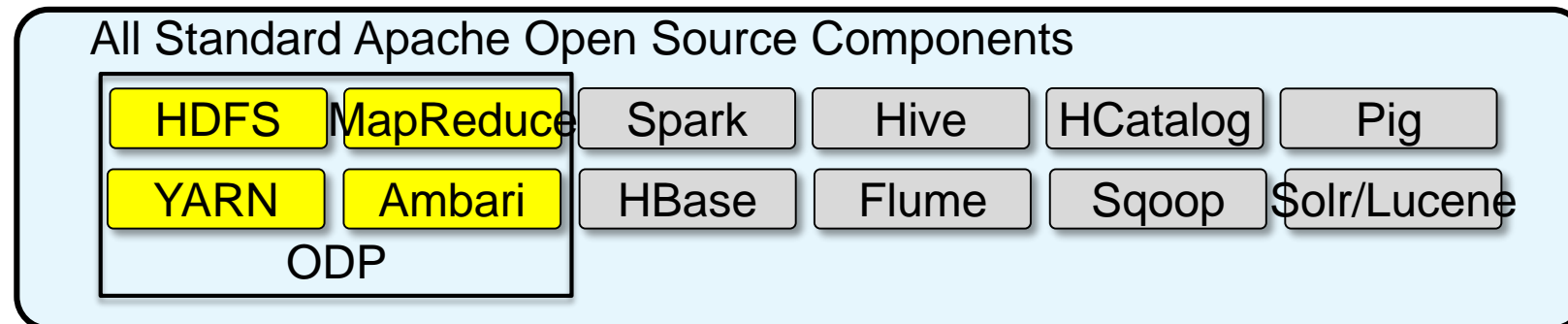
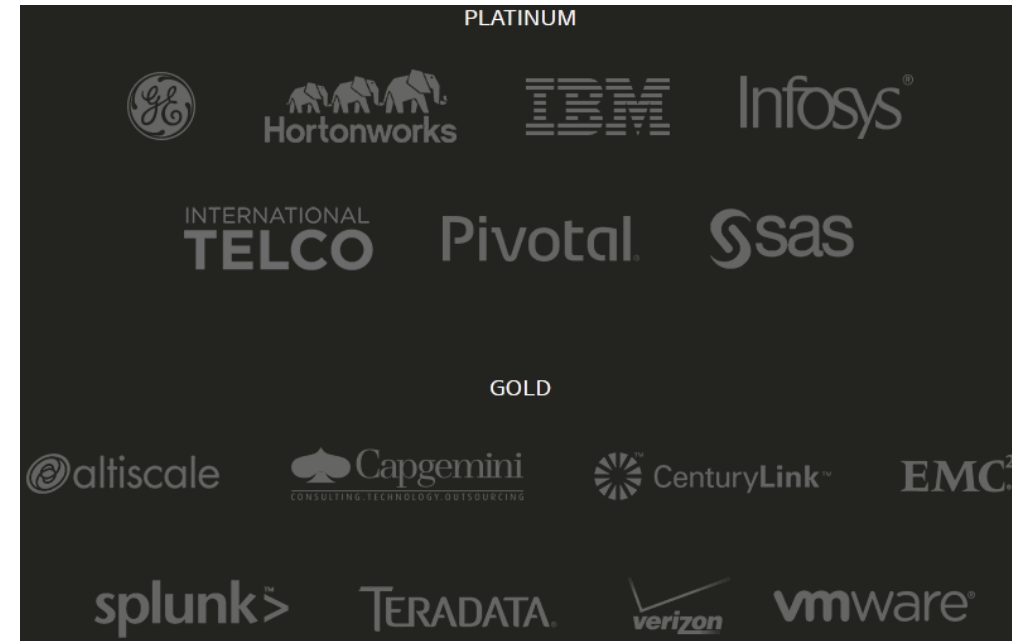


## Why is IBM involved?

- Strong history of leadership in open source & standards
- Supports our commitment to open source currency in all future releases
- Accelerates our innovation within Hadoop & surrounding applications

## Open Data Platform (ODP) vs. Apache Software Foundation (ASF)

- ODP supports the ASF mission
- ASF provides a governance model around individual projects without looking at ecosystem
- ODP aims to provide a vendor-led consistent packaging model for core Apache components as an ecosystem



