Basi Di Dati e di conoscenza

Basi di Dati e Big Data

Contenuti della lezione

- Memorizzare e gestire I Big Data
 - Hadoop & Map-reduce
 - Cloud computing
 - NoSQL DBMS



Goal:

to make effective strategic decisions exploiting the availability of big data

Paola Vocca

Basi di Dati e di Conoscenza - Basi di Dati e Big Data



- selection
- filtering
- Metadata generation
- Managing provenance

Basi di Dati e di Conoscenza - Basi di Dati e Big Data



- transformation
- normalization
- cleaning
- aggregation
- error handling



- standardization
- Conflict management
- reconciliation
- Mapping definition

Basi di Dati e di Conoscenza - Basi di Dati e Big Data



- exploration
- data mining
- Machine learning
- visualization



- Knowledge of the domain
- knowledge of the provenance
- Identification of patterns of interest
- Flexibility of the process

Basi di Dati e di Conoscenza - Basi di Dati e Big Data



- managerial skills
- Continuous improvement of the process

Challenges

- Performance, performance, performance!
- Scalability
- Heterogeneity
- Effectiveness
- Flexibility
- Privacy
- Property
- Human collaboration

Data Volumes

- The amount of data increases every day
- Some numbers (~ 2012):
 - Data processed by Google every day: 100+ PB
 - Data processed by Facebook every day: 10+ PB
- To analyze them, systems that scale with respect to the data volume are needed

From https://docs.microsoft.com/en-us/azure/architecture/data-guide/big-data/

Data volumes: Google Example

- Analyze 10 billion web pages
- Average size of a webpage: **20KB**
- Size of the collection: 10 billion x 20KBs = **200TB**
- HDD hard disk read bandwidth: **550MB/sec**
- Time needed **to read all web pages** (**without analyzing them**): 2 million seconds = more than **15 days**
- A single node architecture is not adequate

Data volumes: Google Example

- Analyze 10 billion web pages
- Average size of a webpage: **20KB**
- Size of the collection: 10 billion x 20KBs = **200TB**
- **SSD** hard disk read bandwidth: **550MB/sec**
- Time needed **to read all web pages** (**without analyzing them**): 2 million seconds = more than 4 **days**
- A single node architecture is not adequate

Failures

- Failures are part of everyday life, especially in data center
 - A single server stays up for 3 years (~1000 days)
 - 10 servers ~ 1 failure every 100 days (~3 months)
 - 100 servers ~ 1 failure every 10 days
 - 1000 servers ~ 1 failure/day
- Sources of failures
 - Hardware/Software
 - Electrical, Cooling, ...
 - Unavailability of a resource due to overload

Failures

- LALN data [DSN 2006]
 - Data for 5000 machines, for 9 years
 - Hardware failures: 60%, Software: 20%, Network 5%
- DRAM error analysis [Sigmetrics 2009]
 - Data for 2.5 years
 - 8% of DIMMs affected by errors
- Disk drive failure analysis [FAST 2007]
 - Utilization and temperature major causes of failures

Failures

- Failure types
 - Permanent
 - E.g., Broken motherboard
- Transient
 - E.g., Unavailability of a resource due to overload

Network bandwidth

- Network becomes the bottleneck if big amounts of data need to be exchanged between nodes/servers
 - Network bandwidth (in a data center): 10Gbps
 - Moving 10 TB from one server to another takes more than 2 hours
 - Data should be moved across nodes only when it is indispensable
 - Usually, codes/programs are small (few MBs)
 - Move code (programs) and computation to data

Network bandwidth

- Network becomes the bottleneck if big amounts of data need to be exchanged between nodes/servers
 - Network bandwidth (in a data center): 10Gbps
 - Moving 10 TB from one server to another takes more than 2 hours
 Data should be moved across nodes only when it is indispensable
 - Usually, codes/programs are small (few MBs)
 - **Move code** (programs) and **computation** to data



Single-node architecture



Server (Single node)

