



UNIVERSITÀ
DEGLI STUDI
DI BERGAMO

Tecnologie Cloud & Mobile

TECNOLOGIE CLOUD
E MOBILE

BIG DATA AND APPLICATIONS IN THE CLOUD

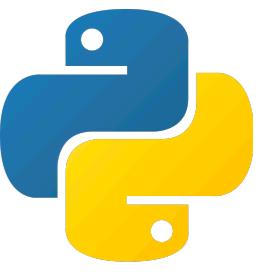
mauro.pelucchi@gmail.com

Mauro Pelucchi

21059 - Tecnologie cloud e mobile

Cosa vedremo?

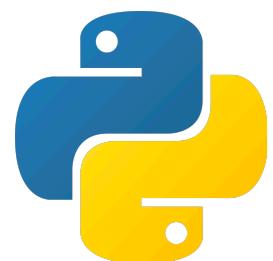
- Python & PySpark
- MongoDB
- Aws Glue



Python

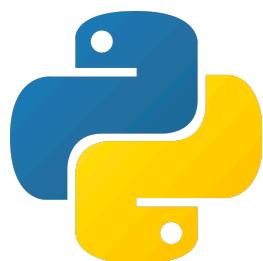
Python

Python in 10 minutes...



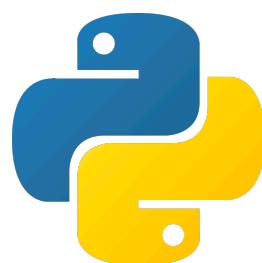
Python

- Interpreted language: work with an evaluator for language expressions
- Dynamically typed: variables do not have a predefined type
- Rich, built-in collection types:
 - Tuples
 - Dictionaries
 - Sets
- Concise
- Open source general-purpose language
- Object Oriented, Procedural, Functional
- Easy to interface with C/ObjC/Java/Fortran
- Great interactive environment



Features

- Indentation instead of braces
- Several sequence types
 - Strings '...': made of characters, immutable
 - Lists [...]: made of anything, mutable
 - Tuples (...) : made of anything, immutable
- Powerful subscripting (slicing)
- Exceptions as in Java
- Simple object system
- Iterators (like Java)



References

https://github.com/Naviden/Python_Introduction

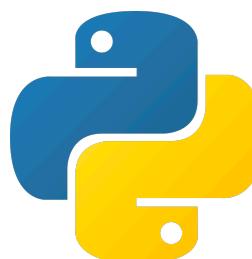
The screenshot shows the GitHub repository page for 'Naviden / Python_Introduction'. At the top, there's a header with the repository name, a 'Watch' button (1 watch), a 'Star' button (0 stars), and a 'Fork' button (0 forks). Below the header is a navigation bar with links for 'Code' (which is selected and highlighted in orange), 'Issues' (0), 'Pull requests' (0), 'Actions', 'Projects' (0), 'Wiki', 'Security', and 'Insights'. Underneath the navigation bar, it says 'Branch: master' and shows the file 'Python_Introduction / Py_intro.ipynb'. To the right of the file name are 'Find file' and 'Copy path' buttons. Below this, there's a commit history section showing a single commit by 'Naviden' with the message 'Rrformatting sections' and the commit hash '7c320ca' dated 'on Feb 23'. It also indicates '1 contributor'. At the bottom of the commit history, there are download and copy buttons. The footer of the page includes a line count ('19833 lines (19833 sloc) | 731 KB') and a set of icons for file operations.

About this course

This introduction course, presented for the first time in BIBDA Master, tends to give you an overview of Python programming language and its general structure as in the future modules of Master like web scrapping, Machine learning and text mining Python will be used as the principal programming language. Please note that due to the vast range of arguments/tools related to Python language we won't/can't cover them all here as this course designed to help you not to become a Python programmers but to learn its basics and enable you to expand your data science skills using it as a modern programming language.

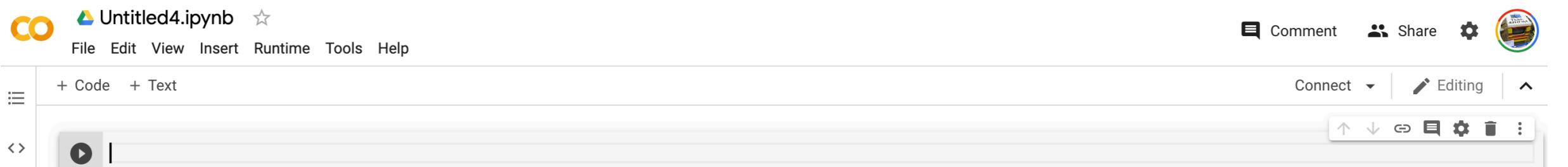
Course structure

Our goal here is not to provide a deep insight into Python and act as an official reference since you can easily can access hundreds of online sources to learn about each and every detail of Python. Given the wide range of students' backgrounds, there is no assumption about prior knowledge about programming or data analysis. We've tried to provide you with simple examples which give you a clear understanding of functionality of different tools. There are also parts called **YOUR TURN** which ask you to put into practice what you have learned during the lesson. These parts are accompanied by three different symbols which show the difficulty of the exercise:



References

- The Python Tutorial (<http://docs.python.org/tutorial/>)
- PEP 8- Style Guide for Python Code
 - <http://www.python.org/dev/peps/pep-0008/>
- Google Colab



Which Python?

- **Python 3!!!!**

Python 2.7 will retire in...

0
Years

0
Months

0
Days

0
Hours

0
Minutes

0
Seconds

[Enable Guido Mode](#) [Huh?](#)

What's all this, then?

Python 2.7 [will not be maintained past 2020](#). Originally, there was no official date. Recently, that date has been updated to [January 1, 2020](#). This clock has been updated accordingly. My original idea was to throw a Python 2 Celebration of Life party at PyCon 2020, to celebrate everything Python 2 did for us. That idea still stands. (If this sounds interesting to you, email pythonclockorg@gmail.com).

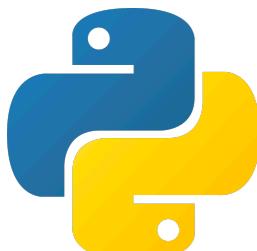
Python 2, thank you for your years of faithful service.

Python 3, your time is now.

How do I get started?

If the code you care about is still on Python 2, that's totally understandable. Most of PyPI's popular packages now [work on Python 2 and 3](#), and more are being added every day. Additionally, a number of critical Python projects have [pledged to stop supporting Python 2 soon](#). To ease the transition, the [official porting guide](#) has advice for running Python 2 code in Python 3.

If you're here, you'd probably enjoy the [full python release schedule](#) (a glitch maintained by [Dustin Ingram](#))



A Code Sample

Open Google Colab

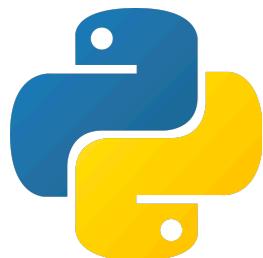
<https://colab.research.google.com/drive/1eoOuTMfnPC3oVJLykRwn388WMp4k6i67#scrollTo=y-EyvIQa2fTA>

```
[1] x = 5  
    y = 6  
    z = x + y  
    print(z)
```

↳ 11

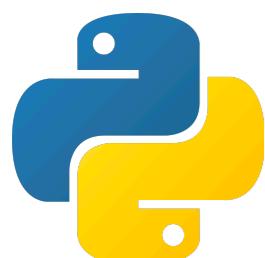
```
▶ if z > 10:  
    print("z > 10")  
else:  
    print("z <= 10")
```

↳ z > 10



A Code sample

- Assignment uses `=` and comparison uses `==`
- For numbers `+` `-` `*` `/` `%` are as expected
- Special use of `+` for string concatenation
- Special use of `%` for string formatting
- The basic printing command is `print(...)`.
- The first assignment to a variable creates it
- Variable types don't need to be declared
 - Python figures out the variable types on its own

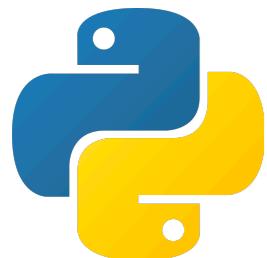


Strings

▼ Strings

```
▶ hello = 'Hello, '  
      name = "Unibg"  
      print(f"{hello} {name}!!!")
```

```
⇨ Hello, Unibg!!!
```



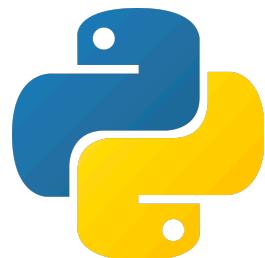
Comments

- Start comments with # – the rest of line is ignored.

- ▼ Comments

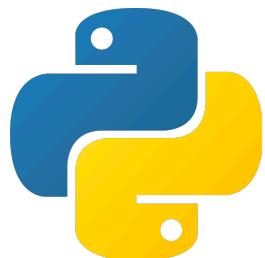
```
# comments  
print("Go")
```

□ Go



Libraries

- numpy
 - Offers Matlab-ish capabilities within Python
 - Fast array operations
 - 2D arrays, multi-D arrays, linear algebra etc.
- matplotlib
 - High quality plotting library
- pandas
 - Data processing
- sklearn
 - Machine Learning and ...
- TensorFlow/keras
 - Deep Learning



Control Flows

```
if x == 3:  
    print "x equals 3."  
elif x == 2:  
    print "x equals 2."  
else:  
    print "x equals something else."  
print "This is outside the 'if'."  
  
assert(number_of_players < 5)
```

```
x = 3  
while x < 10:  
    if x > 7:  
        x += 2  
        continue  
    x = x + 1  
    print "Still in the loop."  
    if x == 8:  
        break  
print "Outside of the loop."
```

```
for x in range(10):  
    if x > 7:  
        x += 2  
        continue  
    x = x + 1  
    print "Still in the loop."  
    if x == 8:  
        break  
print "Outside of the loop."
```

Control Flows

▼ Control Flows

```
[11] my_list = ['a', 'b', 'c']
     for item in my_list:
         print(item)
```

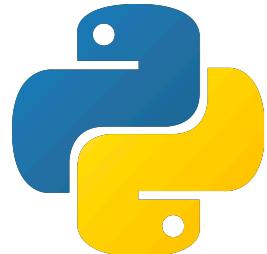
```
⇨ a
    b
    c
```

```
▶ for c in range(0, 10):
    print(c)
```

```
⇨ 0
    1
    2
    3
    4
    5
    6
    7
    8
    9
```

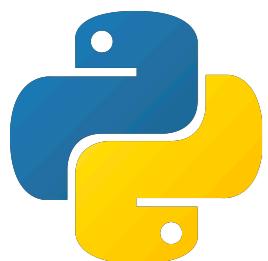
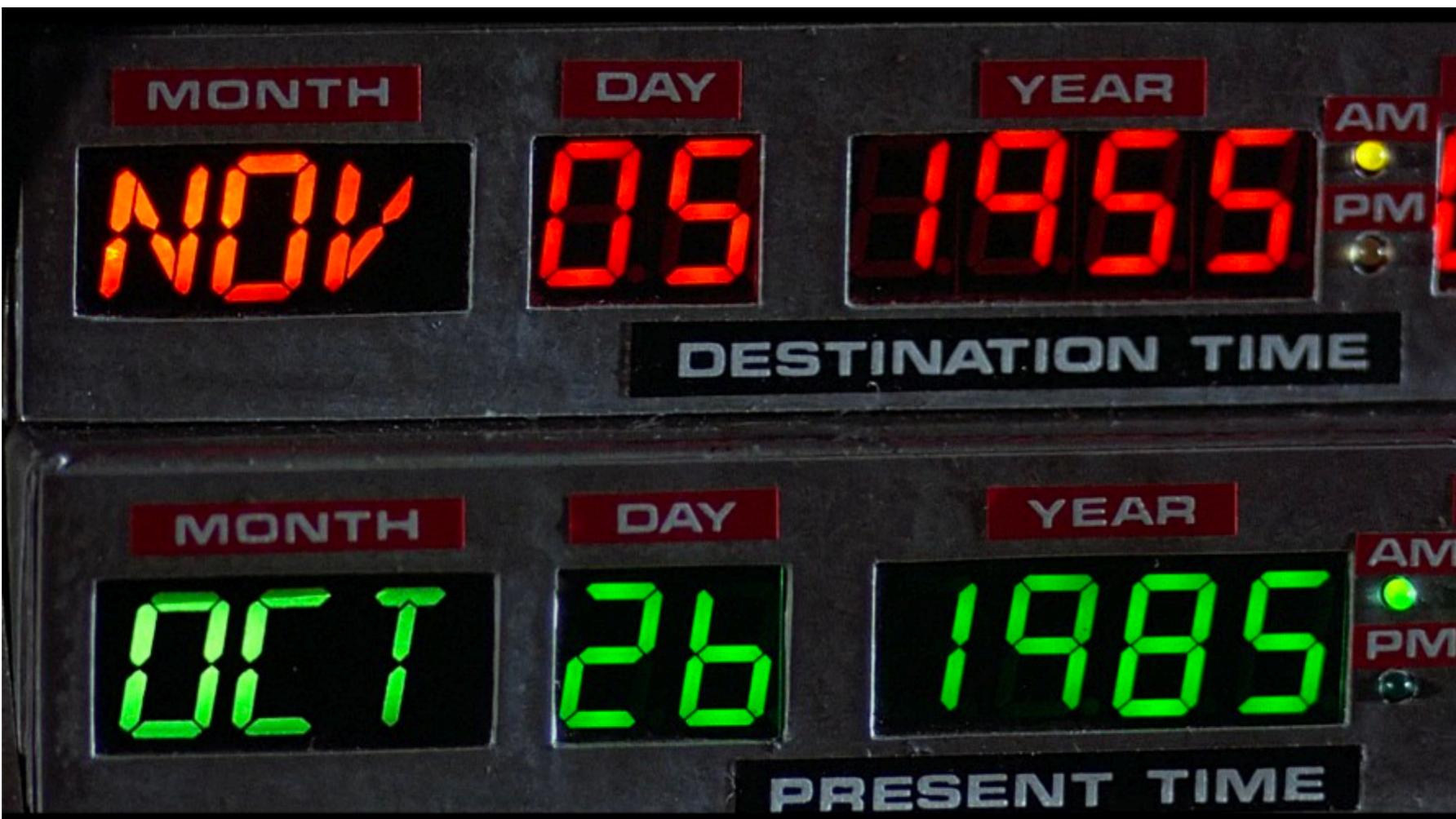
Assignment

- In Google Colab
 - Calcolate in Python in fattoriale di un numero
 - Esempio $5! = 5 * 4 * 3 * 2 * 1 = 120$



Assignment

Ci vendiamo tra 15 minuti...



Solutions

Factorial

```
[9] n=int(input('Digit a number: '))
    factorial=1

    for i in range(1,n+1):
        factorial = factorial*i

    print(factorial)
```

```
↳ Digit a number: 10
3628800
```

```
▶ n=int(input('Digit a number: '))

def factorial(n):
    if n==0:
        return 1
    else:
        return n*factorial(n-1)

print ('Factorial: ', factorial(n))
```

```
↳ Digit a number: 10
Factorial: 3628800
```

PySpark



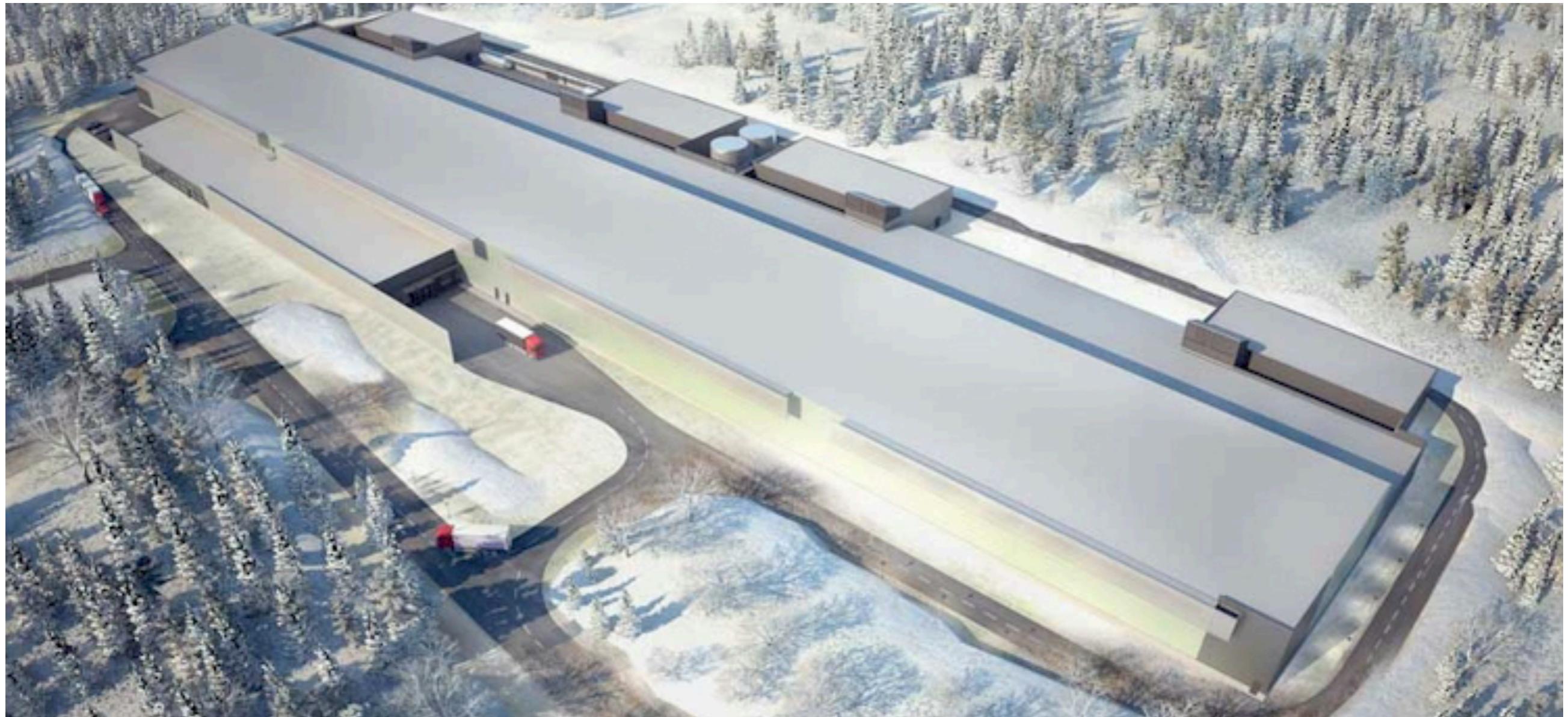
Data is huge



- Banks, city councils, governments, shops, etc.
- Facebook processes > 750TB/day of data
- > 48k iPhones every day
- > 7PB of photo storage / month

This requires computers (a lot of them!)

Facebook datacenter



Scale Up vs Scale Down



- **Scale-up**
 - Increasing the power of your computer (i.e, disk, memory, processor)
- **Scale-out**
 - Use many standard computers and **distribute data** and **computation** over them

Data is fast



- Twitter fire hose
 - In 2011, 1.000 Tweets per second (TPS)
 - In 2014, 20.000 TPS
 - With peaks: 143K TPS
- Services on top
 - Not only internet companies!
 - Stock exchange, sensors in water network, smart-cities, fitness trackers, etc.

Perché Spark e Hadoop?



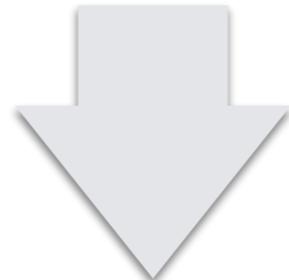
Challenges

Grossi volumi di dati

Dati non strutturati

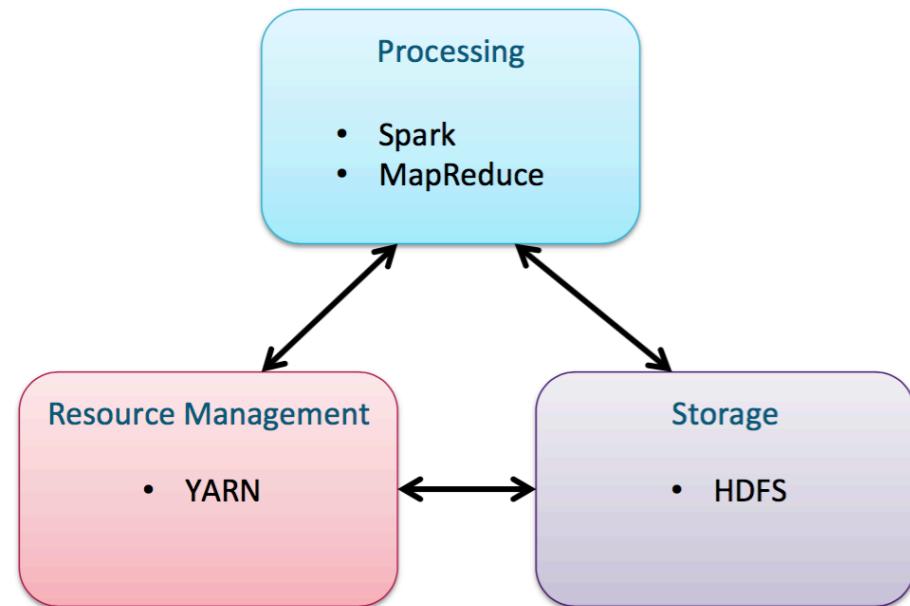
Complessità dei sistemi distribuiti

Applicazioni complesse (ELT, Text mining, modelli predittivi, ...)