## Big SQL – Lightning fast, ANSI compliant, native Hadoop formats

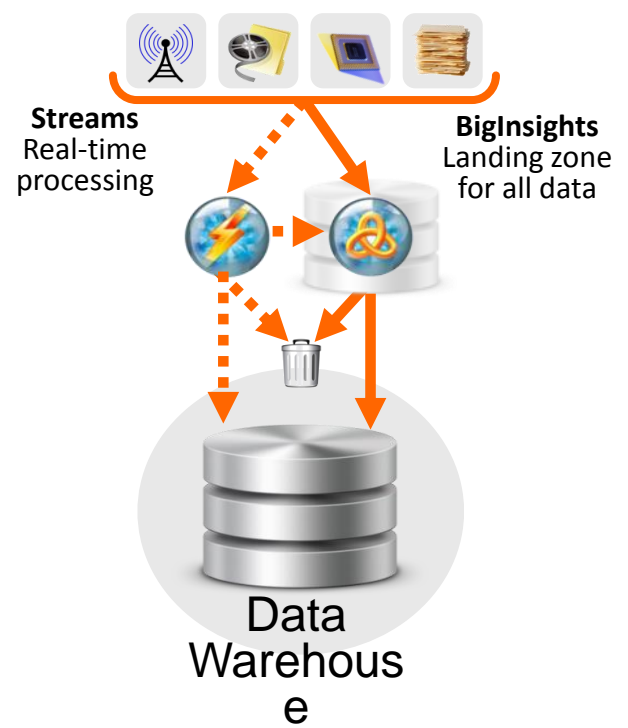
While some Hadoop vendors are building their own SQL-on-Hadoop implementations from scratch, Big SQL is ANSI compliant, lightning fast and runs on native Hadoop file formats. Brought to you by the inventors of SQL, Big SQL provides federated access so customers can query IBM and third party data sources with the same feature rich SQL.



# SQL Access for Hadoop: Why?

- Data warehouse modernization is a leading Hadoop use case

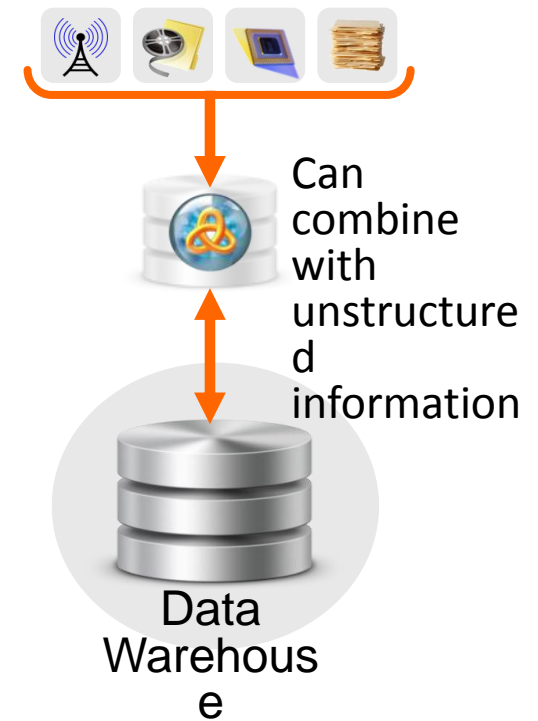
## ● Pre-processing hub



## ● Queryable archive



## ● Exploratory Analysis



- Limited availability of skills in MapReduce, Pig, etc.
- SQL opens the data to a much wider audience
  - Familiar, widely known syntax
  - Common catalog for identifying data and structure

## IBM's SQL for Hadoop

- Makes Hadoop data accessible to a wider audience
- Familiar, widely known syntax
- Leverage native Hadoop data sources