Paola Vocca

Basi di Dati e di Conoscenza - Basi di Dati e Big Data

Single-node architecture



Server (Single node)

Machine Learning, Statistics

• Small data

• Data can be completely loaded in main memory

Basi di Dati e di Conoscenza - Basi di Dati e Big Data

Single-node architecture



Server (Single node)

"Classical" data mining

• Large data

- Data can not be completely loaded in main memory
- Load in main memory one chunk of data at a time
 - Process it and store some statistics
 - Combine statistics to compute the final result

Cluster Architecture

- Cluster of servers (data center)
 - Computation is distributed across servers
 - Data are stored/distributed across servers
- Standard architecture in the Big data context (~ 2012)
 - Cluster of commodity Linux nodes/servers
 - 32 GB of main memory per node
- Gigabit Ethernet interconnection

Commodity Cluster Architecture



Each rack contains 16-64 nodes

Paola Vo

Scalability

- Current systems must scale to address
 - The increasing amount of data to analyze
 - The increasing number of users to serve
 - The increasing complexity of the problems
- Two approaches are usually used to address scalability issues
 - Vertical scalability (scale up)
 - Horizontal scalability (scale out)

Scale up vs. Scale out

- Vertical scalability (scale up)
 - Add more power/resources (main memory, CPUs) to a single node (high-performing server)
 - Cost of super-computers is not linear with respect to their resources

• Horizontal scalability (scale out)

- Add more nodes (commodity servers) to a system
 - The cost scales approximately linearly with respect to the number of added nodes
 - But data center efficiency is a difficult problem to solve

Scale up vs. Scale out

- For data-intensive workloads, a large number of commodity servers is preferred over a small number of high-performing servers
 - At the same cost, we can deploy a system that processes data more efficiently and is more fault-tolerant
- Horizontal scalability (scale out) is preferred for big data applications
 - But distributed computing is hard
 - New systems hiding the complexity of the distributed part of the problem to developers are needed

Cluster computing challenges

• Distributed programming is hard

- Problem decomposition and parallelization
- Task synchronization

• Task scheduling of distributed applications is critical

- Assign tasks to nodes by trying to
 - Speed up the execution of the application
 - Exploit (almost) all the available resources
 - Reduce the impact of node failures

Cluster computing challenges

• Distributed data storage

- How do we store data persistently on disk and keep it available if nodes can fail?
 - **Redundancy** is the solution, but it increases the complexity of the system

Network bottleneck

- Reduce the amount of data send through the network
 - Move computation and code to data

Cluster computing challenges

- Distributed computing is not a new topic
 - HPC (High-performance computing) ~1960
 - Grid computing ~1990
 - Distributed databases ~1990
- Hence, many solutions to the mentioned challenges are already available
- But we are now facing big data driven-problems

The former solutions are not adequate to address big data volumes

Typical Big Data Problem

• Typical Big Data Problem

- Iterate over a large number of records/objects
- Extract something of interest from each record/object
- Aggregate intermediate results
- Generate final output

• The challenges:

- Parallelization
- Distributed storage of large data sets (Terabytes, Petabytes)
- Node Failure management
- Network bottleneck
- Diverse input format (data diversity & heterogeneity)

Parallel Data Processing

Processing Big Data Frameworks to process Big Data

- Relational DBMS: old fashion
- NoSQL data stores and NewSQL databases
- **Batch processing**: store and process data sets at massive scale (especially Volume+Variety)
 - MapReduce and Hadoop
 - Spark
- **Data stream processing**: **process** fast data (in real-time) as data is being generated, **without storing**







Parallel programming: background

- Parallel programming
 - Simultaneous use of multiple computing resources (e.g., processors) to solve a problem
 - How? Break processing into **parts** that can be **executed concurrently** on multiple computing resources