Hadoop is a framework for distributed storage and processing

Perché Spark e Hadoop?

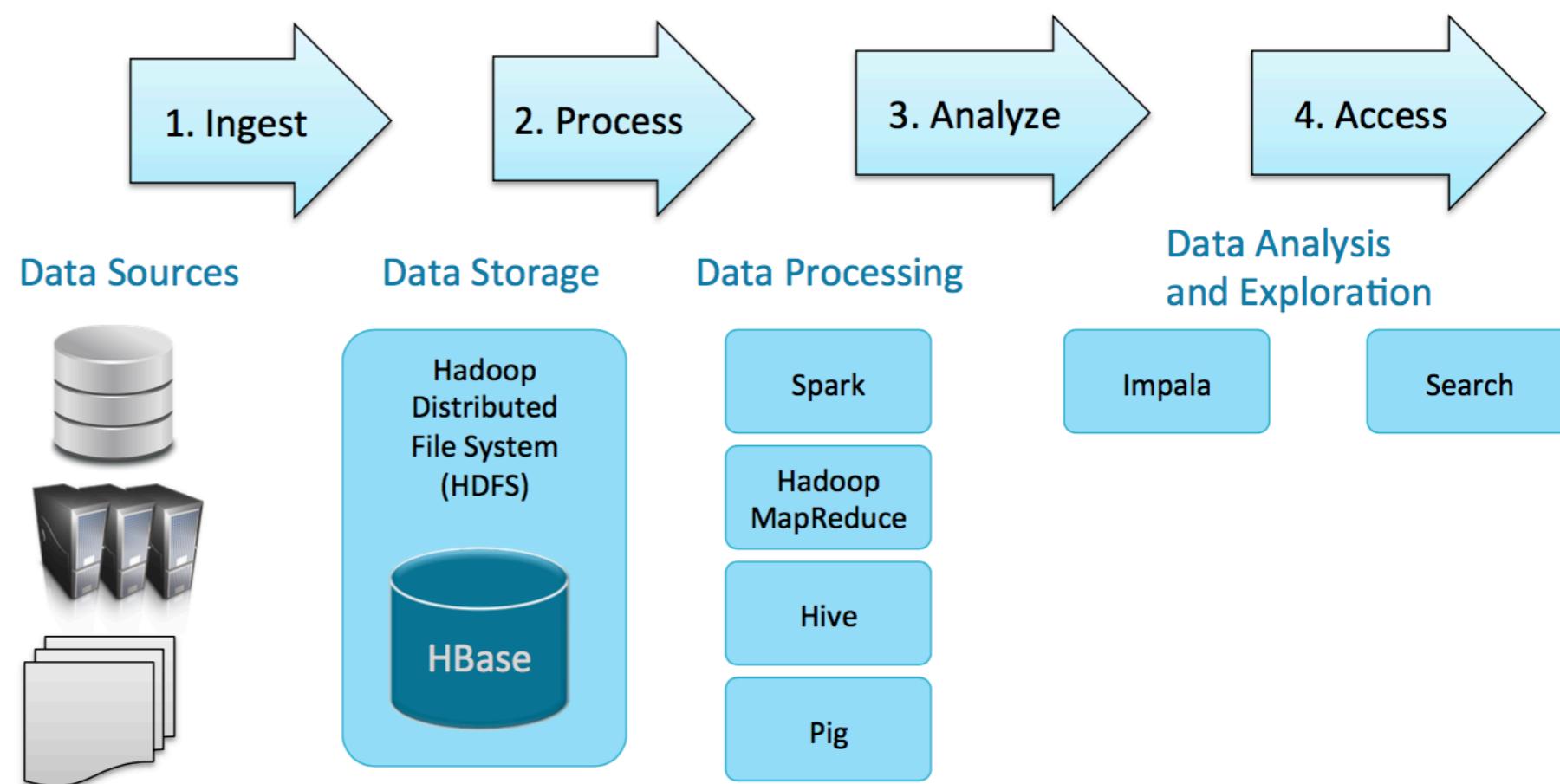


Spark is large-scale data processing engine

- General purpose
- Batch Processing
- Machine learning
- Business intelligence
- Streaming

Hadoop: Bring the program to the data rather than the data to the program

- Data storage
- Scalable
- Distributed
- Fault tolerant



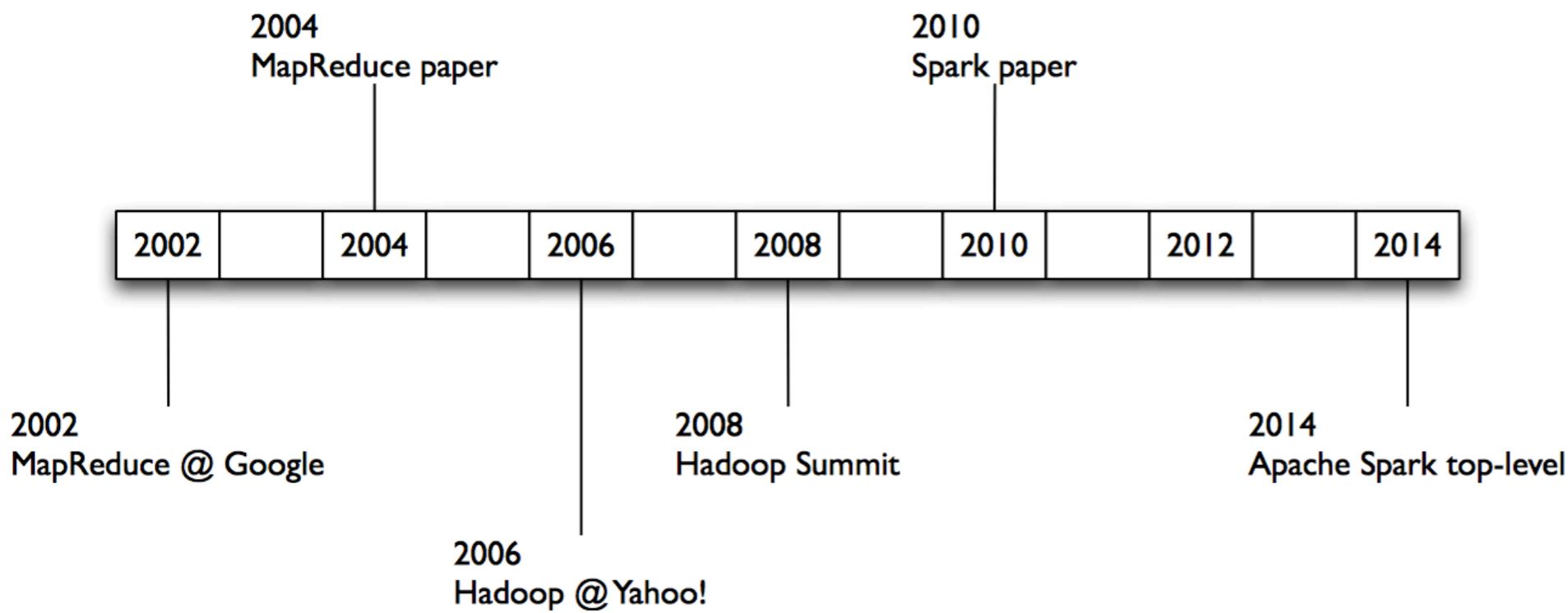
Spark



Apache Spark is an open-source framework for processing **huge volumes of data** (big data) with speed and simplicity.

Spark: Cluster Computing with Working Sets

Matei Zaharia, Mosharaf Chowdhury, Michael J. Franklin, Scott Shenker, Ion Stoica
University of California, Berkeley



Matei Zaharia



Spark



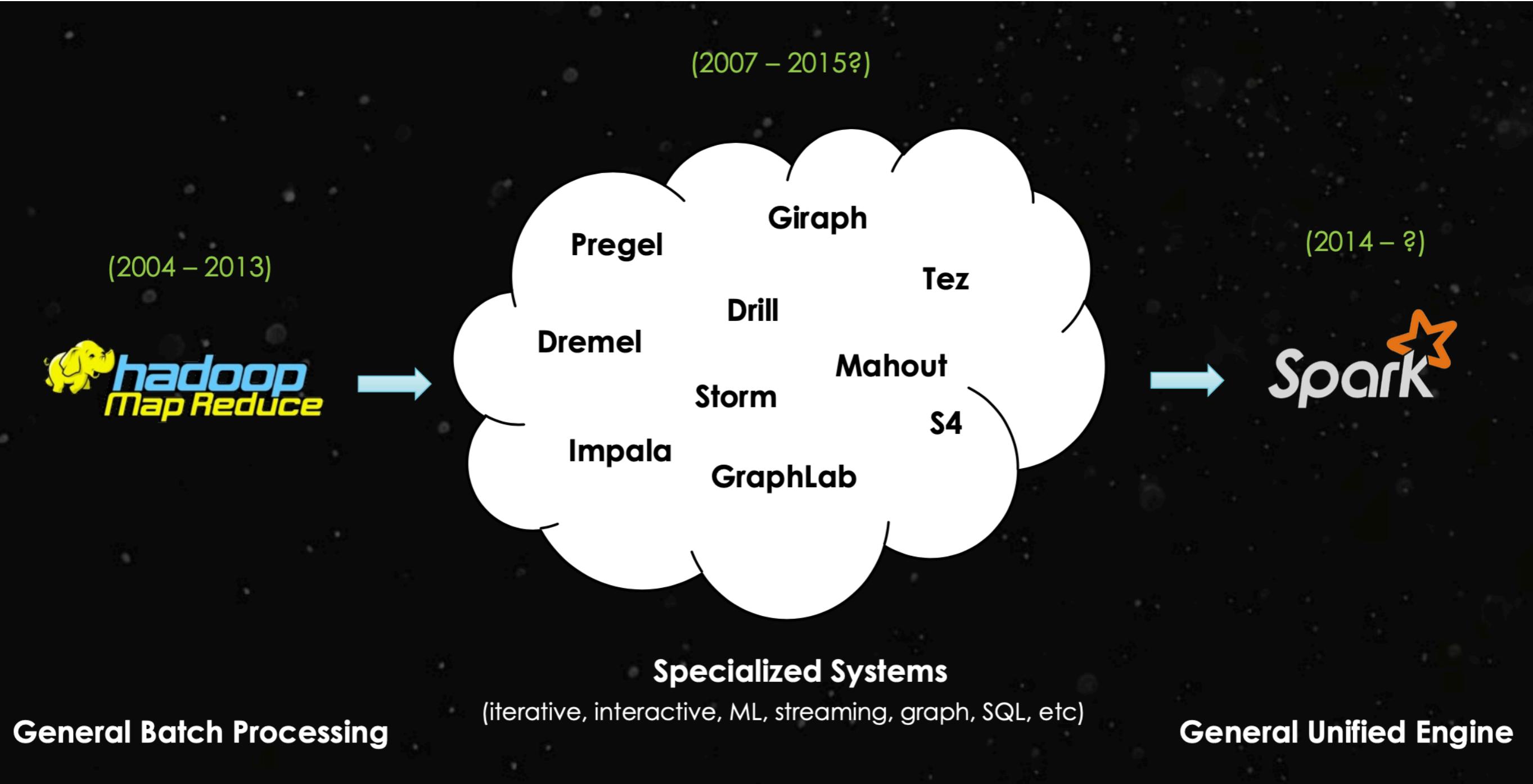
- Sviluppato a UC Berkeley (2010), estende il modello MapReduce per supportare in modo efficace più tipi di operazioni, query interattive e stream processing
 - Scritto in **Scala**
 - Compatibile con **Hadoop**
 - API per **Java**, **Scala**, **Python**
- Semplice, funzionale —> modello di programmazione ad alto livello
- **Veloce**
- **In memory caching, query optimizations**
- **General purpose**, ideale per **ETL** e **Machine Learning**
- Ideale per dati non strutturati

Spark Vs Hadoop



Spark supporta molte tipologie di tasks come interrogazioni **SQL**, applicazioni **streaming**, **machine learning** e **graph operations**.

Hadoop MR è ottimo per lavori pesanti, specifici e complessi su quantità massicce di dati. In ogni caso, molto spesso Spark è una soluzione ottima anche in questi casi perché ottimizza l'uso della memoria e offre un modello di programmazione ad alto livello (semplice).



Why Spark?



Spark has become **the default processing engine** for a myriad of engineering & science problems



Spark Vs Hadoop (MR)



- Generalized patterns → stesso framework per molti casi d'uso
- Lazy evaluation → reduces wait states, better pipelining
- Functional programming → più facile scrivere applicazioni complesse, applicazioni più facili da mantenere

Spark - Lo stack Spark



Spark SQL
structured data

Spark Streaming
real-time

MLlib
Machine Learning

GraphX
Graph processing

Spark Core

Standalone Scheduler

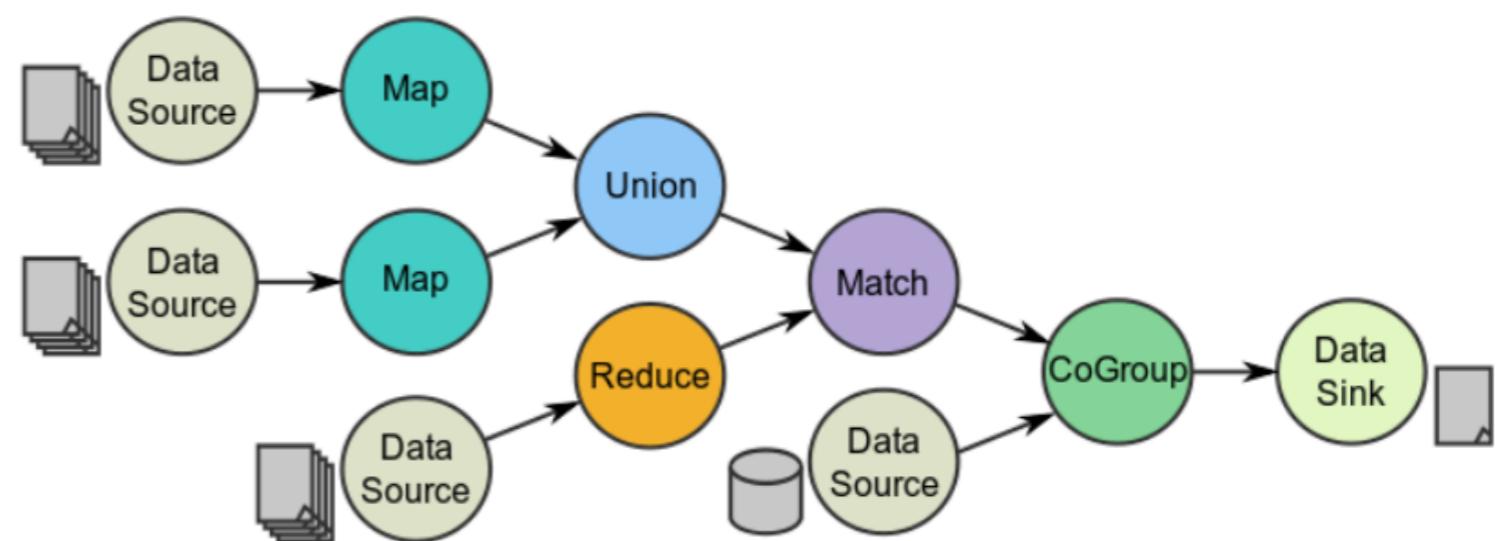
Hadoop / YARN

Mesos

Spark - Programming model



- Spark Programming Model si basa su operatori che possono *andare in parallelo (parallelizable operators)*
- Data flow —> un data flow è composto da sorgenti, operatori e nodi di dati, collegati da link di input e output
- Un JOB è descritto da un grafo aciclico diretto (**directed acyclic graphs o DAG**)

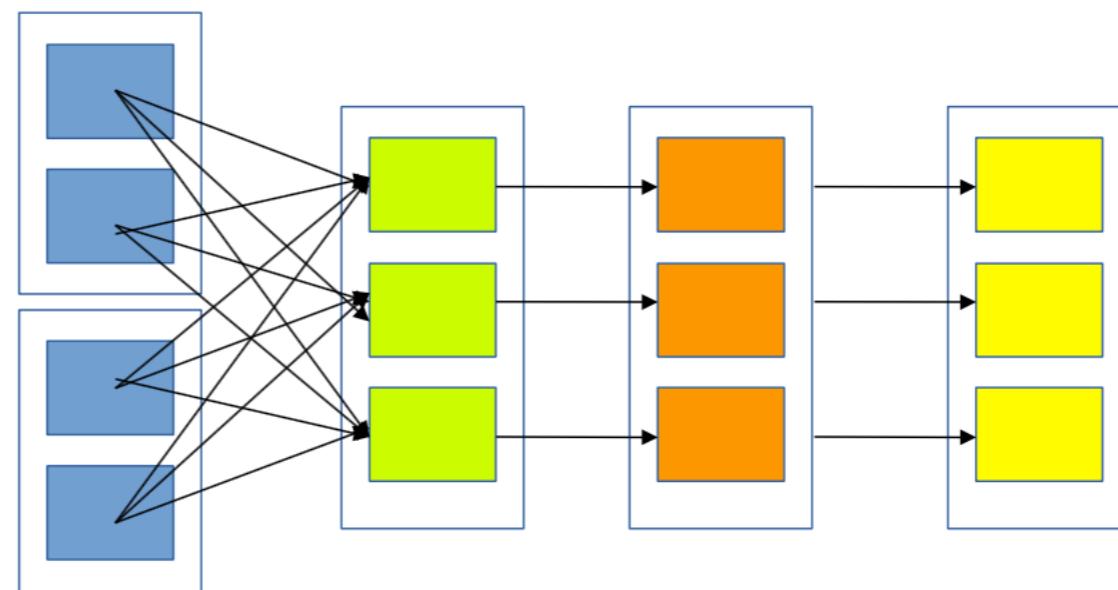
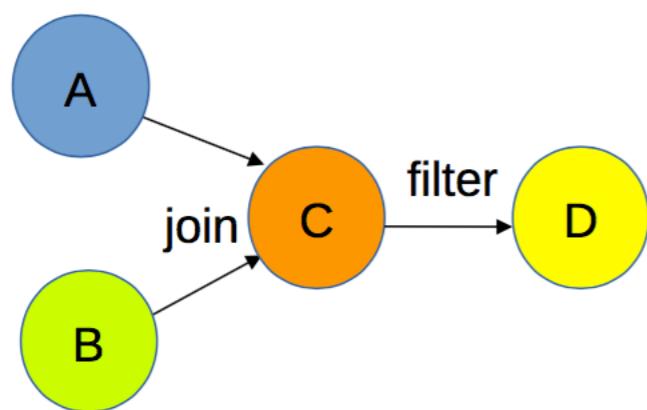