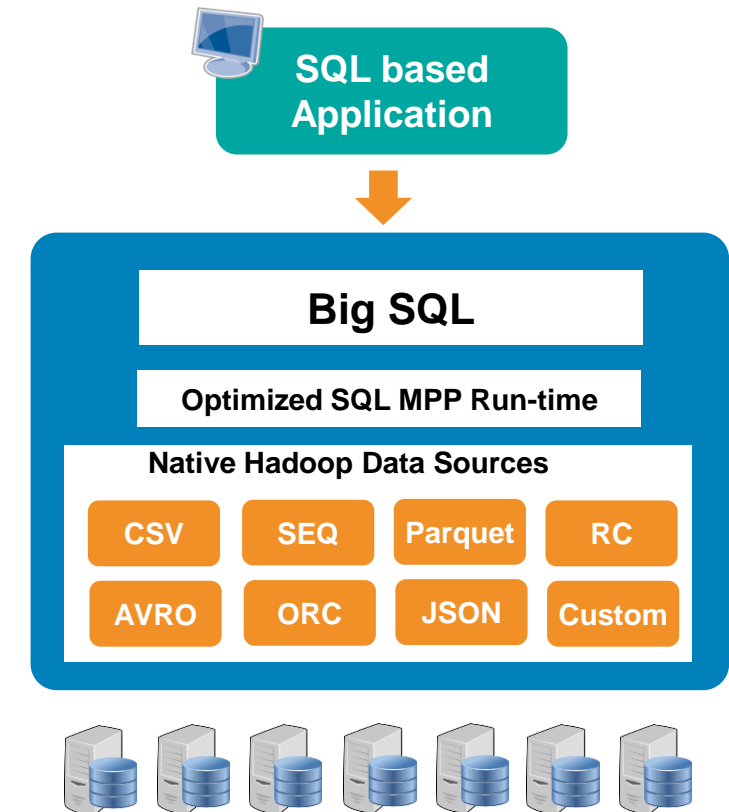
## Complements the Data Warehouse

- Exploratory analytics
- Sandbox, Data Lake

## Included in BigInsights

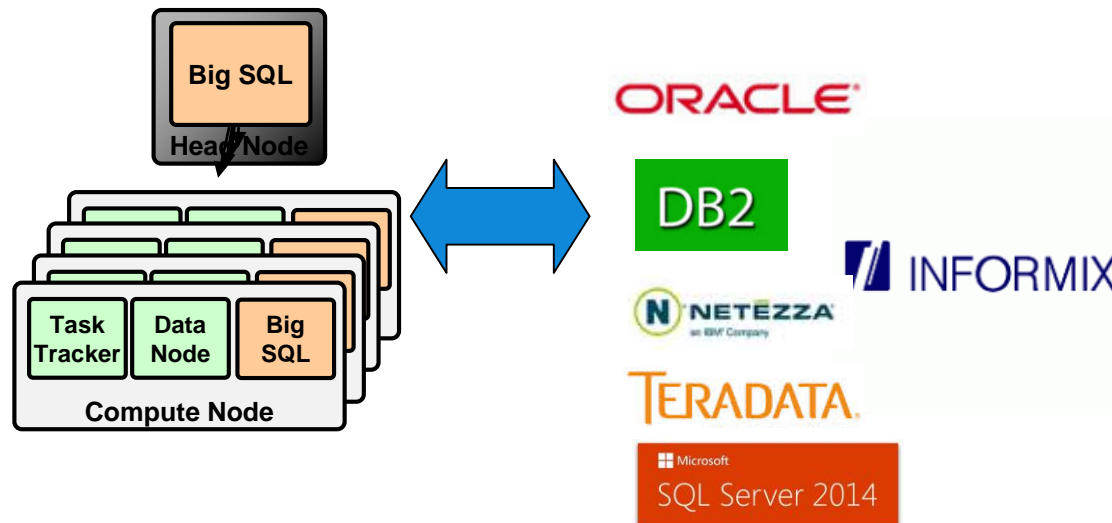
## Use familiar SQL tools

- Cognos, SPSS, Tableau, MicroStrategy




# A word about . . . query federation

- Data rarely lives in isolation
- Big SQL transparently queries heterogeneous systems
  - Join Hadoop to RDBMSs
  - Query optimizer understands capabilities of external system
    - Including available statistics*
  - As much work as possible is pushed to each system to process



- ✦ Created by IBM
- ✦ The Big Data Decision Support Benchmark (*Hadoop-DS*) is inspired by, and is highly compliant with TPC-DS
  - Fully complies with the TPC-DS schema requirement
  - Uses all 99 queries
  - Meets the multi-user requirement
  - Has been audited by a TPC-DS auditor but as a non-TPC benchmark
- ✦ Select deviations from TPC-DS due to Hadoop limitations:
  - No data maintenance operations, referential integrity enforcement, or ACID property validation as these are not feasible with HDFS
  - Additional statistics used
  - Metric adjustments
  - No price/performance measures included
  - Not an official TPC benchmark result**