Paola Vocca

Basi di Dati e di Conoscenza - Basi di Dati e Big Data

Parallel programming: background

- Simplest environment for parallel programming
 - Master/worker architecture
 - Master
 - Gets data and splits it into chunks according to the number of workers
 - Sends each worker equal Number of chunks
 - Receives results from each worker
 - Workers:
 - Receive some chunks of data from master
 - Perform processing
 - Send back results to master



Parallel programming: background

- There are several styles of parallel programming
- Single Program, Multiple Data (SPMD) is the most commonly used
 - **Single Program**: all computing resources execute the same program simultaneously
 - Multiple Data: all computing resources may use different data

Example: Pi estimation

- Estimation algorithm for calculating π
 - Relies on Monte Carlo method
- Let's first consider the sequential version of the algorithm
- Then, how to realize a parallel and faster version
• By definition, π is the area of a circle with radius equal to 1



• How to estimate π ?

- 1. Pick a large number of points randomly inside the circumscribed unit square
 - A certain number of these points will end up inside the area described by the circle, while the remaining number of these points will lie outside of it (but inside the square)
- 2. Count the fraction of points that end up inside the circle out of a total population of points randomly thrown at the circumscribed square

• In formulas: $\frac{\pi}{4} \approx \frac{N_{inner}}{N_{inner}}$

$$\pi pprox 4 rac{N_{inner}}{N_{total}}$$



• The more points generated, the greater the accuracy of the estimation



See animation at https://academo.org/demos/estimating-pi-monte-carlo/

Basi di Dati e di Conoscenza - Basi di Dati e Big Data



Total Number of points: 95770 Points within circle: 75212 Pi estimation: 3.14136

The more points generated, the greater the accuracy of the estimation

- How to get an accurate and faster estimation of π ?
- From sequential to **parallel** computation
- Use master/worker approach
 - Each worker runs the algorithm to generate a set of random points, categorize them, and count how many end up inside the circle
 - The master collects from the workers the total number of generated points and total number of points in the circle. It calculates the ratio of the two numbers and multiplies it by 4 to get a more accurate estimation of π

Key idea behind MapReduce and Spark: Divide and conquer

- Feasible approach to tackle large-data problems
 - Partition a large problem into smaller sub-problems
 - Solve independent sub-problems in parallel
 - Combine intermediate results from each individual worker



Divide and conquer: how?

- **Decompose** the original problem in smaller, parallel tasks
- Schedule tasks on workers distributed in a cluster, keeping into account:
 - Data locality
 - Resource availability
- Ensure workers get the data they need
- Coordinate synchronization among workers
- **Share** partial results
- Handle **failures**

Key idea behind MapReduce and

Spark: scale out, not up!

- For data-intensive workloads, a large number of commodity servers is preferred over a small number of high-end servers
 - Cost of super-computers is not linear
 - Data center efficiency
- **Processing** data is **quick**, I/O is slow
- **Shared nothing** is preferable over sharing
 - **Shared nothing**: each node is completely independent of other nodes in the system, no shared memory or storage
 - Scalability and fault tolerance
 - **Sharing**: nodes share a common/global state that must be managed
 - X Requires synchronization, deadlocks can occur, shared resources can become bottlenecks (e.g., bandwidth to access stored data)

HADOOP – MAP REDUCE



Technology: Hadoop & MapReduce

- What is Hadoop?
 - An open-source software framework (Apache project)
 - Originally developed by Yahoo!
- **Goal**: storage and processing of data-sets at massive scale
- Infrastructure: clusters of commodity hardware
- Core:
 - HDFS, a distributed file system
 - MapReduce, a programming model for large scale data processing
- Includes a number of related projects
 - Apache Pig, Apache Hive, Apache HBase, etc..
- Used in production by Google, Facebook, Yahoo! and many others



The core of Hadoop

• HDFS

- A distributed file systems
- Servers can fail and not abort the computation process
- Data is replicated with redundancy across the cluster
- MapReduce
 - Programming paradigm for expressing distributed computations over multiple servers
 - The powerhouse behind most of today's big data processing
 - Also used in other MPP environments and NoSQL databases (e.g., Vertica and MongoDB)
- Improving programmability: Pig and Hive
- Improving data access: HBase, Sqoop, and Flume