Spark - DAG

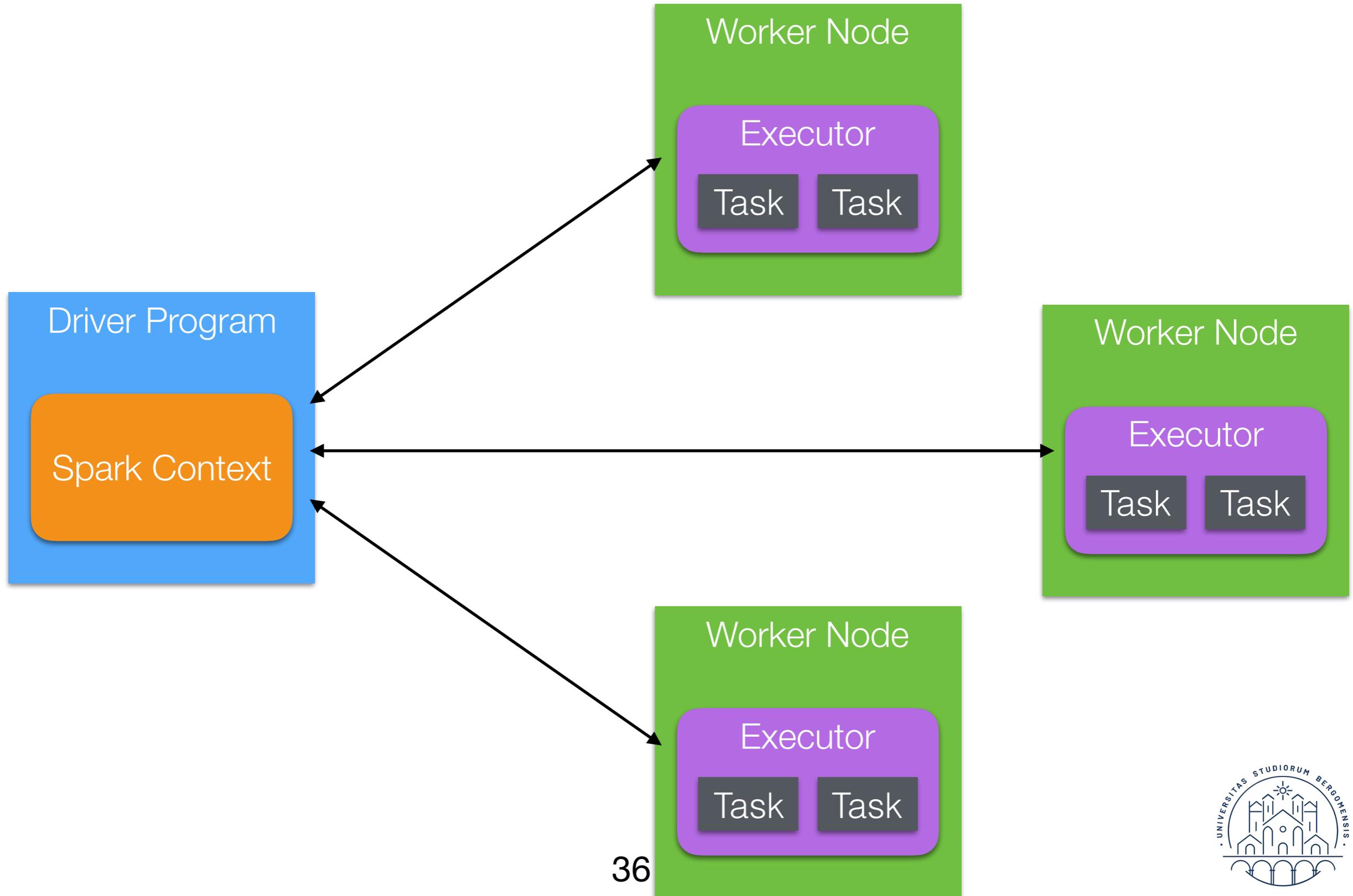


- Directed acyclic graphs o DAG), regole del DAG SCHEDULER

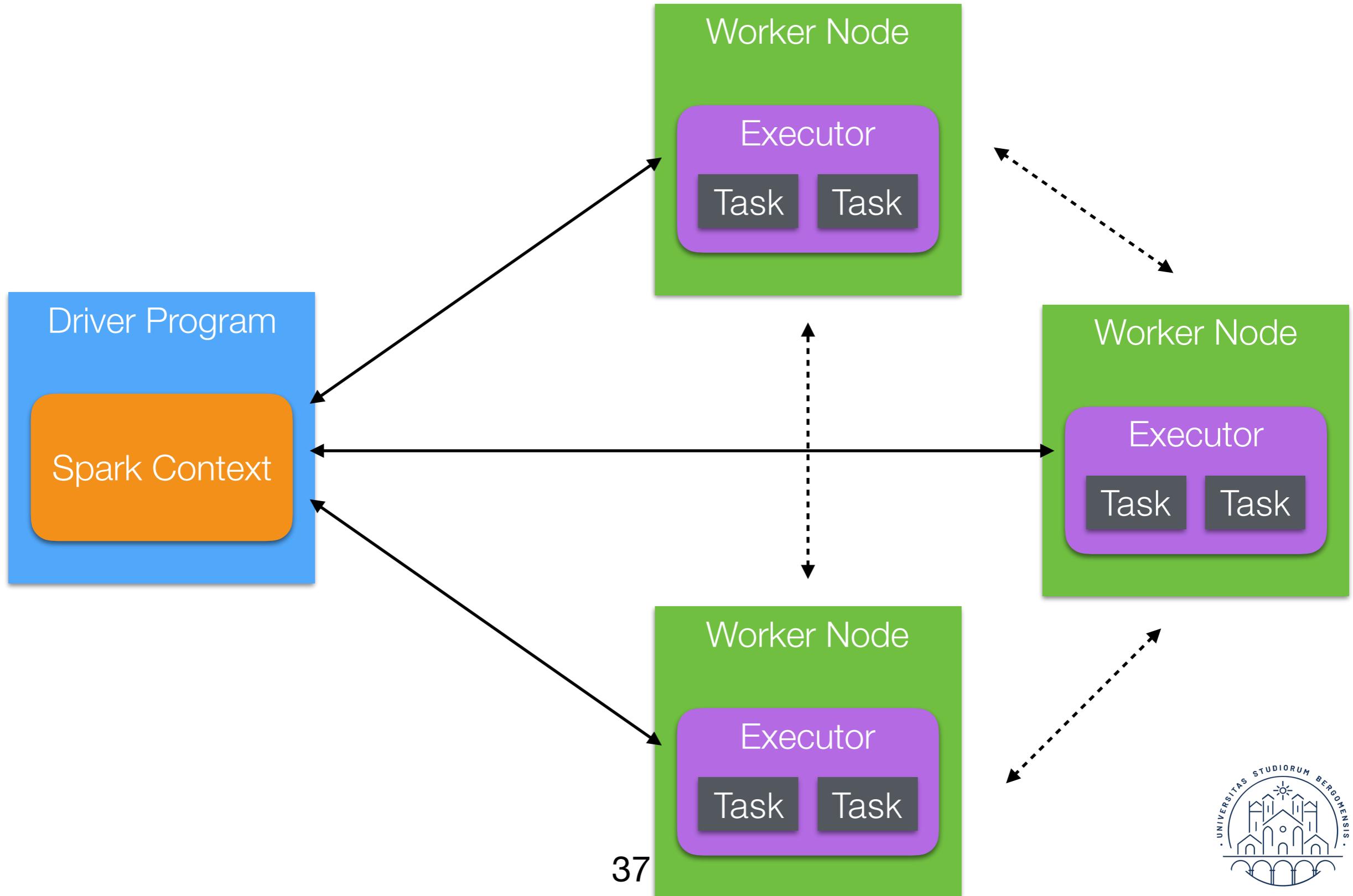
- divide il grafo in una serie di stage di task (insiemi di task)
- sottomette gli stage ai nodi del cluster



Spark - Come funziona?



Spark - Come funziona?



Spark - Come funziona?

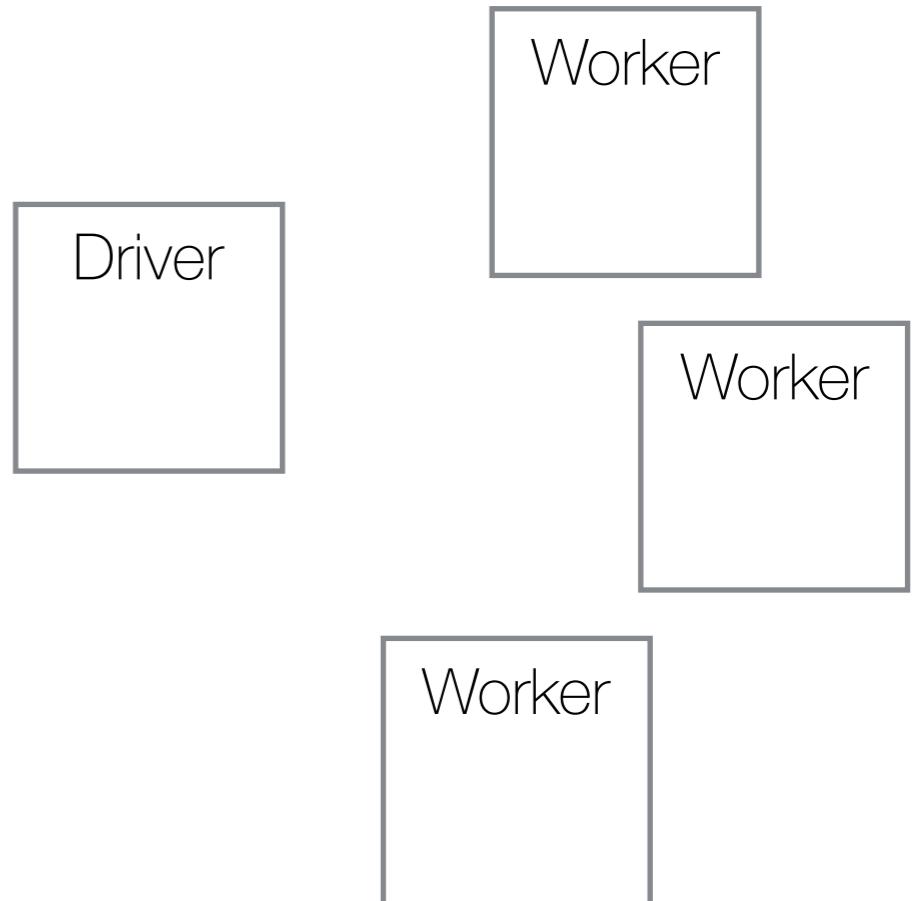


1. L'utente esegue un job attraverso con la funzione `spark-submit`
2. Spark avvia il **driver program** e invoca il metodo `main` specificato dall'utente
3. Il driver program contatta il manager del cluster per chiedere le risorse per avviare gli **executors**
4. Il cluster manager lancia gli **executors** richiesti dal driver program
5. Il driver, in base alle azioni e alle trasformazioni richieste dal job dell'utente, manda il lavoro agli **executors** in forma di **tasks**
6. I tasks sono eseguiti dai processi executor ed il risultato viene salvato
7. Al termine dell'esecuzione vengono rilasciate le risorse

Spark - Esempio



Analizziamo un file di log cercando gli **errori** relativi a **MongoDB**



Spark - Esempio



Carichiamo il file di log da hdfs e prepariamolo in **cache**

```
// leggiamo le righe dal file di log
```

```
val lines = sc.textFile("hdfs://.../var/log/*")
```

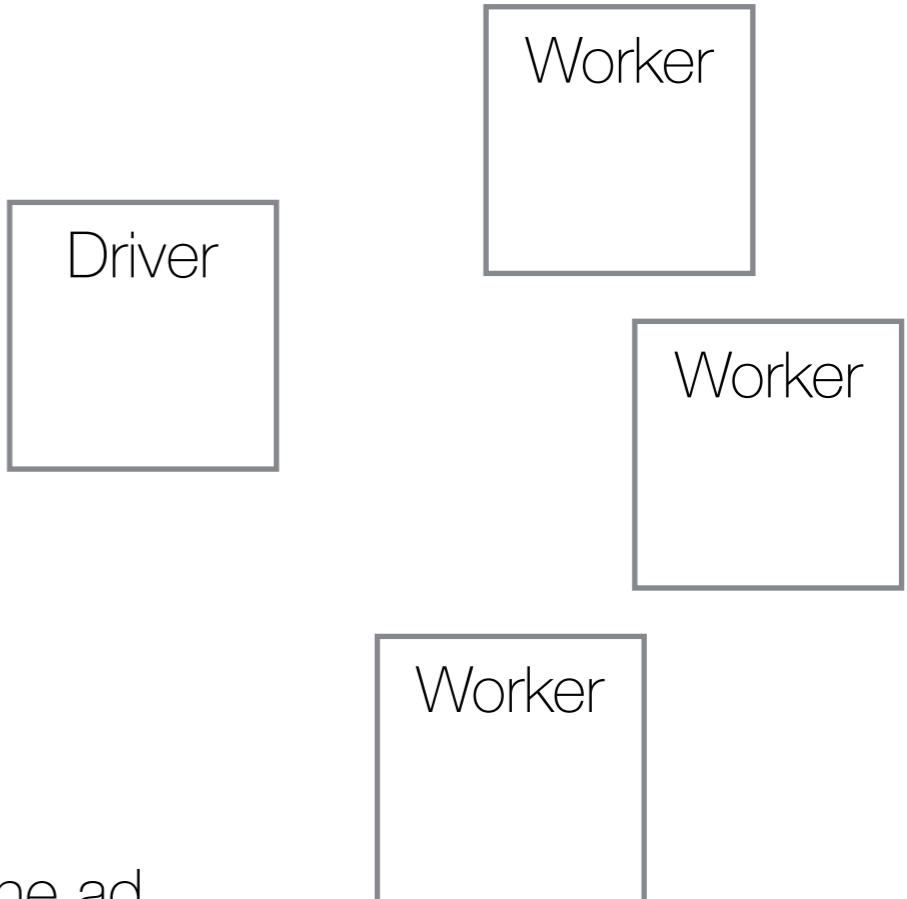
```
// applichiamo un filtro ed una trasformazione
```

```
val errors = lines.filter(_.startsWith("ERROR"))
```

```
val messages = errors.map(_.split("\t")).map(r => r(1))
```

```
// mettiamo in cache i messaggi
```

```
messages.cache()
```



- Spark divide il dataset in 3 partizioni ed assegna ogni partizione ad un nodo
- Non avviene nulla: solo quando richiediamo un'azione sui dati Spark esegue le modifiche

Spark - Esempio



Contiamo le righe che contengono il termine **mongodb**

```
// leggiamo le righe dal file di log
```

```
val lines = sc.textFile("hdfs://.../var/log/*")
```

```
// applichiamo un filtro ed una trasformazione
```

```
val errors = lines.filter(_.startsWith("ERROR"))
```

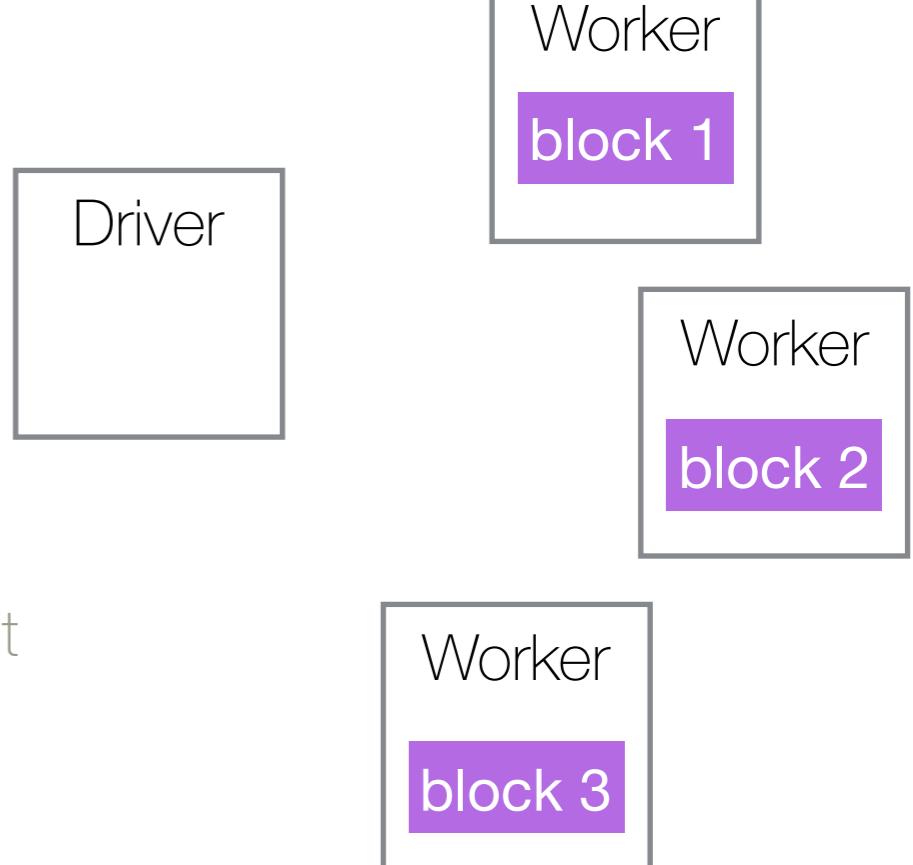
```
val messages = errors.map(_.split("\t")).map(r => r(1))
```

```
// mettiamo in cache i messaggi
```

```
messages.cache()
```

```
// contiamo le righe di errore per MongoDB con l'azione count
```

```
messages.filter(_.contains("mongodb")).count()
```



Spark - Esempio



Contiamo le righe che contengono il termine **mongodb**

```
// leggiamo le righe dal file di log
```

```
val lines = sc.textFile("hdfs://.../var/log/*")
```

```
// applichiamo un filtro ed una trasformazione
```

```
val errors = lines.filter(_.startsWith("ERROR"))
```

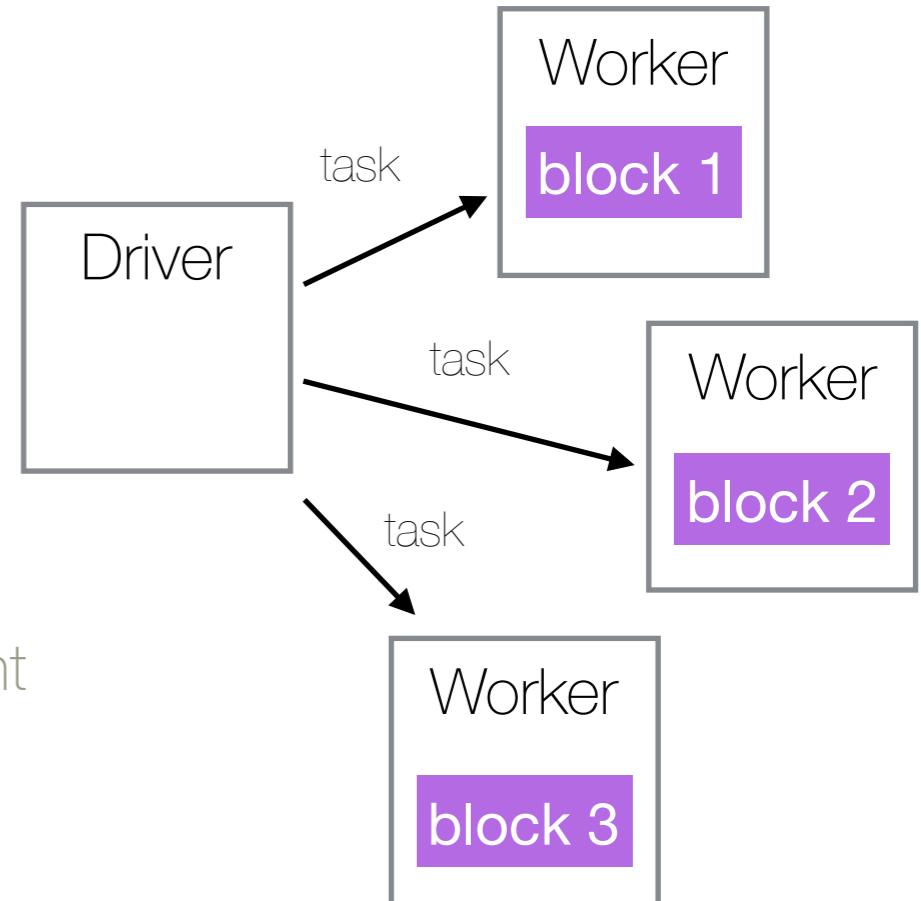
```
val messages = errors.map(_.split("\t")).map(r => r(1))
```

```
// mettiamo in cache i messaggi
```

```
messages.cache()
```

```
// contiamo le righe di errore per MongoDB con l'azione count
```

```
messages.filter(_.contains("mongodb")).count()
```



Spark - Esempio



Contiamo le righe che contengono il termine **mongodb**

```
// leggiamo le righe dal file di log
```

```
val lines = sc.textFile("hdfs://.../var/log/*")
```

```
// applichiamo un filtro ed una trasformazione
```

```
val errors = lines.filter(_.startsWith("ERROR"))
```

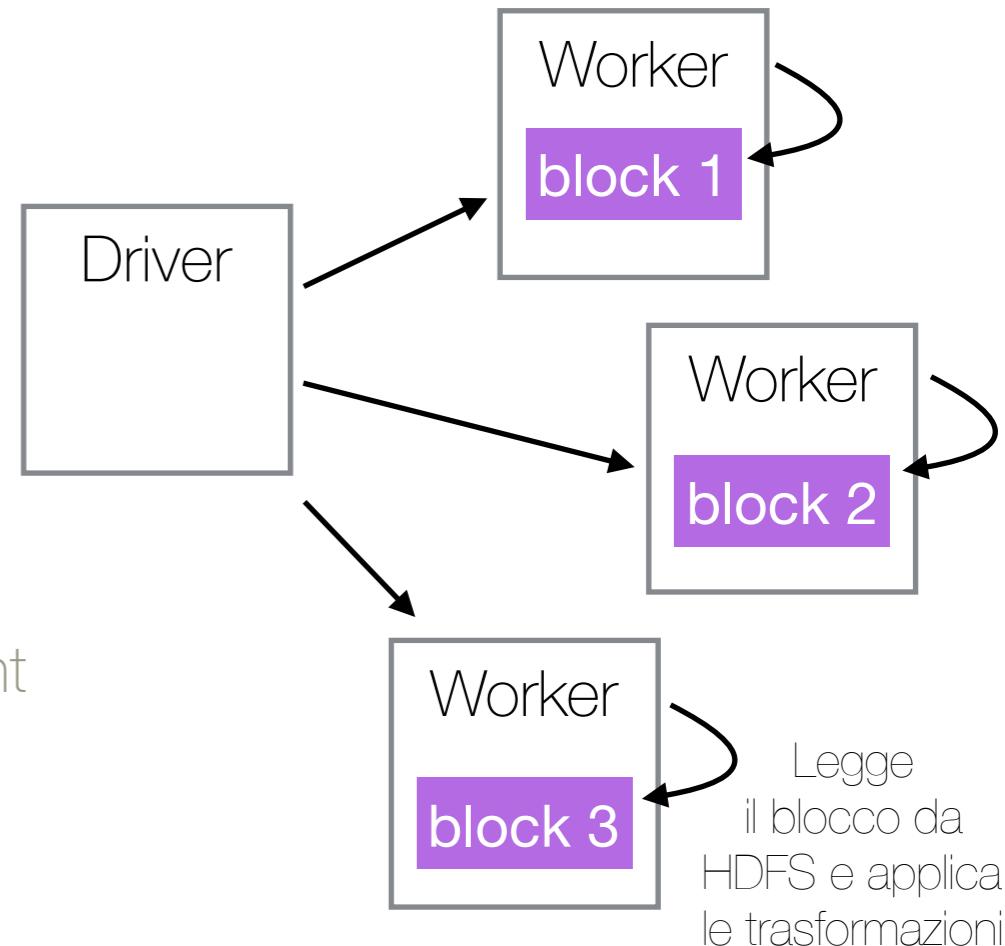
```
val messages = errors.map(_.split("\t")).map(r => r(1))
```

```
// mettiamo in cache i messaggi
```

```
messages.cache()
```

```
// contiamo le righe di errore per MongoDB con l'azione count
```

```
messages.filter(_.contains("mongodb")).count()
```