Competitive environments require significant effort

IBM			cloudera IMPALA			 HIVE		
Query 01	Query 34	Query 67	Query 01	Query 34	Query 67	Query 01	Query 34	Query 67
Query 02	Query 35	Query 68	Query 02	Query 35	Query 68	Query 02	Query 35	Query 68
Query 03	Query 36	Query 69	Query 03	Query 36	Query 69	Query 03	Query 36	Query 69
Query 04	Query 37	Query 70	Query 04	Query 37	Query 70	Query 04	Query 37	Query 70
Query 05	Query 38	Query 71	Query 05	Query 38	Query 71	Query 05	Query 38	Query 71
Query 06	Query 39	Query 72	Query 06	Query 39	Query 72	Query 06	Query 39	Query 72
Query 07	Query 40	Query 73	Query 07	Query 40	Query 73	Query 07	Query 40	Query 73
Query 08	Query 41	Query 74	Query 08	Query 41	Query 74	Query 08	Query 41	Query 74
Query 09	Query 42	Query 75	Query 09	Query 42	Query 75	Query 09	Query 42	Query 75
Query 10	Query 43	Query 76	Query 10	Query 43	Query 76	Query 10	Query 43	Query 76
Query 11	Query 44	Query 77	Query 11	Query 44	Query 77	Query 11	Query 44	Query 77
Query 12	Query 45	Query 78	Query 12	Query 45	Query 78	Query 12	Query 45	Query 78
Query 13	Query 46	Query 79	Query 13	Query 46	Query 79	Query 13	Query 46	Query 79
Query 14	Query 47	Query 80	Query 14	Query 47	Query 80	Query 14	Query 47	Query 80
Query 15	Query 48	Query 81	Query 15	Query 48	Query 81	Query 15	Query 48	Query 81
Query 16	Query 49	Query 82	Query 16	Query 49	Query 82	Query 16	Query 49	Query 82
Query 17	Query 50	Query 83	Query 17	Query 50	Query 83	Query 17	Query 50	Query 83
Query 18	Query 51	Query 84	Query 18	Query 51	Query 84	Query 18	Query 51	Query 84
Query 19	Query 52	Query 85	Query 19	Query 52	Query 85	Query 19	Query 52	Query 85
Query 20	Query 53	Query 86	Query 20	Query 53	Query 86	Query 20	Query 53	Query 86
Query 21	Query 54	Query 87	Query 21	Query 54	Query 87	Query 21	Query 54	Query 87
Query 22	Query 55	Query 88	Query 22	Query 55	Query 88	Query 22	Query 55	Query 88
Query 23	Query 56	Query 89	Query 23	Query 56	Query 89	Query 23	Query 56	Query 89
Query 24	Query 57	Query 90	Query 24	Query 57	Query 90	Query 24	Query 57	Query 90
Query 25	Query 58	Query 91	Query 25	Query 58	Query 91	Query 25	Query 58	Query 91
Query 26	Query 59	Query 92	Query 26	Query 59	Query 92	Query 26	Query 59	Query 92
Query 27	Query 60	Query 93	Query 27	Query 60	Query 93	Query 27	Query 60	Query 93
Query 28	Query 61	Query 94	Query 28	Query 61	Query 94	Query 28	Query 61	Query 94
Query 29	Query 62	Query 95	Query 29	Query 62	Query 95	Query 29	Query 62	Query 95
Query 30	Query 63	Query 96	Query 30	Query 63	Query 96	Query 30	Query 63	Query 96
Query 31	Query 64	Query 97	Query 31	Query 64	Query 97	Query 31	Query 64	Query 97
Query 32	Query 65	Query 98	Query 32	Query 65	Query 98	Query 32	Query 65	Query 98
Query 33	Query 66	Query 99	Query 33	Query 66	Query 99	Query 33	Query 66	Query 99

## Key points

- With competing solutions, many queries needed to be re-written, some significantly
- Owing to various restrictions, some queries could not be re-written or failed at run-time
- Re-writing queries in a benchmark scenario where results are known is one thing – doing this against real databases in production is another

	Works without modification
	Minor modification
	Extensive modification
	Not working