Apache Hadoop

- Scalable fault-tolerant distributed system for Big Data
 - Distributed Data Storage
 - Distributed Data Processing
 - Borrowed concepts/ideas from the systems designed at Google (Google File System for Google's MapReduce)
 - Open source project under the Apache license
 - But there are also many commercial implementations (e.g., Cloudera, Hortonworks, MapR)

Hadoop History

- Dec 2004 Google published a paper about GFS
- July 2005 Nutch uses MapReduce
- Feb 2006 Hadoop becomes a Lucene subproject
- Apr 2007 Yahoo! runs it on a 1000-node cluster
- Jan 2008 Hadoop becomes an Apache Top Level Project
- Jul 2008 Hadoop is tested on a 4000 node cluster

Hadoop History

- Feb 2009 The Yahoo! Search Webmap is a Hadoop application that runs on more than 10,000 core Linux cluster
- June 2009 Yahoo! made available the source code of its production version of Hadoop
- In 2010 Facebook claimed that they have the largest Hadoop cluster in the world with 21 PB of storage
- On July 27, 2011 they announced the data has grown to 30 PB.

- Core components of Hadoop:
 - Distributed Big Data Processing Infrastructure based on the MapReduce programming paradigm
 - Provides a high-level abstraction view
 - Programmers do not need to care about task scheduling and synchronization
 - Fault-tolerant
 - Node and task failures are automatically managed by the Hadoop system
 - HDFS (Hadoop Distributed File System)
 - High availability distributed storage
 - Fault-tolerant



Paola Vocca

Example with nounder of replication per schunder 2 Big Data



Paola Vocc Example with number of replicas per chunk = 2



Paola Vocca Example with number of Pepile Some rendered Unit di Dati e Big Data

Distributed Big Data Processing Infrastructure

• Separates the what from the how

- Hadoop programs are based on the MapReduce programming paradigm
- MapReduce abstracts away the "distributed" part of the problem (scheduling, synchronization, etc)
 - Programmers focus on what
- The distributed part (scheduling, synchronization, etc) of the problem is handled by the framework
 - The Hadoop infrastructure focuses on how

Distributed Big Data Processing Infrastructure

- But an in-depth knowledge of the Hadoop framework is important to develop efficient applications
 - The design of the application must exploit data locality and limit network usage/data sharing

- Each file is split in "chunks/blocks" that are spread across the servers
 - Each chuck is replicated on different servers (usually there are 3 replicas per chunk)
 - Ensures persistence and availability
 - To increase persistence and availability, replicas are stored in different racks, if it is possible
- Typically each chunk is 64-128MB
- A special file (the master node) stores, for each file, the positions of its
- Chunks
 - The master node is itself replicated
 - A directory for the file system knows where to find the master node
 - The directory itself can be replicated



Example with number of replicas per chunk = 2

Paola Vocca Basi di Dati e di Conoscenza - Basi di Dati e Big Data

- The Master node, a.k.a. Name Nodes in HDFS, is a special node/server that
 - Stores HDFS metadata
 - E.g., the mapping between the name of a file and the location of its chunks
 - Might be replicated
- **Client applications:** file access through HDFS APIs
 - Talk to the master node to find data/chuck servers associated with the file of interest
 - Connect to the selected chunk servers to access data

• An HDFS cluster has two types of nodes:

- Multiple DataNodes
- The NameNode



- The **datanodes** just store and retrieve the blocks when they are told to (by clients or the namenode)
- The **namenode**:
 - Manages the filesystem tree and the metadata for all the files and directories
 - Knows the datanodes on which all the blocks for a given file are located
- Without the namenode HDFS cannot be used
- It is important to make the namenode resilient to failure