Spark - Esempio



Contiamo le righe che contengono il termine **mongodb**

```
// leggiamo le righe dal file di log
```

```
val lines = sc.textFile("hdfs://.../var/log/*")
```

```
// applichiamo un filtro ed una trasformazione
```

```
val errors = lines.filter(_.startsWith("ERROR"))
```

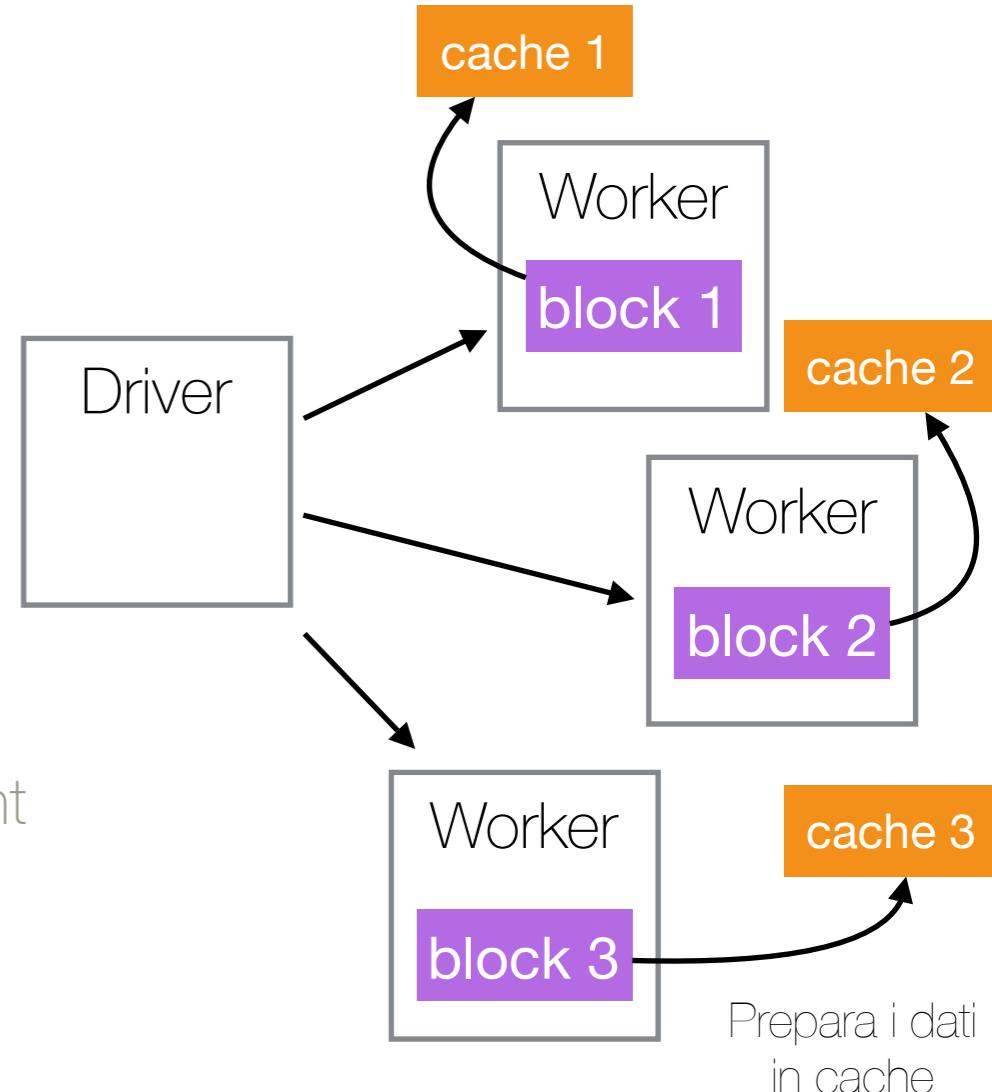
```
val messages = errors.map(_.split("\t")).map(r => r(1))
```

```
// mettiamo in cache i messaggi
```

```
messages.cache()
```

```
// contiamo le righe di errore per MongoDB con l'azione count
```

```
messages.filter(_.contains("mongodb")).count()
```



Spark - Esempio

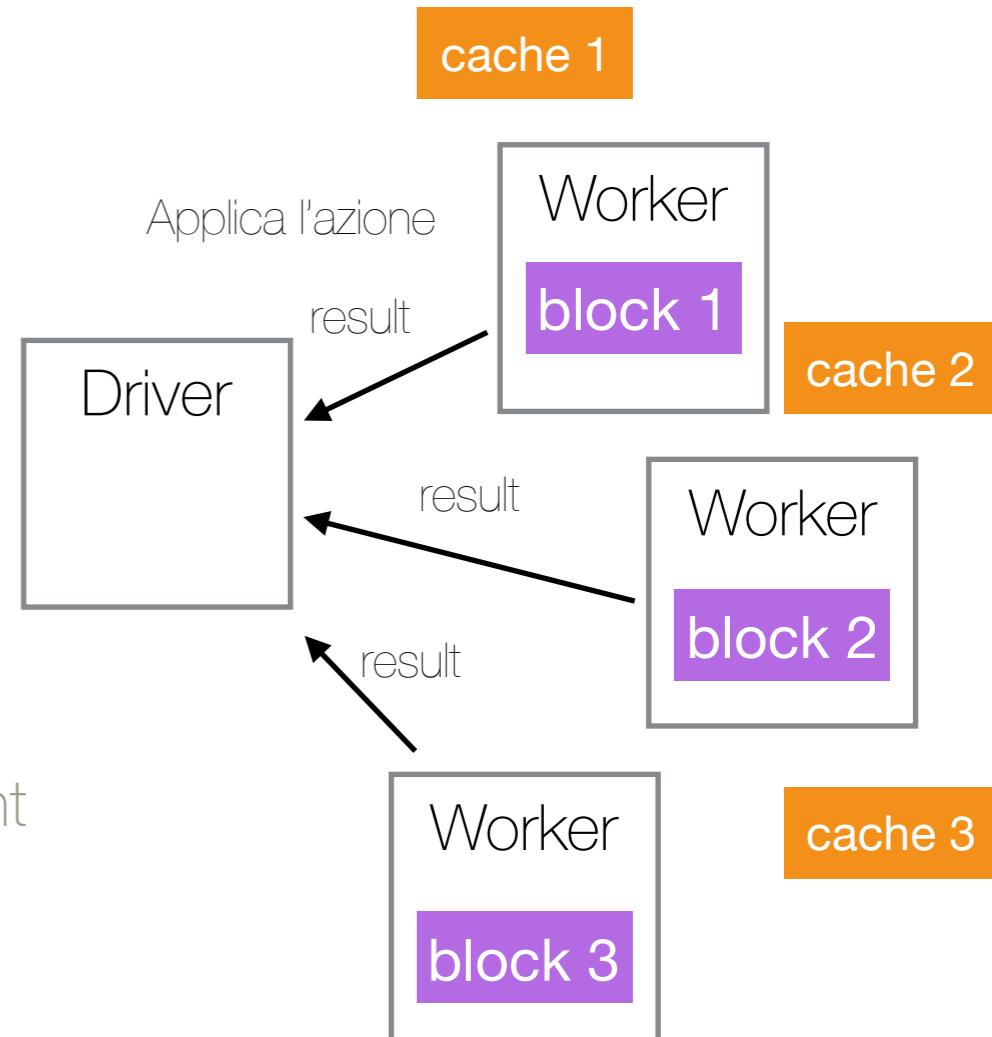
Contiamo le righe che contengono il termine **mongodb**

```
// leggiamo le righe dal file di log
val lines = sc.textFile("hdfs://.../var/log/*")

// applichiamo un filtro ed una trasformazione
val errors = lines.filter(_.startsWith("ERROR"))
val messages = errors.map(_.split("\t")).map(r => r(1))

// mettiamo in cache i messaggi
messages.cache()

// contiamo le righe di errore per MongoDB con l'azione count
messages.filter(_.contains("mongodb")).count()
```



Spark - Esempio



Contiamo le righe che contengono il termine mysql

```
// leggiamo le righe dal file di log
```

```
val lines = sc.textFile("hdfs://.../var/log/*")
```

```
// applichiamo un filtro ed una trasformazione
```

```
val errors = lines.filter(_.startsWith("ERROR"))
```

```
val messages = errors.map(_.split("\t")).map(r => r(1))
```

```
// mettiamo in cache i messaggi
```

```
messages.cache()
```

```
// contiamo le righe di errore per MongoDB con l'azione count
```

```
messages.filter(_.contains("mongodb")).count()
```

```
// contiamo le righe di errore per Mysql con l'azione count
```

```
messages.filter(_.contains("mysql")).count()
```

Applica l'azione

Driver

cache 1

Worker
block 1

cache 2

Worker
block 2

cache 3

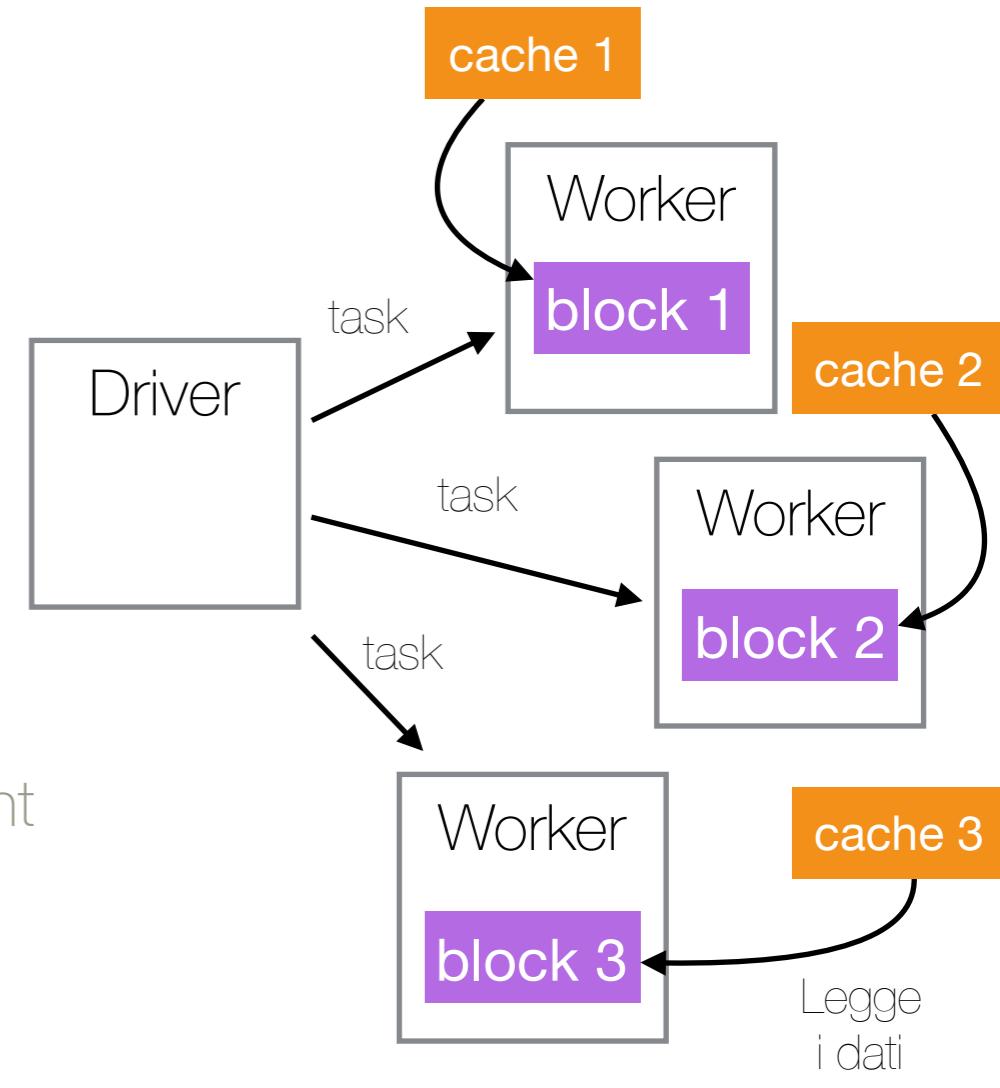
Worker
block 3

Spark - Esempio



Contiamo le righe che contengono il termine mysql

```
// leggiamo le righe dal file di log  
val lines = sc.textFile("hdfs://.../var/log/*")  
  
// applichiamo un filtro ed una trasformazione  
val errors = lines.filter(_.startsWith("ERROR"))  
val messages = errors.map(_.split("\t")).map(r => r(1))  
  
// mettiamo in cache i messaggi  
messages.cache()  
  
// contiamo le righe di errore per MongoDB con l'azione count  
messages.filter(_.contains("mongodb")).count()  
  
// contiamo le righe di errore per Mysql con l'azione count  
messages.filter(_.contains("mysql")).count()
```

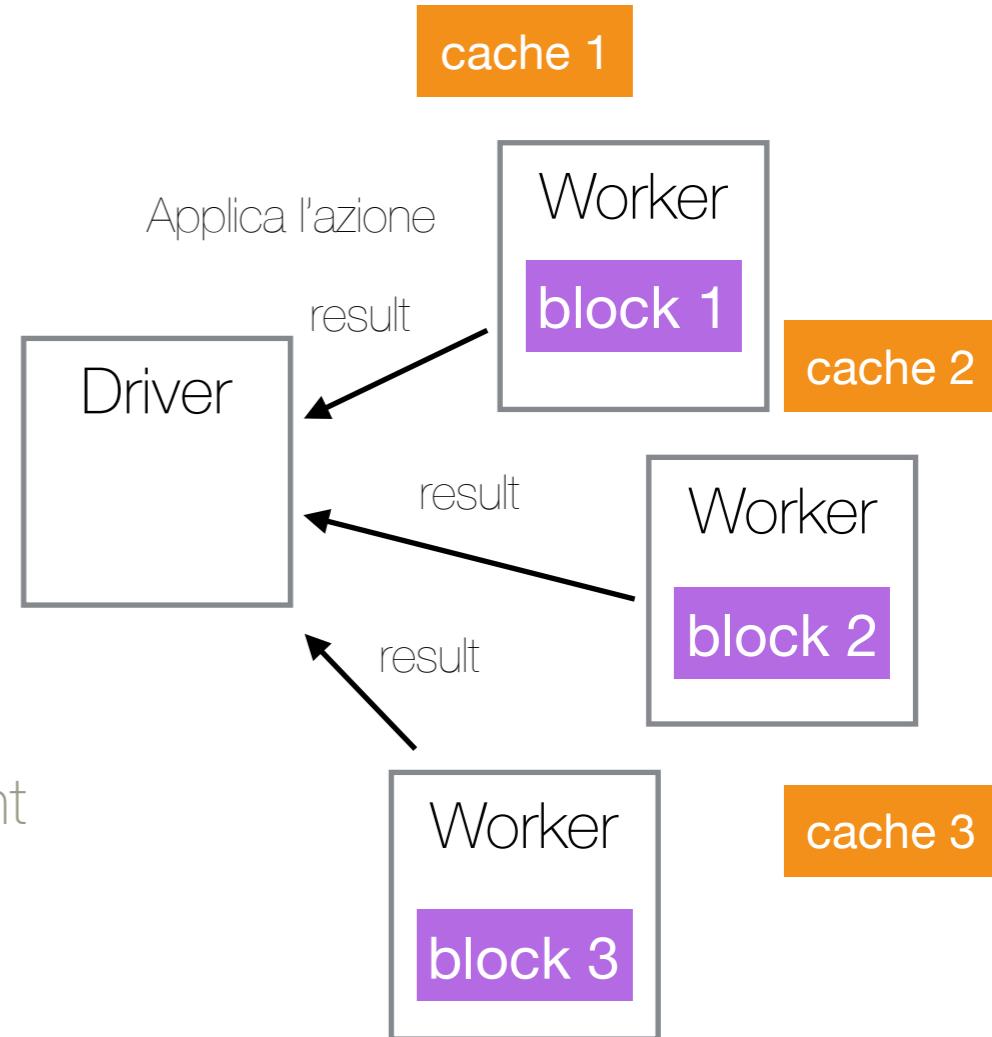


Spark - Esempio



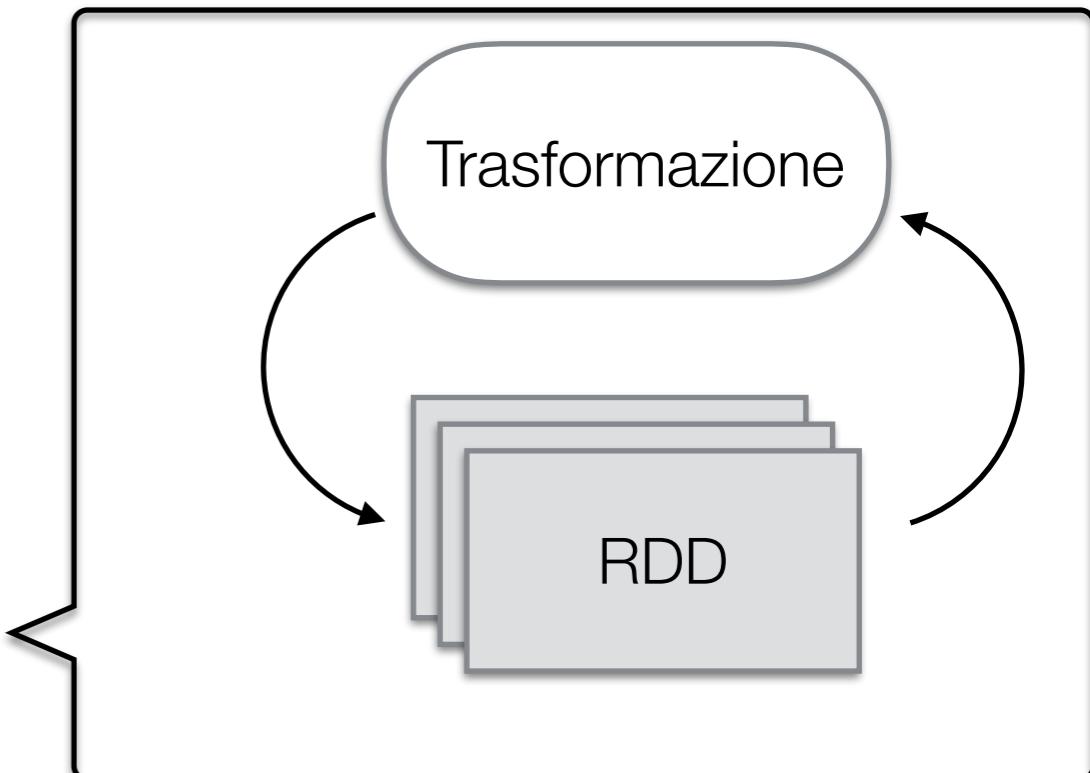
Contiamo le righe che contengono il termine mysql

```
// leggiamo le righe dal file di log  
val lines = sc.textFile("hdfs://.../var/log/*")  
  
// applichiamo un filtro ed una trasformazione  
val errors = lines.filter(_.startsWith("ERROR"))  
val messages = errors.map(_.split("\t")).map(r => r(1))  
  
// mettiamo in cache i messaggi  
messages.cache()  
  
// contiamo le righe di errore per MongoDB con l'azione count  
messages.filter(_.contains("mongodb")).count()  
  
// contiamo le righe di errore per Mysql con l'azione count  
messages.filter(_.contains("mysql")).count()
```