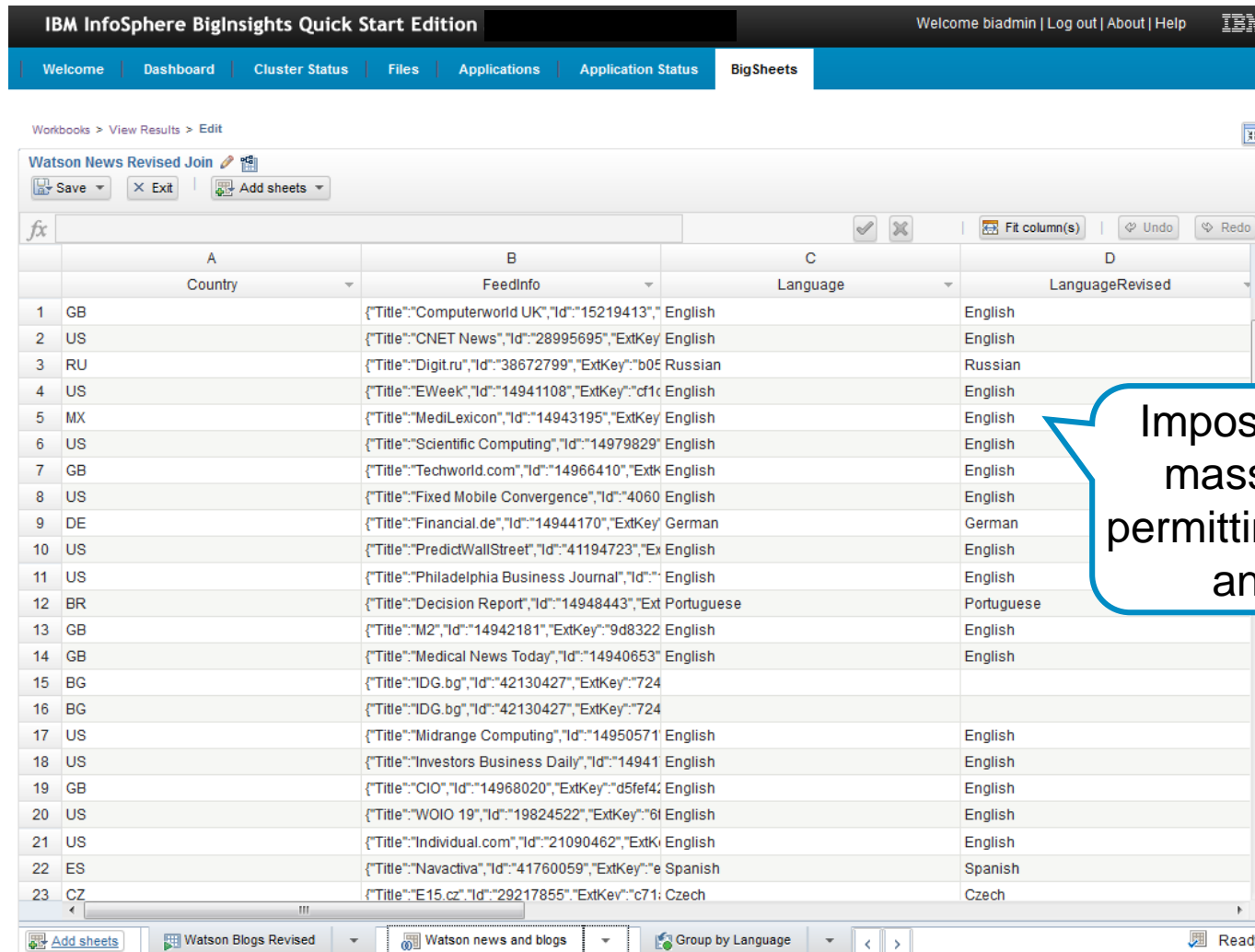
## BigSheets – Spreadsheet-like data access for business users

Shockingly, there are those who prefer not to code in Java or write Pig scripts to get data from Hadoop. BigSheets provides an easy-to-use spreadsheet interface allowing business users to extract, manipulate and visualize data from a variety of Hadoop and non-Hadoop data sources.