HDFS I/O

- An application client wishing to read a file (or a portion thereof) must first contact the **namenode** to determine where the actual data is stored
- In response to the client request the namenode returns:
 - the relevant block id
 - the location where the block is held (i.e., which datanode)
- The client then contacts the datanode to retrieve the data.
- Blocks are themselves stored on standard single-machine file systems
 - HDFS lies on top of the standard OS stack
- Important feature of the design:
 - data is never moved through the namenode
 - all data transfer occurs directly between clients and datanodes
 - communications with the namenode only involves transfer of metadata

Hadoop ecosystem

- Many Hadoop-related projects/systems are available
 - Hive
 - A distributed relational database, based on MapReduce, for querying data stored in HDFS by means of a query language based on SQL
 - HBase
 - A distributed column-oriented database that uses HDFS for storing data
 - Pig
 - A data flow language and execution environment, based on MapReduce, for exploring very large datasets

Hadoop ecosystem

• Sqoop

• A tool for efficiently moving data from traditional relational databases and external flat file sources to HDFS

• ZooKeeper

• A distributed coordination service. It provides primitives such as distributed locks

•

• Each project/system addresses one specific class of problems

MapReduce: introduction

"Big Ideas" of large-scale computing

- Scale "out", not "up"
 - Limits of Symmetric Multi-Processing and large shared-memory machines
- Hide system-level details from the application developer
 - Concurrent programs are difficult to reason about and harder to debug
- Move processing to the data
 - Cluster have limited bandwidth
- Process data sequentially, avoid random access
 - Seeks are expensive, disk throughput is reasonable
- Seamless scalability
 - From the mythical man-month to the tradable machine-hour

Warmup: Word count

• Input

- A large textual file of words
- Problem
 - Count the number of times each distinct word appears in the file

• Output

- A list of pairs <word, counter>
- counting the number of occurrences of each specific word in the input file

Word count

• Case 1: Entire file fits in main memory

- A traditional single node approach is probably the most efficient solution in this case
 - The complexity and overheads of a distributed system affects the performance when files are "small"
 - "small" depends on the resources you have

• Case 2: File too large to fit in main memory

- How can we split this problem in a set of (almost) **independent subtasks**, and
- execute them in parallel on a cluster of servers?

Word count

Suppose that

- 1. The cluster has **3 servers**
- 2. The content of the input file is
 - "Toy example file for Hadoop. Hadoop running example."
- 3. The input file is split into **2 chunks**
- 4. The number of **replicas** is **1**

Word count with a very large file

CPU Mem

Toy example file for Hadoop. Hadoop running example.








- The problem can be easily parallelized
 - 1. Each server processes its chunk of data and counts the number of times each word appears in its own chunk
 - Each server can execute its sub-task independently from the other servers of the cluster

→ synchronization is not needed in this phase

• The output generated from each chunk by each server represents a **partial result**







- Each server sends its local (partial) list of pairs <word, number of occurrences in its chunk> to a server that is in charge of aggregating all local results and computing the global result
 - The server in charge of computing the global result needs to receive all the local (partial) results to compute and emit the final list
 - → A synchronization operation is needed in this phase

Word Count: a more realistic example

Case 2: File too large to fit in main memory

- Suppose that
 - The file size is 100 GB and the number of distinct words occurring in it is at most 1,000
 - The cluster has 101 servers
 - The file is spread across 100 servers and each of these servers contains one (different) chunk of the input file
 - i.e., the file is optimally spread across 100 servers (each server contains 1/100 of the file in its local hard drives)

Word Count: complexity

- Each server reads 1 GB of data from its local hard drive (it reads one chunk from HDFS)
 - Few seconds
- Each local list consists of at most 1,000 pairs (because the number of distinct words is 1,000)
 - Few MBs
- **The maximum amount of data sent** on the network is 100 x size of local list (number of servers x local list size)
 - Some MBs

Word Count: scalability

- We can define **scalability** along two dimensions
 - In terms of data:
 - Given twice the amount of data, the word count algorithm takes approximately no more than twice as long to run
 - Each server processes 2 x data => 2 x execution time to compute local list
 - In terms of resources
 - Given twice the number of servers, the word count algorithm takes approximately no more than half as long to run
 - Each server processes $\frac{1}{2}$ x data => $\frac{1}{2}$ x execution time to compute local list