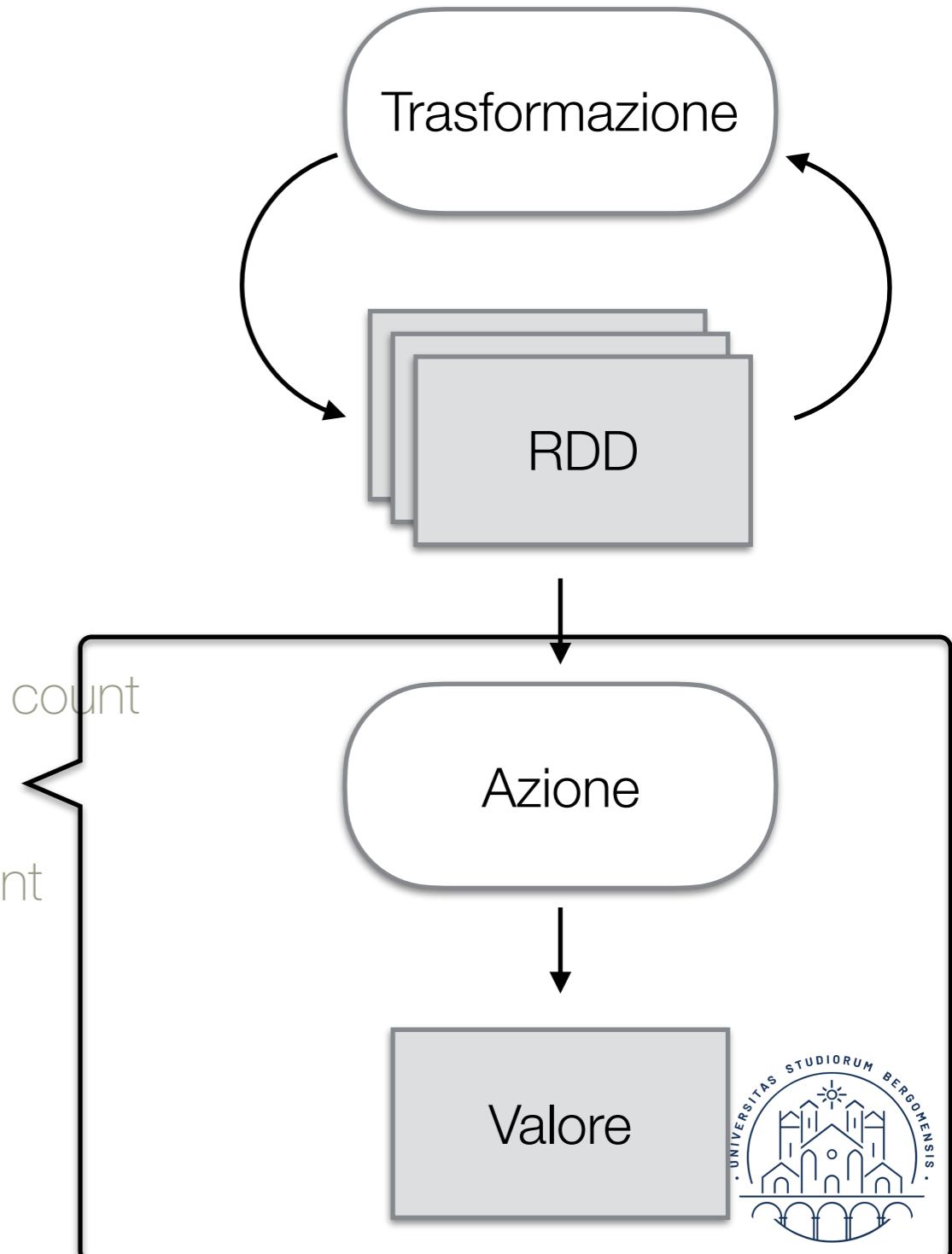
RDD - Resilient Distributed Datasets

```
// leggiamo le righe dal file di log  
val lines = sc.textFile("hdfs://.../var/log/*")  
  
// applichiamo un filtro ed una trasformazione  
val errors = lines.filter(_.startsWith("ERROR"))  
val messages = errors.map(_.split("\t")).map(r => r(1))  
  
// mettiamo in cache i messaggi  
messages.cache()  
  
// contiamo le righe di errore per MongoDB con l'azione count  
messages.filter(_.contains("mongodb")).count()  
  
// contiamo le righe di errore per Mysql con l'azione count  
messages.filter(_.contains("mysql")).count()
```



RDD - Resilient Distributed Datasets

```
// leggiamo le righe dal file di log  
val lines = sc.textFile("hdfs://.../var/log/*")  
  
// applichiamo un filtro ed una trasformazione  
val errors = lines.filter(_.startsWith("ERROR"))  
val messages = errors.map(_.split("\t")).map(r => r(1))  
  
// mettiamo in cache i messaggi  
messages.cache()  
  
// contiamo le righe di errore per MongoDB con l'azione count  
messages.filter(_.contains("mongodb")).count()  
  
// contiamo le righe di errore per Mysql con l'azione count  
messages.filter(_.contains("mysql")).count()
```



Spark - SparkContext



// dalla spark shell verifichiam

// la presenza dello spark context

scala>sc

```
res0: org.apache.spark.SparkContext = org.apache.spark.SparkContext@23b087c3
```

- Ogni applicazione deve inizializzare uno Spark Context
 - E' il punto di accesso per le API di Spark
 - Contiene tutte le informazioni e la configurazione
 - Accesso al cluster
 - Parametri ed informazioni per la connessione con il data warehouse (Mysql, MongoDB, Elasticsearch, Solr,)
 - Parametri dell'applicazione
 - Parametri per Spark

Spark - Cluster Manager



sc.master → determina quale cluster utilizzare

// verifichiamo come è impostato il master

scala>sc.master

res1: String = local[*]

Master URL	Meaning
local	Run Spark locally with one worker thread (i.e. no parallelism at all).
local[K]	Run Spark locally with K worker threads (ideally, set this to the number of cores on your machine).
local[*)	Run Spark locally with as many worker threads as logical cores on your machine.
spark://HOST:PORT	Connect to the given Spark standalone cluster master. The port must be whichever one your master is configured to use, which is 7077 by default.
mesos://HOST:PORT	Connect to the given Mesos cluster. The port must be whichever one your is configured to use, which is 5050 by default. Or, for a Mesos cluster using ZooKeeper, use mesos://zk://.... To submit with --deploy-mode cluster, the HOST:PORT should be configured to connect to the MesosClusterDispatcher .
yarn	Connect to a YARN cluster in client or cluster mode depending on the value of --deploy-mode. The cluster location will be found based on the HADOOP_CONF_DIR or YARN_CONF_DIR variable.

RDD - Resilient Distributed Datasets

Un RDD è una raccolta di elementi che possono essere utilizzati in parallelo

Resilient: possono essere ricreati

Distributed: vengono elaborati in tutto il cluster

Dataset: contengono dati provenienti da un file, da un data warehouse, ...

Un RDD è una collezione **distribuita immutabile** di oggetti: gli elementi di un RDD sono partizionati su un insieme di macchine.

Se una partizione viene persa può sempre essere ricostruita.

Spark - Dataframe



- Limiti e problemi degli RDD
 - Non possono essere ottimizzati da Spark
 - Risultano lenti in Python e Java
 - E' molto facile costruire una catena di trasformazioni non efficienti
- Dataframe e Dataset superano questi problemi
- Le API DataFrame forniscono un livello di astrazione superiore: possiamo manipolare dati attraverso un linguaggio per query, anche SQL
 - Le query su DataFrame sono ottimizzate da Spark
 - Non sono type safety

Spark - Dataset

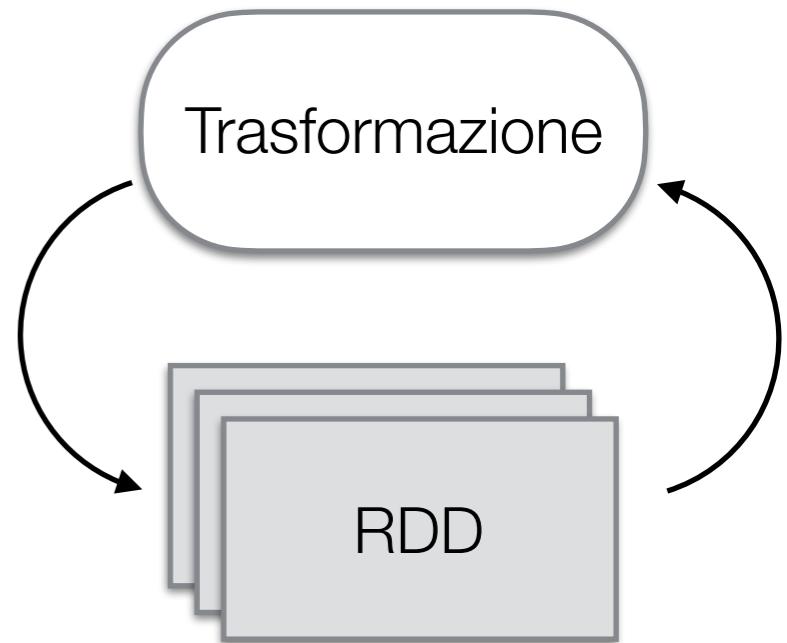
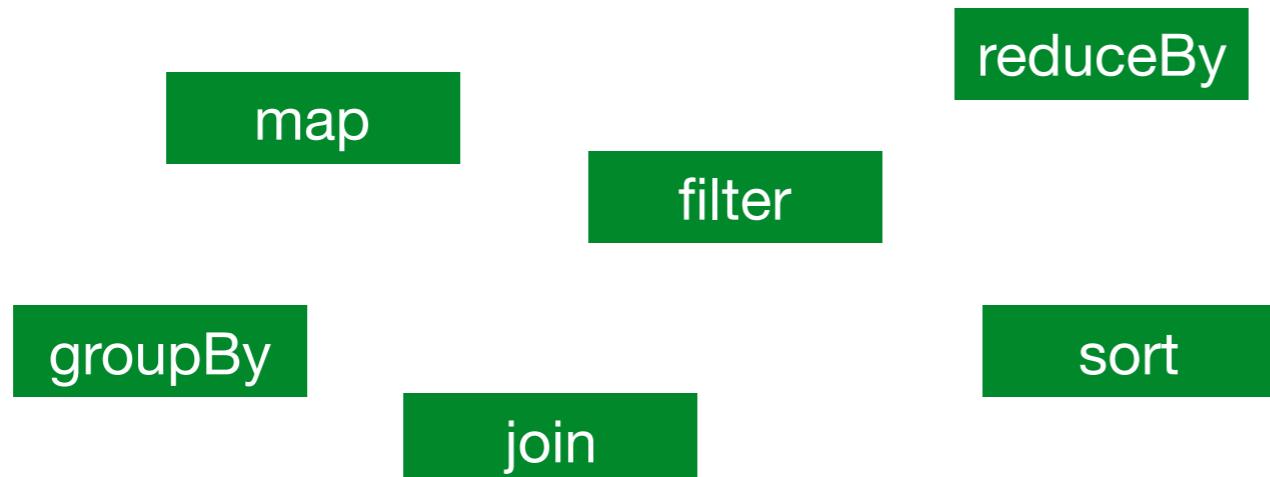


- **Dataset**
 - Simili concettualmente a un RDD
 - Usano il framework Tungsten (encoding in-memory): migliori performance in termini di spazio utilizzato e velocità
 - Più facili dei Dataframe
 - Disponibili con Spark 2.*

Spark - Trasformazioni



Le **trasformazioni** creano un nuovo dataset da uno esistente



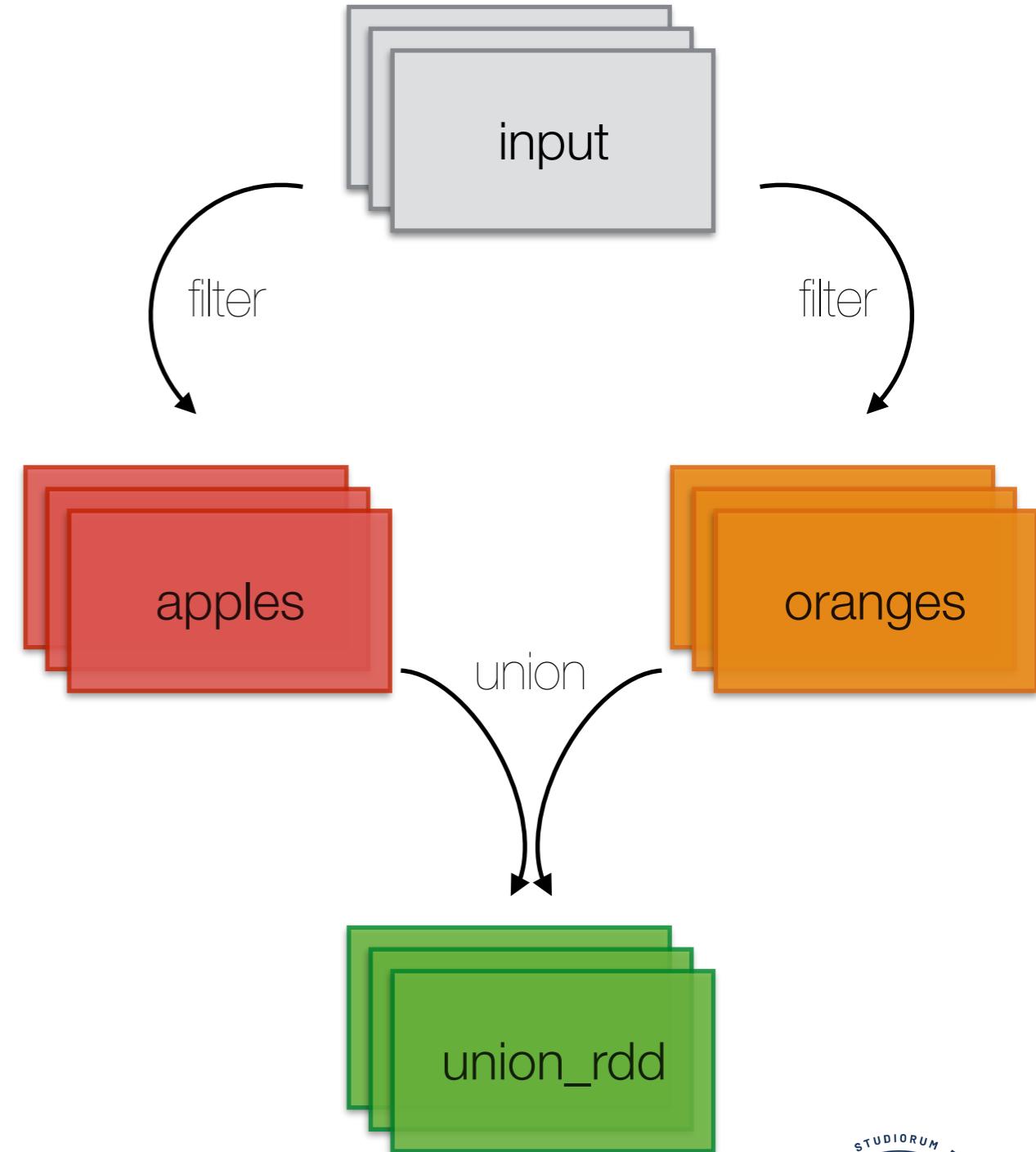
Tutte le trasformazioni sono **lazy**!!!

Non elaborano subito i loro risultati, in questo modo ottimizzano i calcoli e se necessario possono recuperare le partizioni di dati perse

Spark - Trasformazioni



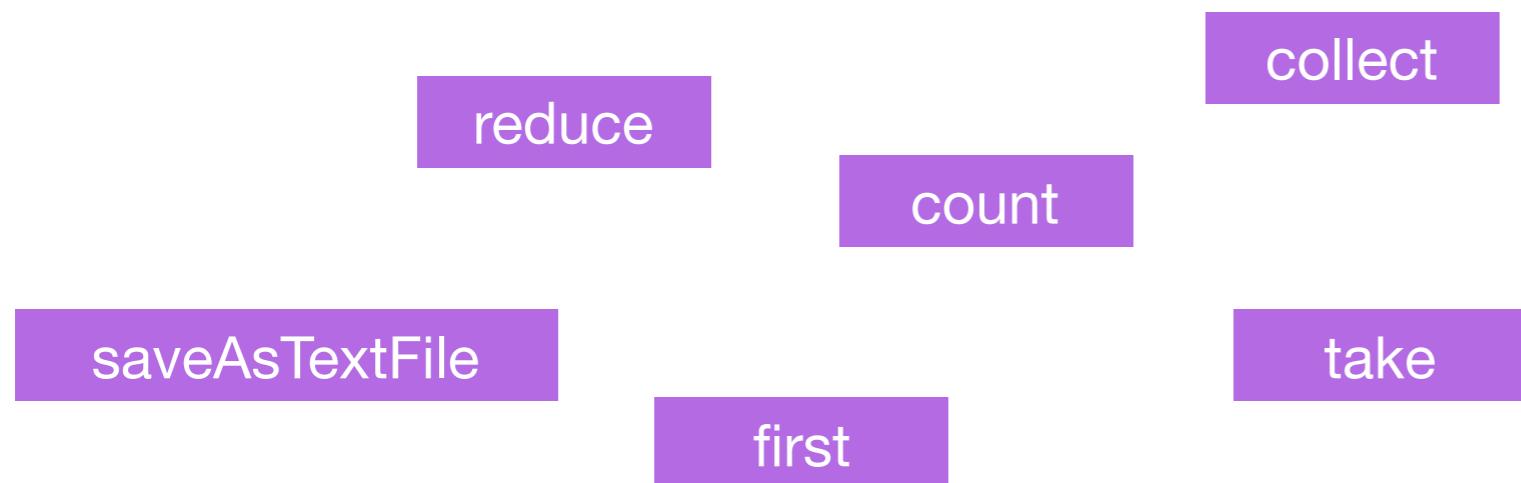
```
// leggiamo le righe dal file di log  
val input = sc.textFile("myfile.txt")  
  
// cerco le linee che contengono "apple"  
val apples = input.filter(x => x.contains("apple"))  
  
// cerco le linee che contengono "orange"  
val oranges = input.filter(x => x.contains("orange"))  
  
// unisco i due dataset  
val union_rdd = apples.union(oranges)
```



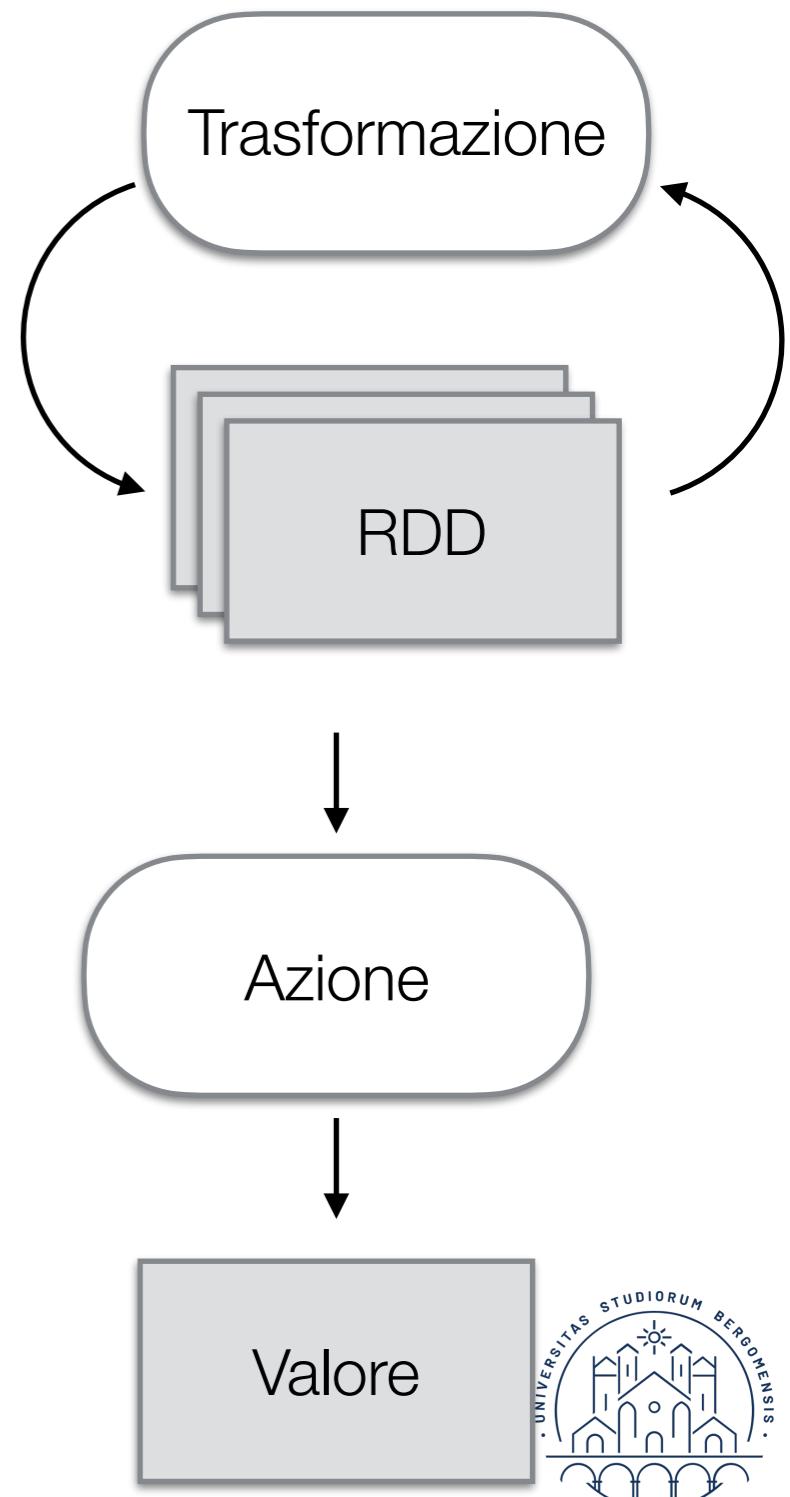
Spark - Azioni



Un **azione** indica che qualcosa deve essere effettuato



Tutte le azioni vengono eseguite subito da Spark e producono un risultato.



Spark - Lazy Evaluation



Tutte le trasformazioni sono **pigre**... perché?