## Spreadsheet style analysis tool for business users

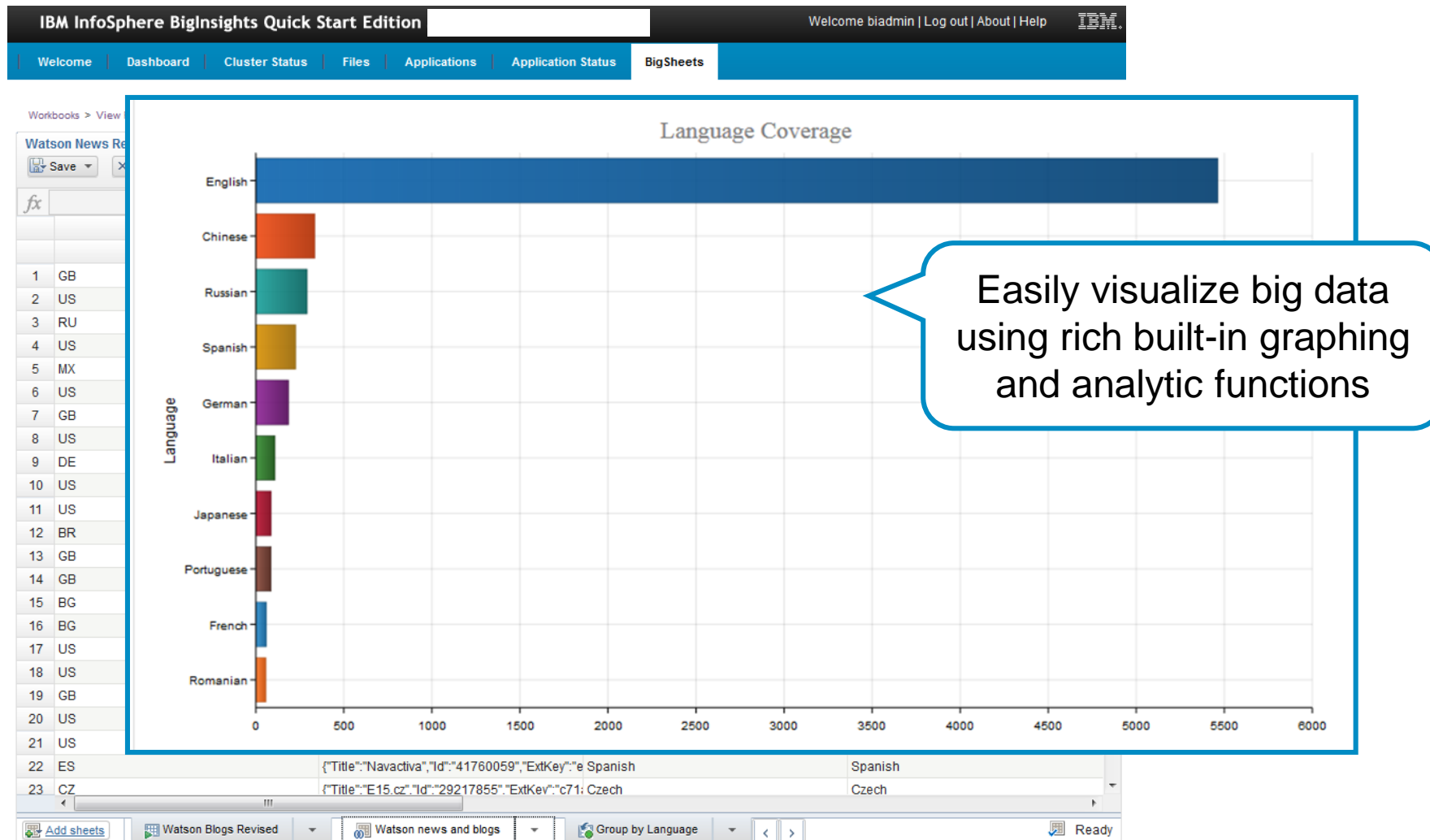


The screenshot shows the IBM InfoSphere BigInsights Quick Start Edition interface. The top navigation bar includes 'Welcome', 'Dashboard', 'Cluster Status', 'Files', 'Applications', 'Application Status', and 'BigSheets'. The main area displays a spreadsheet titled 'Watson News Revised Join'. The spreadsheet has columns for 'Country', 'FeedInfo', 'Language', and 'LanguageRevised'. The data rows show various news items with their respective country codes and language details.

	A	B	C	D
	Country	FeedInfo	Language	LanguageRevised
1	GB	{\"Title\":\"Computerworld UK\",\"Id\":\"15219413\",	English	English
2	US	{\"Title\":\"CNET News\",\"Id\":\"28995695\",	English	English
3	RU	{\"Title\":\"Digit.ru\",\"Id\":\"38672799\",	Russian	Russian
4	US	{\"Title\":\"EWeek\",\"Id\":\"14941108\",	English	English
5	MX	{\"Title\":\"MediLexicon\",\"Id\":\"14943195\",	English	English
6	US	{\"Title\":\"Scientific Computing\",\"Id\":\"14979829\",	English	English
7	GB	{\"Title\":\"Techworld.com\",\"Id\":\"14966410\",	English	English
8	US	{\"Title\":\"Fixed Mobile Convergence\",	English	English
9	DE	{\"Title\":\"Financial.de\",	German	German
10	US	{\"Title\":\"PredictWallStreet\",	English	English
11	US	{\"Title\":\"Philadelphia Business Journal\",	English	English
12	BR	{\"Title\":\"Decision Report\",	Portuguese	Portuguese
13	GB	{\"Title\":\"M2\",	English	English
14	GB	{\"Title\":\"Medical News Today\",	English	English
15	BG	{\"Title\":\"IDG.bg\",		
16	BG	{\"Title\":\"IDG.bg\",		
17	US	{\"Title\":\"Midrange Computing\",	English	English
18	US	{\"Title\":\"Investors Business Daily\",	English	English
19	GB	{\"Title\":\"CIO\",	English	English
20	US	{\"Title\":\"WOIO 19\",	English	English
21	US	{\"Title\":\"Individual.com\",	English	English
22	ES	{\"Title\":\"Navactiva\",	Spanish	Spanish
23	CZ	{\"Title\":\"E15.cz\",	Czech	Czech

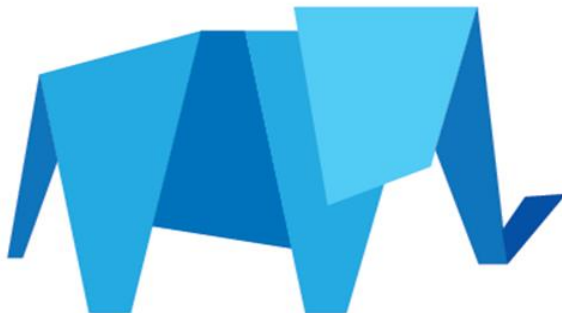
Impose structure on massive datasets permitting manipulation and analysis

## Spreadsheet style analysis tool for business users



## Big Text – Simplify text analytics and natural language processing

Building applications able to extract meaning from text is hard. Rather than build your own solution, stand on the shoulders of the inventors of Watson, the reigning Jeopardy champ. Build your own advanced analytic applications, parse jargon-laden text in multiple languages, and do what your competitors cannot.