Word Count: scalability

- The time needed to send local results to the node in charge of computing the final result and the computation of the final result are considered negligible in this running example
- Frequently, this assumption is not true
 - It depends
 - on the complexity of the problem
 - on the ability of the developer to limit the amount of data sent on the network

MapReduce-approach key ideas

• Scale "out", not "up"

- Increase the number of servers, avoiding to upgrade the resources (CPU, memory) of the current ones
- Move processing to data
 - The network has a limited bandwidth
- Process data sequentially, avoid random access
 - Seek operations are expensive
 - Big data applications usually read and analyze all input records/objects
 - Random access is useless

Data locality

- Traditional distributed systems (e.g., HPC) move data to computing nodes (servers)
 - This approach cannot be used to process TBs of data
 - The network bandwidth is limited
- Hadoop moves code to data
 - Code (few KB) is copied and executed on the servers where the chunks of data are stored
 - This approach is based on "data locality"

Hadoop and MapReduce

- Hadoop/MapReduce is designed for
 - Batch processing involving (mostly) full scans of the input data
 - Data-intensive applications
 - Read and process the whole Web (e.g., PageRank computation)
 - Read and process the whole Social Graph (e.g., LinkPrediction, a.k.a. "friend suggestion")
 - Log analysis (e.g., Network traces, Smart-meter data, ..)

What can we do with MapReduce?

- Solving complex problems is difficult
- However, there are several important problems that can be adapted to MapReduce
 - Log analysis
 - PageRank computation
 - Social graph analysis
 - Sensor data analysis
 - Smart-city data analysis
 - Network capture analysis

Hadoop and MapReduce

- Hadoop/MapReduce is not the panacea for all Big Data problems
- Hadoop/MapReduce does not feet well
 - Iterative problems
 - Recursive problems
 - Stream data processing
 - Real-time processing

MapReduce and Functional programming

- The MapReduce programming paradigm is based on the basic concepts of Functional programming
- MapReduce "implements" a subset of functional programming
 - The programming model appears quite limited and strict
 - Everything is based on two "functions" with predefined signatures



Building blocks: Map and Reduce

- MapReduce is based on two main "building blocks"
 - Map and Reduce functions
- Map function
 - It is applied over each element of an input data set and emits a set of (key, value) pairs
- Reduce function
 - It is applied over each set of (key, value) pairs (emitted by the map function) with the same key and emits a set of (key, value) pairs → Final result

• Input

- A textual file (i.e., a list of words)
- Problem
 - Count the number of times each distinct word appears in the file
- Output
 - A list of pairs <word, number of occurrences in the input file>

• The input textual file is considered as a list of words L

L = [toy, example, toy, example, hadoop]

[...] denotes a list. (k, v) denotes a key-value pair.

Paola Vocca

L = [toy, example, toy, example, hadoop]

Apply a function on each element

L_m =[(toy, +1),(example, +1), (toy, +1), (example, +1), (hadoop, +1)]

[...] denotes a list. (k, v) denotes a key-value pair.

Paola Vocca

L = [toy, example, toy, example, hadoop]

 $L_m = [(toy, +1), (example, +1), (toy, +1), (example, +1), (hadoop, +1)]$

Group by key

(toy, [+1, +1]) (example, [+1, +1]) (hadoop, [+1])



Paola Vocca

L = [toy, example, toy, example, hadoop]

L_m =[(toy, +1),(example, +1), (toy, +1), (example, +1), (hadoop, +1)]

(toy, [+1, +1]) (example, [+1, +1]) (hadoop, [+1]) Apply a function on each group

[(toy, 2), (example, 2), (hadoop, 1)]

[...] denotes a list. (k, v) denotes a key-value pair.