Spark - Lazy Evaluation



- Big Picture: se guardiamo dall'alto possiamo vedere più alternative
- Non eseguiamo operazioni non necessarie

Esempio:

```
val input = sc.textFile("myfile.txt")
val apples = input.filter(x => x.contains("apple"))
val oranges = input.filter(x => x.contains("orange"))
val union_rdd = apples.union(oranges)
// prendo le prime 5 righe
val first5rows = union_rdd.take(5)
```

union_rdd.count()



Dopo la union dei due dataset prendo solo le prime 5 righe:

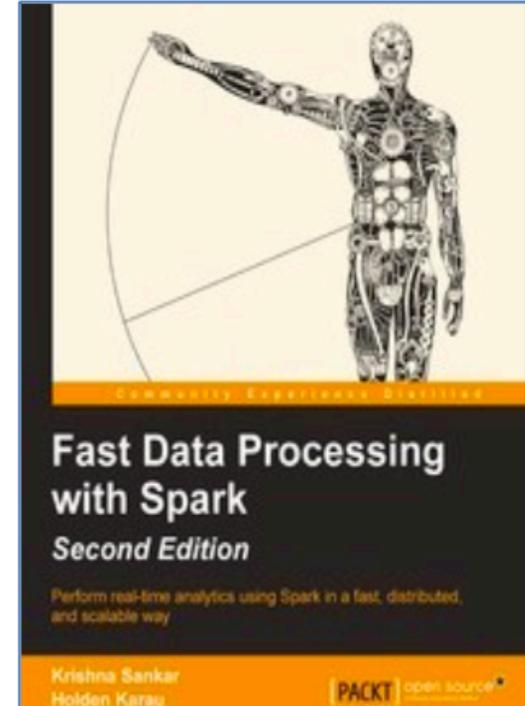
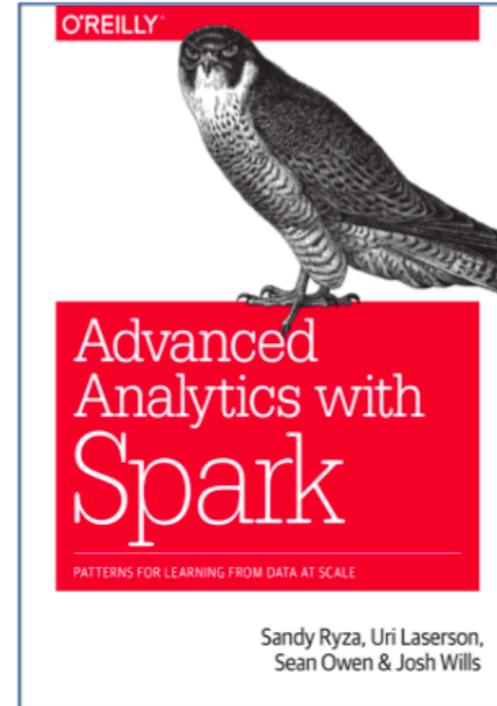
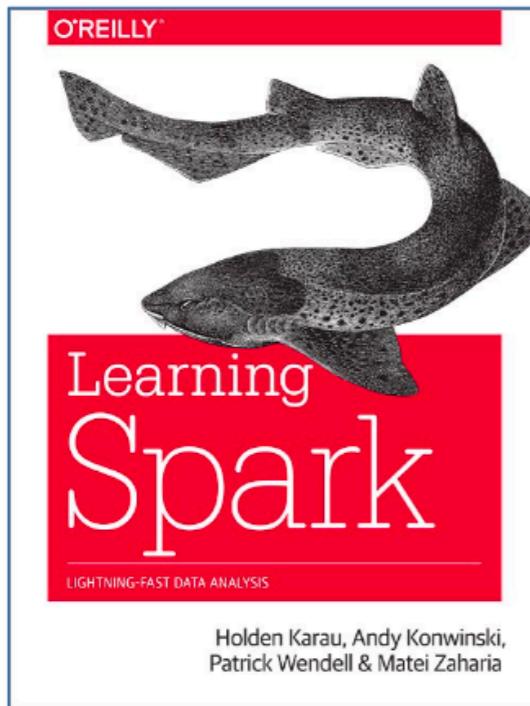
- se Spark fosse **eager** farebbe la union dei due dataset, trasformerebbe tutto in una collezione di elementi e prenderebbe le prime 5 righe
- ma Spark è **pigro!** Non calcola l'intera unione, ma prepara solo i primi 5 elementi della collezione.



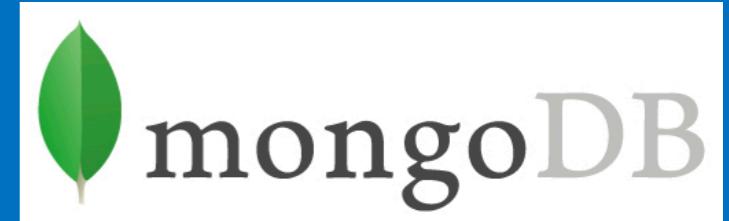
Resources



<https://spark.apache.org>



MongoDB



Quale è la differenza fra database e DBMS?

uhmm...

C'è una differenza?

Quali sono le proprietà di un DBMS relazionale?

Caratteristiche di NoSql



- Non usano SQL (Hbase, Cassandra, Redis,...)
- Sono generalmente progetti open source
- Gran parte di essi lavorano su cluster
- I RDBMS usano le transazioni ACID per gestire la consistenza, NoSQL utilizzano altri metodi
- Sono schema free o schema less
- NoSQL è un movimento piuttosto che una tecnologia
- I RDBMS non sono destinati a scomparire, ma la novità è che ora esistono altre possibilità

CAP Theorem



«Of three properties of shared data systems, data **Consistency**, system **Availability** and tolerance to network **Partitions**, only two can be achieved at any given moment in time»

University of California at Berkeley, Eric Brewer

- CA: relazione tradizionale
- AP: le richieste vengono eseguite su ogni nodo anche se violano la consistenza (Cassandra, DynamoDB)
- CP: le richieste vengono eseguite solo sui nodi che garantiscono la consistenza (Hbase, MongoDB, Google BigTables)

In pratica...



If you cannot limit the number of faults and requests can be directed to any server and you insist on serving every request you receive then you cannot possibly be consistent

CAP Theorem

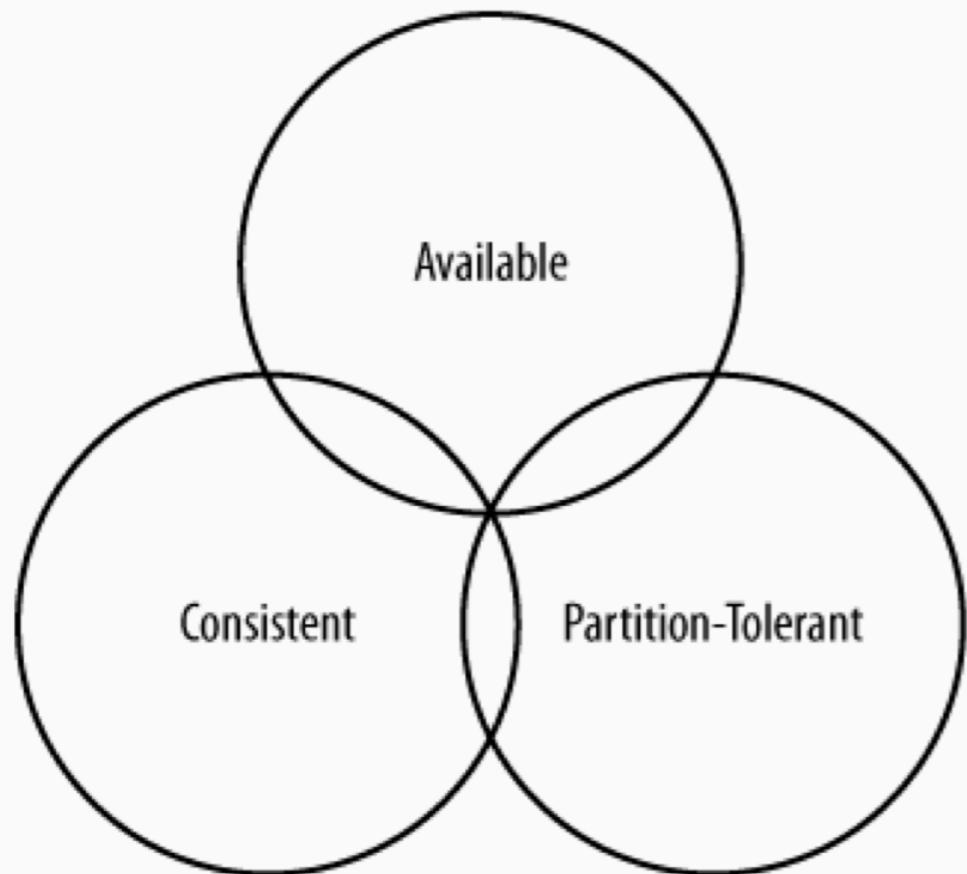


Figure 2-1. CAP Theorem indicates that you can realize only two of these properties at once

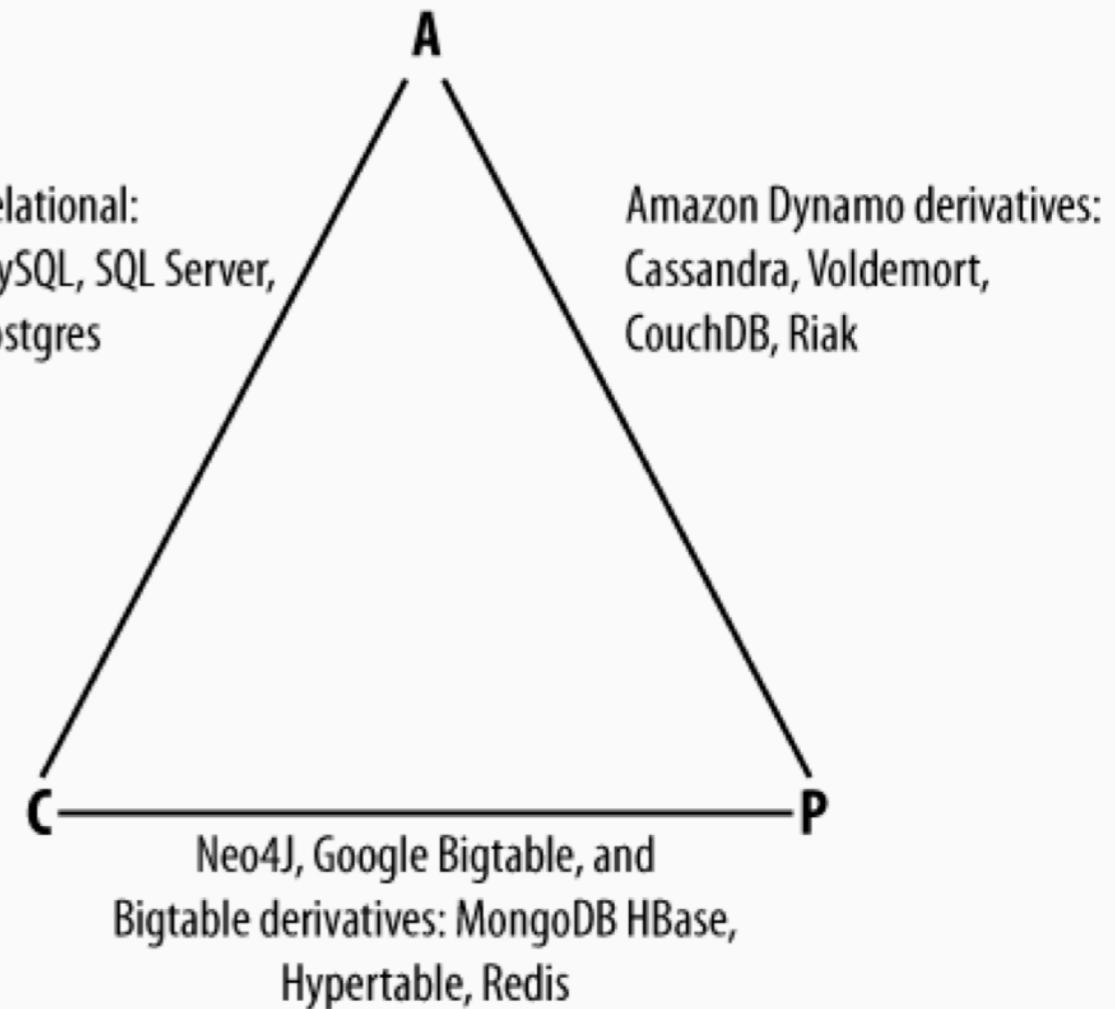
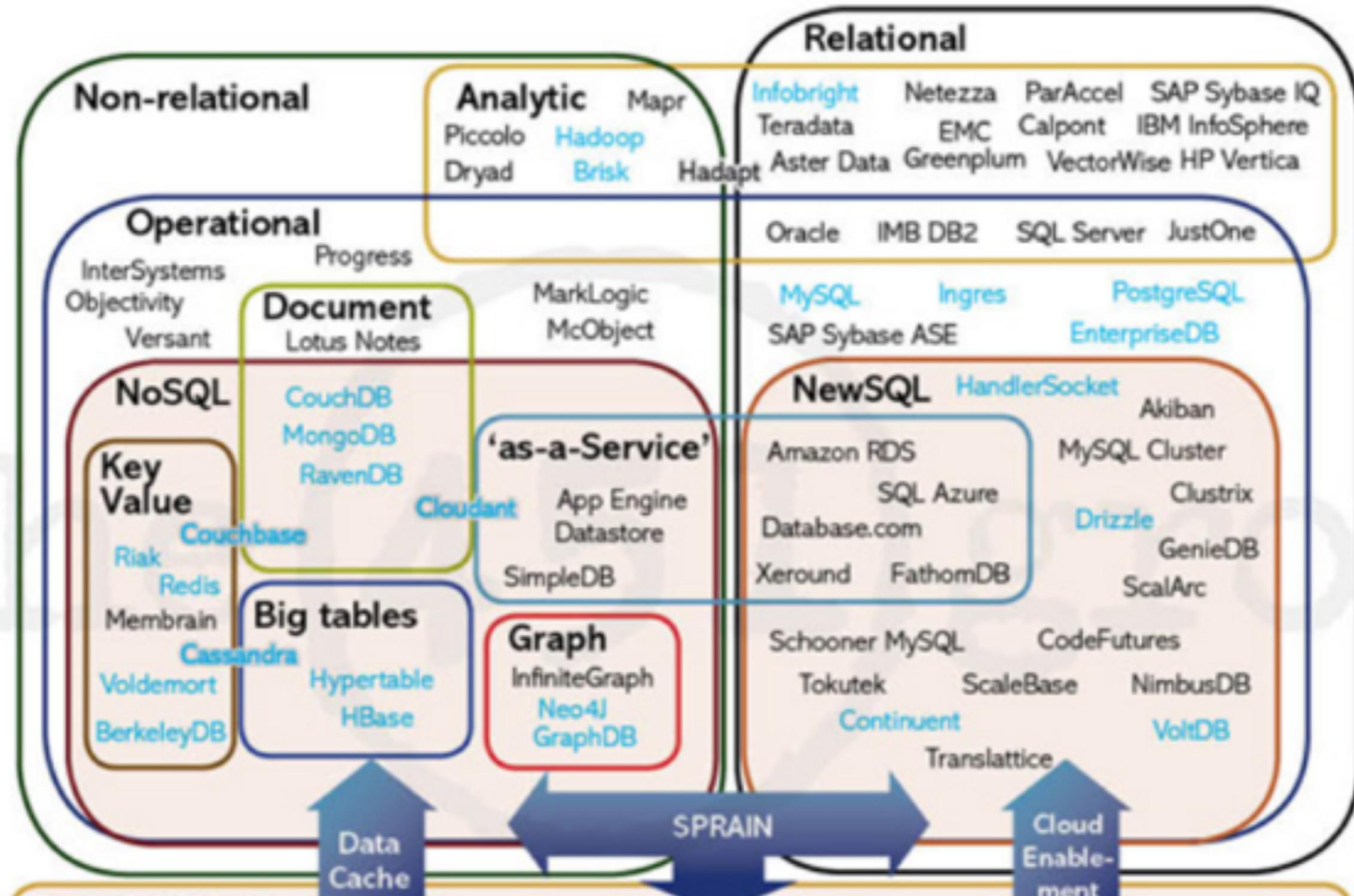


Figure 2-2. Where different databases appear on the CAP continuum

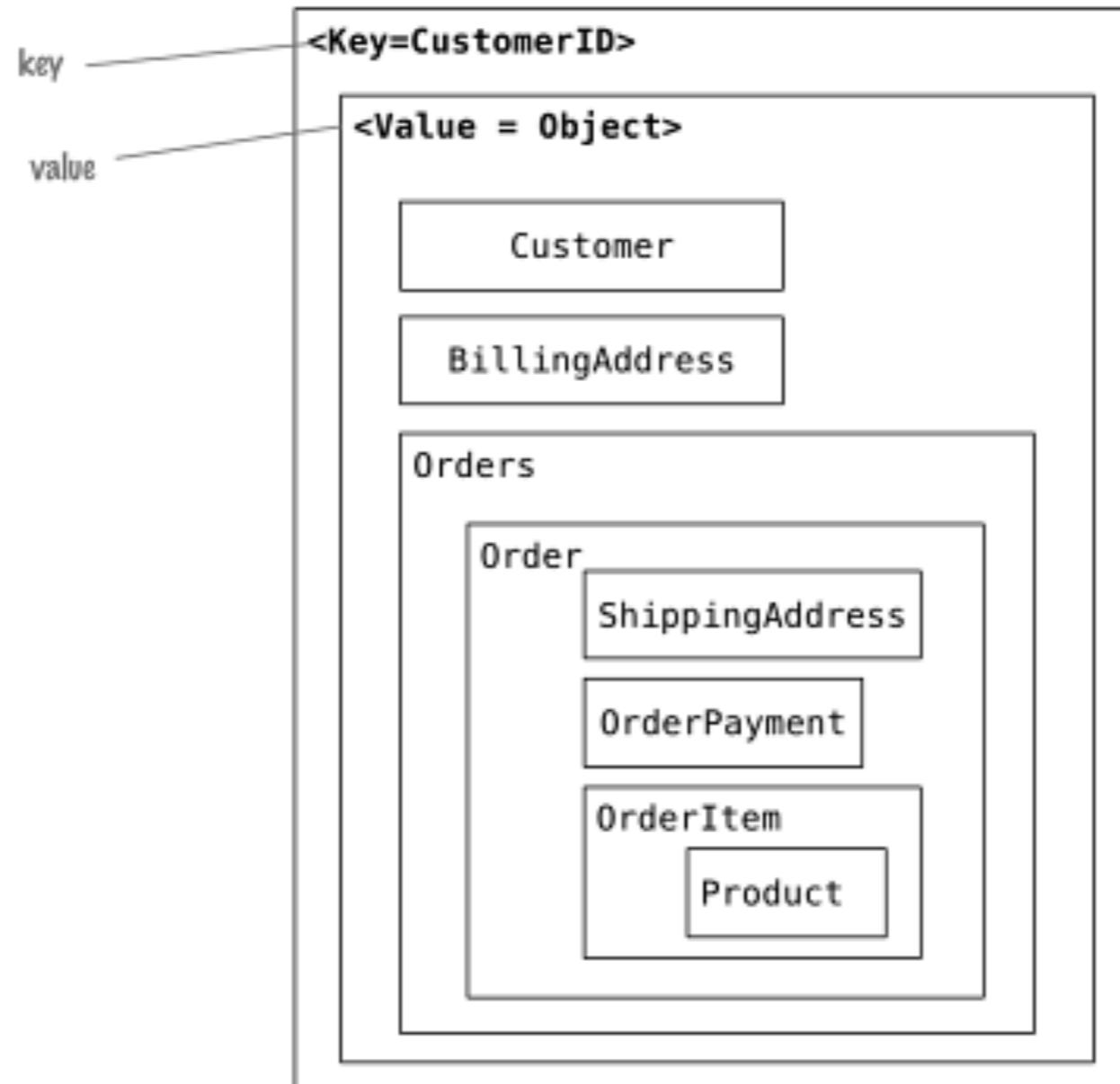


KeyValue Stores



- Key values stores
 - i dati vengono immagazzinati in un elemento che contiene una chiave assieme ai dati veri e propri
 - più semplice da implementare
 - più inefficiente se la maggior parte delle operazioni riguardano soltanto una parte di un elemento
 - Una chiave, un valore, nessun duplicato
 - Il valore è un oggetto binario chiamato «blob»: il database non sa cosa sia e non ha bisogno di saperlo

KeyValue Stores



Redis



WEBSITE

<http://redis.io/>

DATABASE MODEL

key-value store

INITIAL RELEASE

2009



LICENSE

Open source, in-memory data structure store, used as a database, cache and message broker.