- Distills structured info from unstructured text
  - Sentiment analysis
  - Consumer behavior
  - Illegal or suspicious activities
  - ...
- Parses text and detects meaning with annotators
- Understands the context in which the text is analyzed
- Features pre-built extractors for names, addresses, phone numbers, etc.

## Unstructured text (document, email, etc)

Football **World Cup 2010**, one team distinguished themselves well, losing to the eventual champions 1-0 in the Final. Early in the second half, **Netherlands' striker, Arjen Robben**, had a breakaway, but the **keeper for Spain, Iker Casillas** made the save. **Winger Andres Iniesta** scored for **Spain** for the win.



## Classification and Insight

World Cup 2010 Highlights

Name	Position	Country
Arjen Robben	Striker	Netherlands
Iker Casillas	Keeper	Spain
Andres Iniesta	Winger	Spain

## Big R – Deep R language integration on Hadoop

It seems that everyone has a solution for R on Hadoop. What's unique about Big R is that it implements standard R-language functions for parallel processing, supports routines available from the Comprehensive R Archive Network (CRAN), it uses standard R developer tools, and supports advanced machine analytic functions in R.

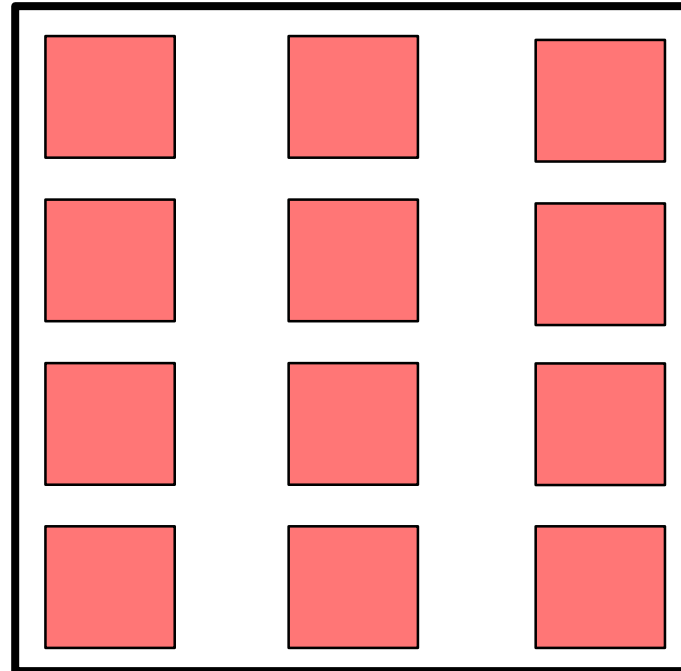


# Big R Data Structures: Proxy to entire dataset

Appears and acts like all of the data is on your laptop

```
data <- bigr.frame(...)
```

You



## In-Hadoop Analytics – Deploy the analytics to the data

The whole point of Hadoop is to move the compute to the data, but IBM takes this to the next level with native analytic functions accessible from R, SQL and other languages. Provide seamless access to advanced in Hadoop statistical and machine learning functions avoiding the cost and complexity of software development and additional tooling.



# Link utili



Try it for free!

Non-production, no-limit version of IBM  
InfoSphere BigInsights

<http://IBM.co/QuickStart>



**Hadoop, welcome  
to the enterprise.**

Download free now at [IBM.co/QuickStart](http://IBM.co/QuickStart)

**New!** InfoSphere BigInsights Quick Start Edition

The advertisement features a central image of an open cardboard box with a bright yellow light emanating from it, set against a background of radiating yellow lines. The text is in a bold, black, sans-serif font.