Paola Vocca



[(toy, 2), (example, 2), (hadoop, 1)]

[...] denotes a list. (k, v) denotes a key-value pair.

Paola Vocca

- The input textual file is considered as a list of words L
- A key-value pair (w, 1) is emitted for each word w in L
 - i.e., the map function is

 $\mathbf{m}(\mathbf{w}) = (\mathbf{w}, 1)$

• A new list of (key, value) pairs \mathbf{L}_{m} is generated

- The key-value pairs in L_m are aggregated by key (i.e., by word w in our example)
 - One group G_{w} is generated for each word w
 - Each group **G**_w is a key-list pair (**w**, [list of values]) where [list of values] contains all the values of the pairs associated with the word **w**
 - i.e., [list of values] is a list of [1, 1, 1, ...] in our example
 - Given a group **G**_w the number of ones [1, 1, 1, ...] is equal to the occurrences of word w in the input file

- A key-value pair (w, sum G_w.[list of values]) is emitted for each group G_w
 - i.e., the **reduce** function is

 $r(G_w) = (w, sum(G_w.[list of values]))$

- The list of emitted pairs is the result of the word count problem
 - One pair (word **w**, num. of occurrences) for each word in our running example

MapReduce: Map

- The Map phase can be viewed as a transformation over each element of a data set
 - This transformation is a function **m** defined by developers
 - **m** is invoked one time for each input element
 - Each invocation of **m** happens in **isolation**
 - The application of **m** to each element of a data set can be parallelized in a straightforward manner

MapReduce: Reduce

- The **Reduce** phase can be viewed as an **aggregate** operation
 - The aggregate function is a function **r** defined by developers
 - **r** is invoked one time for each distinct key and aggregates all the values associated with it
 - Also the reduce phase can be performed in parallel and in isolation
 - Each group of key-value pairs with the same key can be processed in isolation

MapReduce: Shiffle and Sort

- The shuffle and sort phase is always the same
 - i.e., group the output of the map phase by key
 - It does not need to be defined by developers
 - It is already provided by the Hadoop system

Data Structure

• Key-value pair is the basic data structure in MapReduce

- Keys and values can be: integers, float, strings, ...
- They can also be (almost) arbitrary data structures defined by the designer
- Both input and output of a MapReduce program are lists of key-value pairs
 - Note that also the input is a list of key-value pairs

Data Structure

- The design of MapReduce involves
 - Imposing the key-value structure on the input and output data sets
 - E.g., for a collection of Web pages, input keys may be URLs and values may be their HTML content

Formal definition of Map and Reduce functions

- The map and reduce functions are formally defined as follows:
 - map: $(k1, v1) \rightarrow [(k2, v2)]$
 - reduce: (k2, [v2]) \rightarrow [(k3, v3)]
- Since the input data set is a list of key-value pairs, the argument of the map function is a key-value pair



Formal definition of Map and Reduce functions

- Map function
 - map: $(k1, v1) \rightarrow [(k2, v2)]$
- The argument of the map function is a key-value pair
- Note that the map function emits a list of key-value pairs for each input record
 - The list can also be empty



Formal definition of Map and Reduce functions

- Reduce function
 - reduce: (k2, [v2]) \rightarrow [(k3, v3)]
- Note that the reduce function
 - Receives the complete list of values [v2] associated with a specific key k2
 - Emits a list of key-value pairs


MapReduce Algorithms

- In many applications, the key part of the input data set is ignored
 - i.e., usually the map function does not consider the key of its key-value pair argument
 - E.g., word count problem
- Some specific applications exploit also the keys of the input data
 - E.g., keys can be used to uniquely identify records/objects

WordCount using MapReduce: Pseudocode

- Input file: a textual document with one word per line
- The map function is invoked over each word of the input file

```
map(key, value):
// key: offset of the word in the file
// value: a word of the input document
emit(value, 1)
reduce(key, values):
// key: a word; values: a list of integers
occurrences = 0
for each c in values:
occurrences = occurrences + c
emit(key, occurrences)
```