DESCRIPTION



flickr GitHub

 **stackoverflow**

Voldemort



WEBSITE
[voldemort.com/](http://www.project-voldemort.com/)

DATABASE MODEL

[http://www.project-](http://www.project-voldemort.com/)

key-value store

INITIAL RELEASE

2009

LICENSE

It is basically just a big, distributed, persistent, fault-tolerant hash table

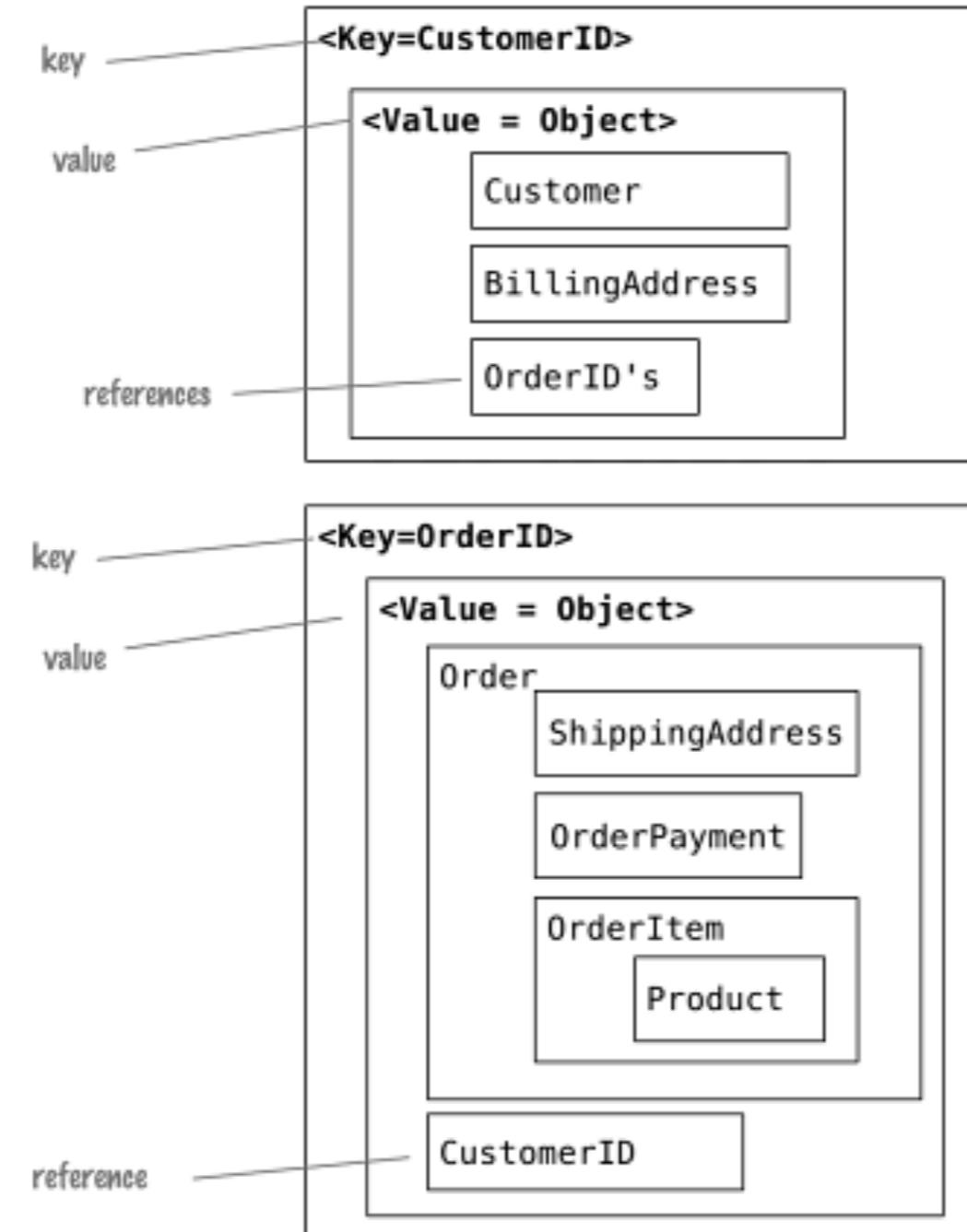
DESCRIPTION



- Document stores
 - è l'evoluzione del metodo key/value
 - il valore è di solito strutturato e compreso dal DB
 - è possibile interrogare i dati (in modi diversi dall'uso della semplice chiave)

Rispetto ai normali database relazionali invece che immagazzinare i dati in tabelle con dei campi fissi, questi vengono messi in un documento che può contenere illimitati campi di illimitata lunghezza, così se ad esempio di una persona conosciamo solo nome e cognome, ma magari di un'altra persona anche indirizzo, data di nascita e codice fiscale, si evita che per il primo nominativo ci siano campi inutilizzati che occupano inutilmente spazio

Document stores



MongoDB



WEBSITE

<https://www.mongodb.com/>
document store

INITIAL RELEASE

2009

LICENSE



DESCRIPTION

MongoDB is a document database with the scalability and flexibility that you want with the querying and indexing that you need



DynamoDB



WEBSITE	https://aws.amazon.com/dynamodb/
DATABASE MODEL	document store
INITIAL RELEASE	2012
LICENSE	commercial
DESCRIPTION	È un database cloud interamente gestito che supporta sia i modelli di storage document store sia quelli di tipo chiave-valore.
WHO USES	

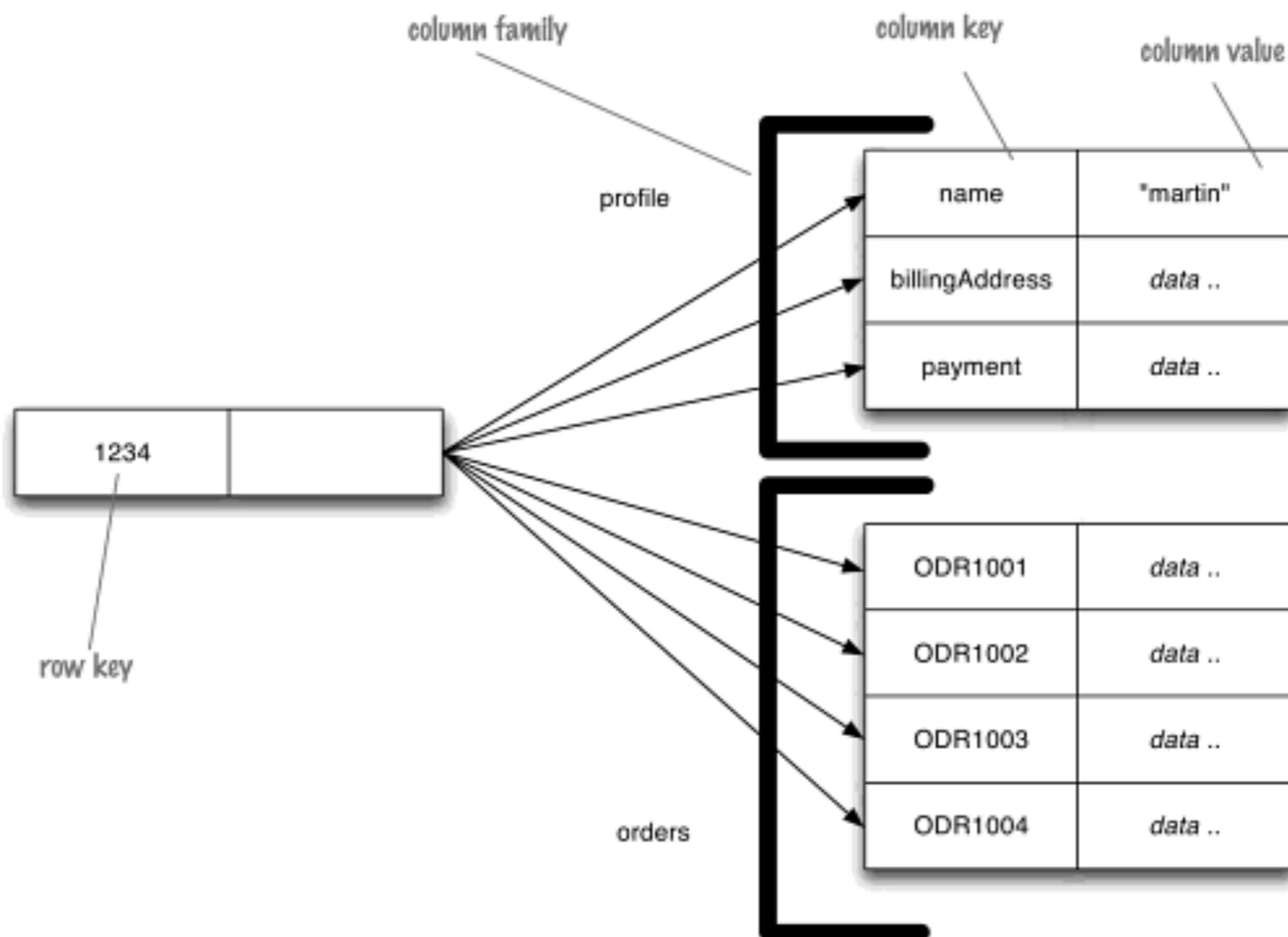


Column family stores



- Google's BigTable
 - Il nome deriva dalla sua struttura composta da colonne sparse e nessuno schema
- Non dobbiamo pensare alla sua struttura come una tabella ma come una mappa a due livelli
- Molti database utilizzano le righe come unità di storage
- Tuttavia in molti casi è necessario accedere a poche colonne di molte righe

Column family stores



Cassandra



WEBSITE

<http://cassandra.apache.org/>

DATABASE MODEL

wide column store

INITIAL RELEASE

2008



LICENSE

Wide-column store based on ideas of
BigTable and DynamoDB

DESCRIPTION



ebay

Instagram

GitHub

NETFLIX

Google Big Table



WEBSITE

<https://cloud.google.com/bigtable/>

DATABASE MODEL wide column store

INITIAL RELEASE 2015

LICENSE commercial

DESCRIPTION Google's NoSQL Big Data database service. It's the same database that powers many core Google services, including Search, Analytics, Maps, and Gmail.

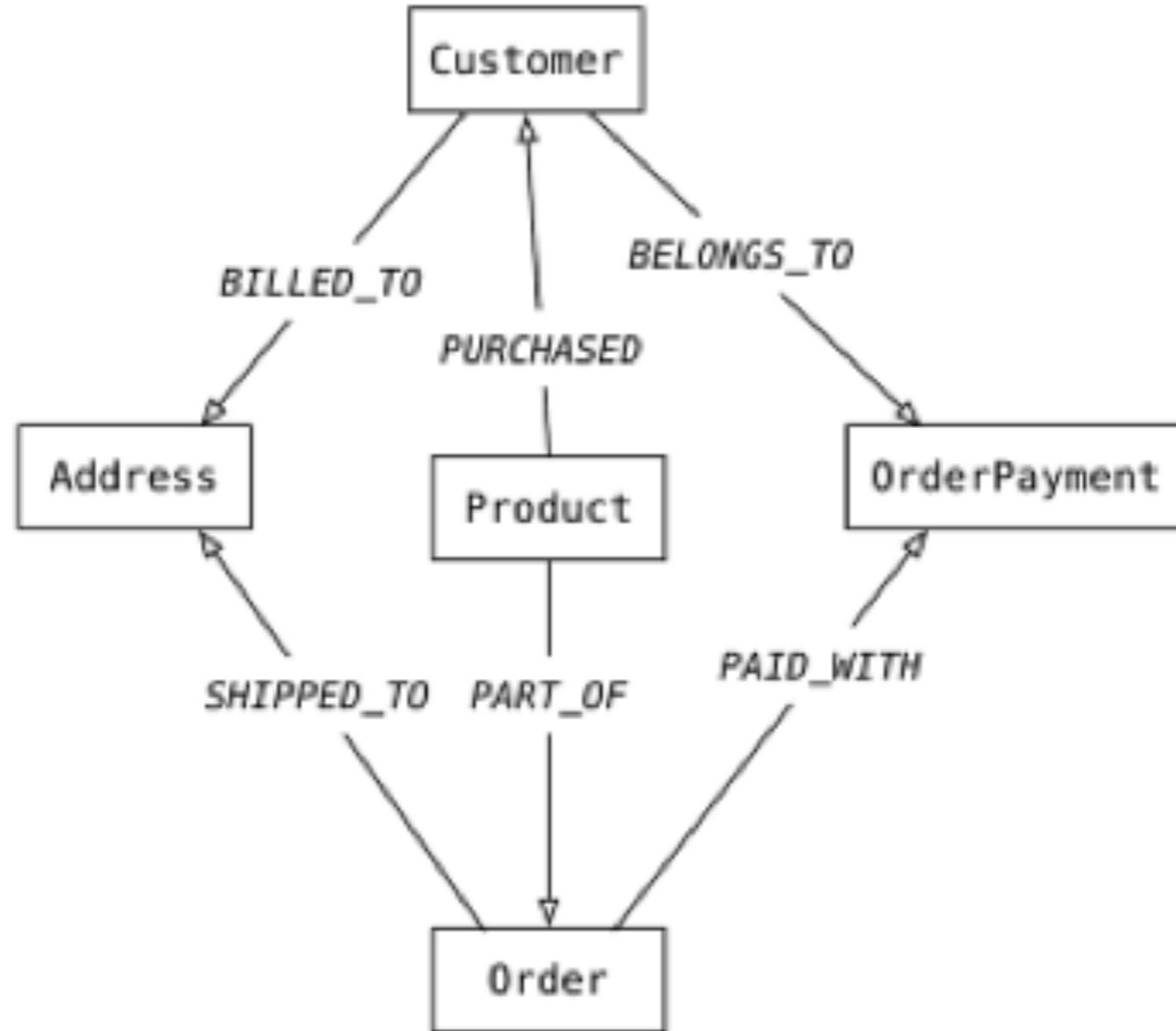


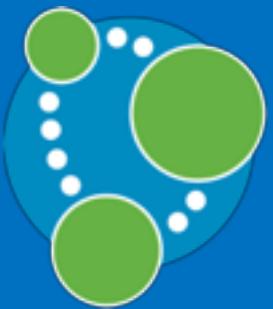
Graph databases



- Graph databases
 - i dati vengono immagazzinati sotto forma di strutture a grafi, rendendo più performante l'accesso a questi da applicativi orientati agli oggetti

Graph databases





WEBSITE	https://neo4j.com/
DATABASE MODEL	graph DBMS
INITIAL RELEASE	2007
LICENSE	
DESCRIPTION	Highly scalable native graph database that leverages data relationships as first-class entities



Time series



- Un Time series DBMS è un sistema ottimizzato per la gestione dei dati delle serie temporali: ogni voce è associata a un timestamp
- I dati delle serie temporali possono essere prodotti da sensori, smart meter o RFID nel mondo IoT o possono rappresentare i titoli di borsa di un sistema di negoziazione di titoli ad alta frequenza
- Sono progettati per raccogliere, archiviare e interrogare in modo efficiente serie temporali con elevati volumi di transazioni

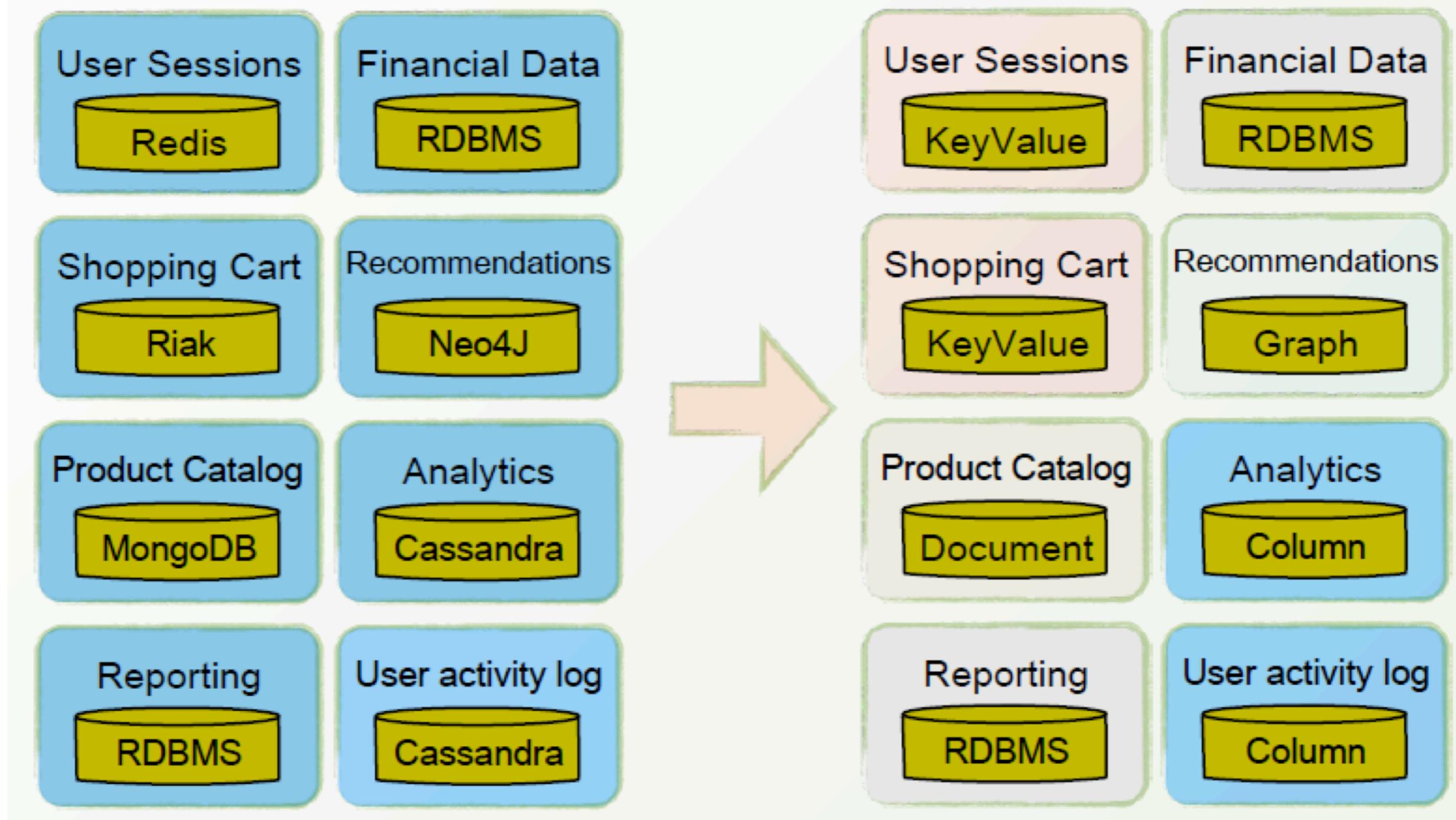


Quando usare un Document Store



- Quando i dati rappresentano collezioni di oggetti simili
- Semi strutturati o sparsi piuttosto che tabulari
- Quando i campi degli oggetti hanno valori multipli

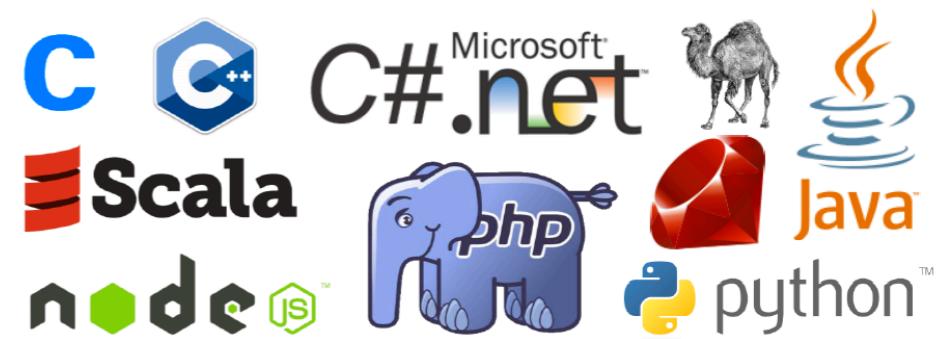
Polyglot persistence



Cosa è MongoDB?



- <https://www.mongodb.com/>
- Nasce nel 2007 (10gen)
 - La prima versione per produzione è pronta nel 2010
 - Ora siamo alla versione 4
- **Document Oriented NoSQL Database**
 - Hash-based, schema-less database
 - No Data Definition Language
 - Ogni documento è identificato da un identificativo (**_id**)
 - BSON format (Based on JSON –B stands for Binary)
- Le applicazioni gestiscono lo schema dei dati
- Scritto in C++
- Supporta API (driver) per molti diversi linguaggi
 - JavaScript, Python, Ruby, Perl, Java, JavaScala, C#, C++, Haskell, Erlang
 - Off-the-shelf drivers (Hadoop, BI tools, ...)
 - REST interface
- MultiOS (Windows, Linux, macOS)
- Cloud (AWS, Azure, GCP,...) + MongoDB Atlas



Funzionalità



- Dynamic schema
 - No DDL
- Document-based database
- Secondary indexes
- Query language via an API
- Atomic writes and fully-consistent reads
 - If system configured that way
- Master-slave replication with automated failover (replica sets)
- Built-in horizontal scaling via automated range-based partitioning of data (sharding)
- No joins nor transactions

Perché usare MongoDB?