## A few essential links



**IBM Analytics**  
<http://www.ibm.com/analytics/>



**IBM Big Data**  
<http://www-01.ibm.com/software/data/bigdata/>



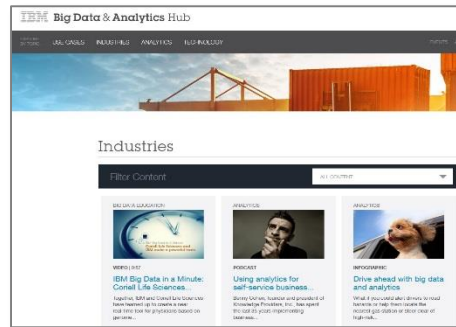
**IBM Watson**  
<http://www.ibm.com/smarterplanet/us/en/ibmwatson/>

# For more info

## Other relevant links



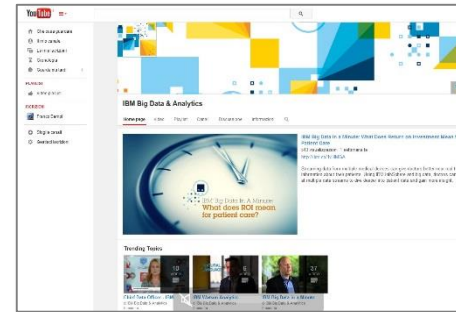
Big Data university  
<http://bigdatauniversity.com/>



IBM Big Data Hub  
<http://www.ibmbigdatahub.com>



IBM Analytics Zone  
<https://www.analyticszone.com>



IBM Big Data Youtube Channel

<https://www.youtube.com/user/ibmbigdata>



IBM Big Data case studies

<http://www.ibm.com/big-data/us/en/big-data-and-analytics/case-studies.html#filter2=customerinsight>

# Due buoni libri ...free...sui Big Data



## Big Data Beyond the Hype

A Guide to Conversations for Today's Data Center

- Expand what you know about Big Data, putting you on track to go beyond the hype
- Stay on top of the latest news, including the cloud, performance optimization, streaming analytics, containerization, including Big SQL, NoSQL, integration, governance, and more
- See how IBM Watson and other services make sense of the IBM FI series and essential dimensions in the Big Data story
- Learn about the Big Data Cloud Made easy brings a new approach to managing data faster to deploy, faster to insights, and with less risk
- Gain confidence in your Big Data projects and learn about the importance of governance in a Big Data world

Paul Zikopoulos   Dirk deRooz  
Chris Sienko   Rick Buglio   Marc Andrews

[https://www.ibm.com/developerworks/community/blogs/SusanVisser/entry/big\\_data\\_beyond\\_the\\_hype\\_a\\_guide\\_to\\_conversations\\_for\\_today\\_s\\_data\\_center?lang=en](https://www.ibm.com/developerworks/community/blogs/SusanVisser/entry/big_data_beyond_the_hype_a_guide_to_conversations_for_today_s_data_center?lang=en)



## The Power of Now

Real-Time Analytics  
and IBM InfoSphere Streams

- Discover real-time analytics and multiple technologies
- Analyze big data and structure that can process data, both at rest and in motion
- Learn how new data distributed analysis can provide insights
- Discover how to build a real-time analytics platform with a new approach
- Discover the ways to implement a real-time analytics solution that works for your business
- Explore the characteristics that make a good real-time analytics platform
- Learn how to apply InfoSphere Streams technology to your business problems to get solutions faster

Jacques Roy

[https://www.ibm.com/developerworks/community/blogs/SusanVisser/entry/the\\_power\\_of\\_now\\_real\\_time\\_analytics\\_and\\_ibm\\_infoSphere\\_streams?lang=en](https://www.ibm.com/developerworks/community/blogs/SusanVisser/entry/the_power_of_now_real_time_analytics_and_ibm_infoSphere_streams?lang=en)

# Uno degli ultimi COURSES...su SPARK



**Courses**  
Pick from our vast selection of courses

Search courses

Home / Courses / All Courses / Spark Fundamentals

**All Courses** >  
Featured >  
Just Released! >  
Big Data >  
Cloud Computing >  
Programming >  
Database Technologies >

**Spark Fundamentals**  
with Henry Quach

**Audience:** Data scientists, engineers, or anyone who is interested in learning about Spark.

**Time to complete:** 03:00

**Available in:** English

Apache Spark is an open source processing engine built around speed, ease of use, and analytics. If you have large amounts of data that requires low latency processing that a typical Map Reduce program cannot provide, Spark is the alternative. Spark performs at speeds up to 100 times faster than Map Reduce for iterative algorithms or interactive data mining. Spark provides in-memory cluster computing for lightning fast speed and supports Java, Scala, and Python APIs for ease of development.

Spark combines SQL, streaming and complex analytics together seamlessly in the same application to handle a wide range of data processing. Spark runs on top of Hadoop, Mesos, standalone, or is available as a service to access diverse data sources such as...

Mi fermo qui  
**GRAZIE**  
a tutti voi