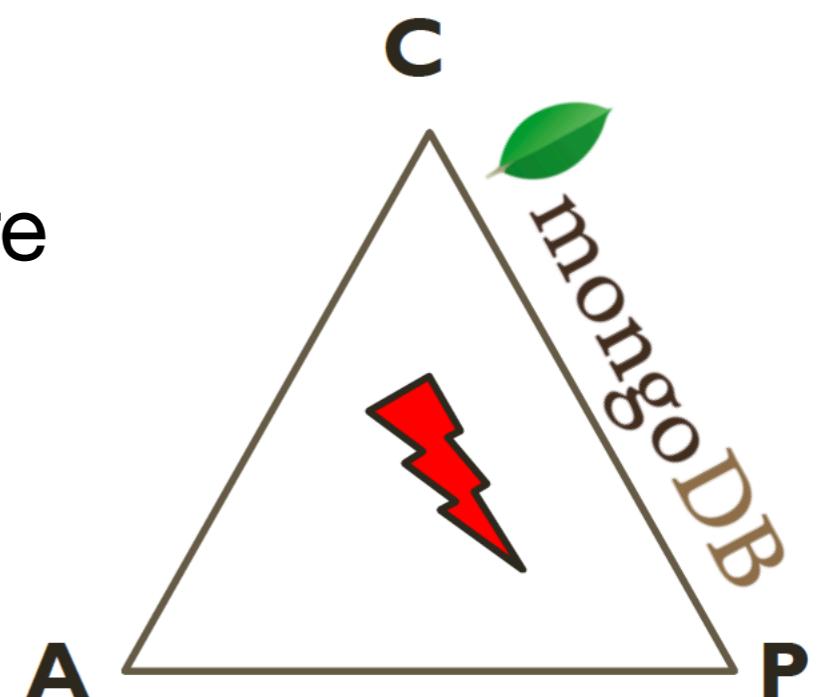


- Query molto semplici
 - Linguaggio di scripting molto potente
 - Scale out
- Integrazione dei dati veloci
- Non adatto per transazione complesse

MongoDB e CAP



- Consistenza
 - Tutte le repliche contengono la stessa versione dei dati
- Disponibilità
 - il sistema rimane attivo in caso di failure
- Partition tolerance
 - entry points multipli
 - il sistema rimane attivo in caso di split



Hierarchical Objects



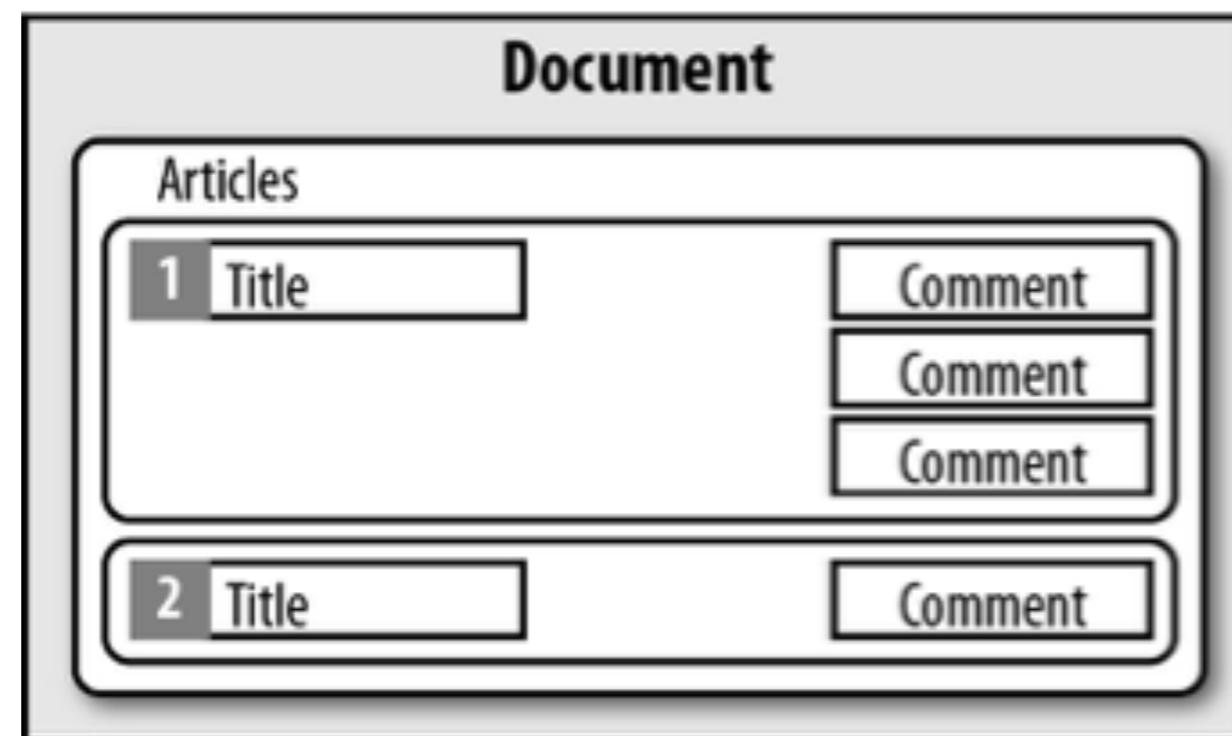
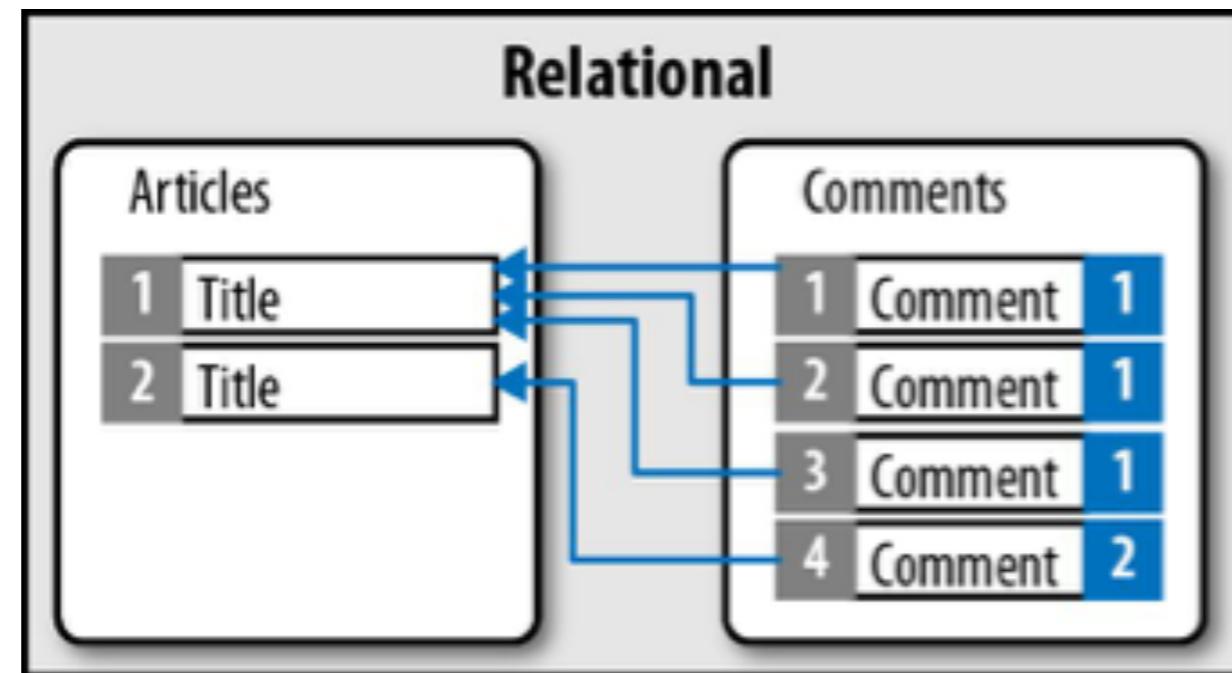
- A MongoDB instance may have zero or more ‘databases’
- A database may have zero or more ‘collections’
- A collection may have zero or more ‘documents’
- A document may have one or more ‘fields’
- MongoDB ‘Indexes’ function much like their RDBMS counterparts

RDB vs MongoDB

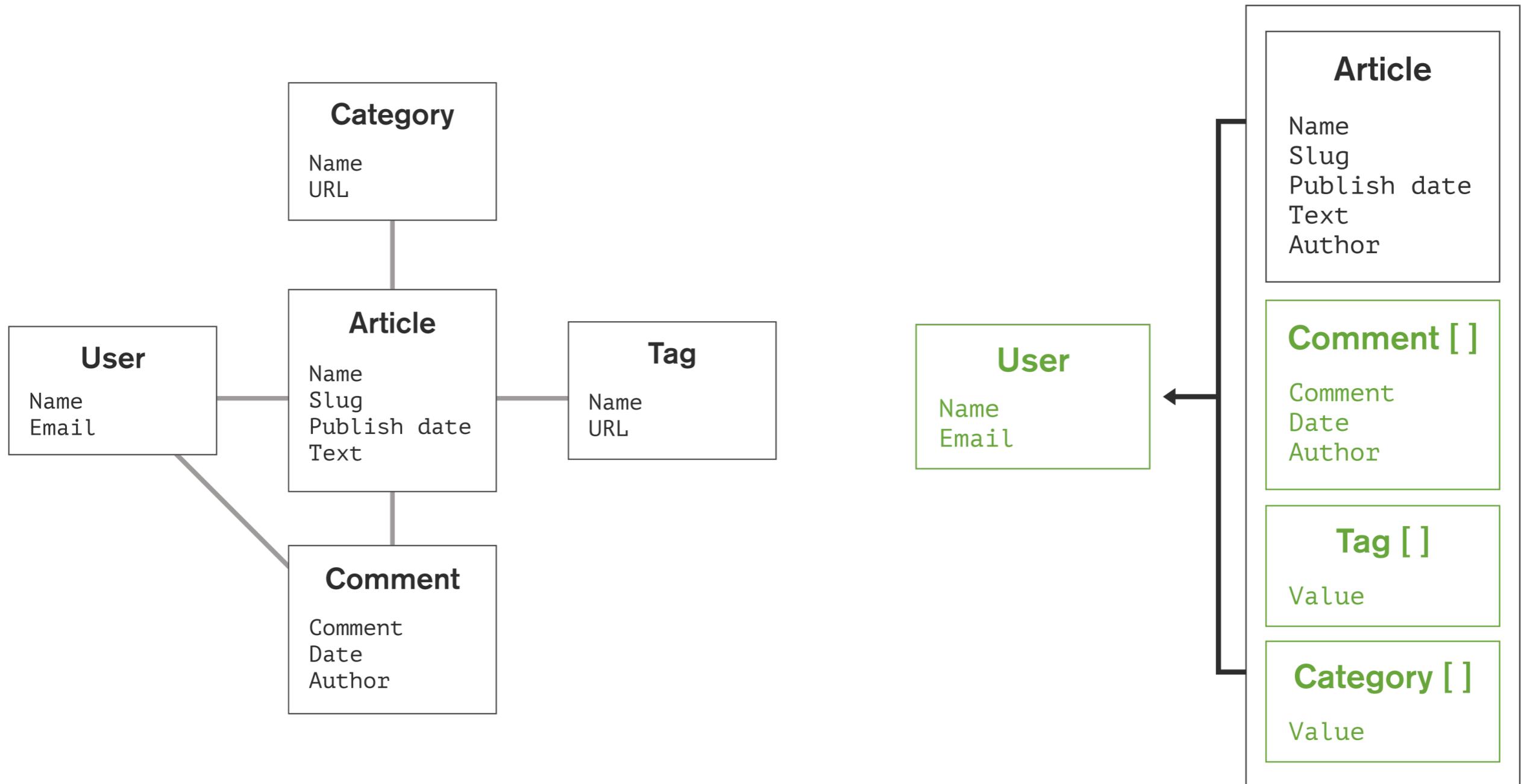


RDBMS	MongoDB
Database	Database
Table, view	Collection
Row	Document
Column	Field
Index	Index
Join	Embedded Document
Foreign key	Reference
Partition	Shard

RDB vs MongoDB



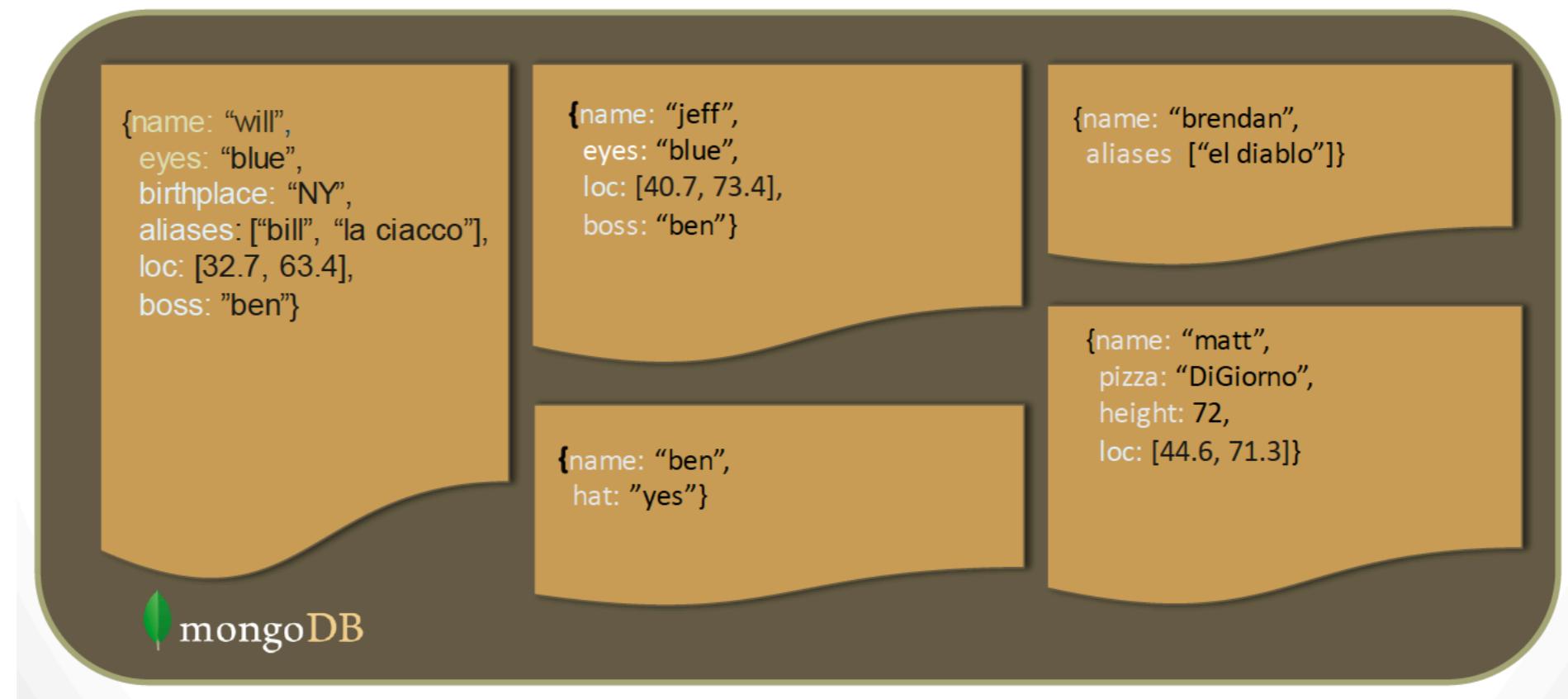
Data As Documents



What is a document?



- MongoDB, such as almost all NoSQL DBs, stores data in json format (BSON format)
- A single document has it's own json and any json has it's own structure. So multiple documents in a collection can have totally different structures (or schemas). This is schema-free paradigm.
- The maximum BSON document size is 16 MB. The maximum document size helps ensure that a single document cannot use excessive amount of RAM or, during transmission, excessive amount of bandwidth.



Mongo Document



{

```
name: 'Mauro Pelucchi',
```

```
address: {
```

```
street: 'via Finazzi 47',
```

```
city: 'Carvico'
```

}

}

Mongo Document



```
{
```

```
  _id: ObjectId("4efa8d2b7d284dad101e4bc9"),
```

```
  name: 'Mauro Pelucchi',
```

```
  address: {
```

```
    street: 'via Finazzi 47',
```

```
    city: 'Carvico'
```

```
}
```

```
}
```

Mongo Document



```
{
```

```
  _id: 12345,
```

```
  name: 'Mauro Pelucchi',
```

```
  address: {
```

```
    street: 'via Finazzi 47',
```

```
    city: 'Carvico'
```

```
}
```

```
}
```

Come posso inserire gli indirizzi di spedizione?

Mongo Document



{

```
_id: 12345,
```

```
name: 'Mauro Pelucchi',
```

```
address: {
```

```
    street: 'via Finazzi 47',
```

```
    city: 'Carvico'
```

},

```
shipAdresses: [
```

```
    {strett: 'via Finazzi 47', city: 'Carvico'},
```

```
    {strett: 'via Cavour 23', city: 'Carvico', note: 'suonare a Gianpietro'},
```

```
    {strett: 'Edificio U7 - Università degli Studi di Milano Bicocca', city: 'Milano'}
```

]

}

MongoDB Atlas



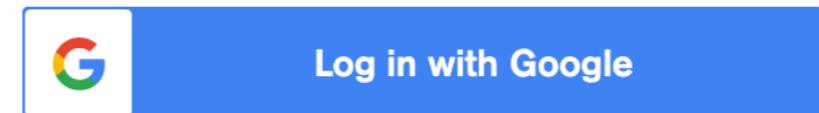
Database as a Service for MongoDB

- Scalable back-end for your application on-demand
- Secure by default
- Highly available, even while scaling
- Patch maintenance performed for you
- Your own MongoDB cluster in the cloud

<https://account.mongodb.com/account>



Log in to your account



or

Email Address i

A text input field for entering an email address, with a small "..." icon on the right side.

Next

Don't have an account? [Sign Up](#)

MongoDB Atlas



mongoDB Atlas All Clusters

Please set your time zone Usage This Month:\$0.00 details Preferences Mauro Pelucchi ▾

CONTEXT mauropelucchi.com ▾

ATLAS

Clusters

Find a cluster...

SANDBOX

MyCluster Version 4.2.5

CONNECT METRICS ...

CLUSTER TIER M0 Sandbox (General)

REGION AWS / Ireland (eu-west-1)

TYPE Replica Set - 3 nodes

LINKED STITCH APP None Linked

Operations R: 0 W: 0 100.0/s Last 6 Hours

Logical Size 3.5 MB 512.0 MB max 0.0 B Last 30 Days

Connections 0 500 max Last 6 Hours

Enhance Your Experience
For dedicated throughput, richer metrics and enterprise security options,
upgrade your cluster now!

Upgrade

Integrations

Settings

Services

Charts

Stitch

Triggers

Help

Docs

Support

System Status: All Good Last Login: 91.255.53.33
©2020 MongoDB, Inc. Status Terms Privacy Atlas Blog Contact Sales

A red arrow points to the 'METRICS' button on the MyCluster card.

MongoDB Atlas



mongoDB Atlas All Clusters

! Please set your time zone Usage This Month:\$0.00 [details](#)

[Preferences](#)

Mauro Pelucchi ▾

CONTEXT

mauropelucchi.com ▾

MAUROPELUCCHI > MAUROPELUCCHI.COM

Clusters

! Find a cluster...

[Build a New Cluster](#)

ATLAS

Clusters

Data Lake BETA

SECURITY

Database Access

Network Access

Advanced

PROJECT

Access Management

Activity Feed

Alerts 0

Integrations

Settings

SERVICES

Charts

Stitch

Triggers

HELP

Docs

Support !

SANDBOX

MyCluster

Version 4.2.5

[CONNECT](#) [METRICS](#) [COLLECTIONS](#) [...](#)

CLUSTER 1

M0 Sandbox

REGION

AWS / Ireland

TYPE

Replica Set - 3 nodes

LINKED STITCH APP

None Linked

Operations R: 0 W: 0

100.0/s

Last 6 Hours

Logical Size 3.5 MB

512.0 MB max

Last 30 Days

Connections 0

500 max

Last 6 Hours

Enhance Your Experience

For dedicated throughput, richer metrics and enterprise security options, upgrade your cluster now!

[Upgrade](#)

System Status: All Good Last Login: 91.255.53.33

©2020 MongoDB, Inc. [Status](#) [Terms](#) [Privacy](#) [Atlas Blog](#) [Contact Sales](#)



MongoDB Atlas



mongoDB Atlas All Clusters

Please set your time zone Usage This Month:\$0.00 details Preferences Mauro Pelucchi ▾

CONTEXT mauropelucchi.com ▾ MAUROPELUCCHI > MAUROPELUCCHI.COM

ATLAS

Clusters

Find a cluster...

SANDBOX

MyCluster Version 4.2.5

CONNECT METRICS COLLECTIONS ...

CLUSTER TIER M0 Sandbox (General)

REGION AWS / Ireland (eu-west-1)

SIZE 3.5 MB 512.0 MB max

Operations R: 0 W: 0 100.0/s Last 6 Hours

Logical Size 3.5 MB 0.0 B Last 30 Days

Connections 0 500 max Last 6 Hours

Enhance Your Experience
For dedicated throughput, richer metrics and enterprise security options,
upgrade your cluster now!

Upgrade

Red arrow pointing to the MyCluster card.

Clusters Data Lake BETA

SECURITY Database Access Network Access Advanced

PROJECT Access Management Activity Feed Alerts 0 Integrations Settings

SERVICES Charts Stitch Triggers

HELP Docs Support

System Status: All Good Last Login: 91.255.53.33
©2020 MongoDB, Inc. Status Terms Privacy Atlas Blog Contact Sales

Build a New Cluster

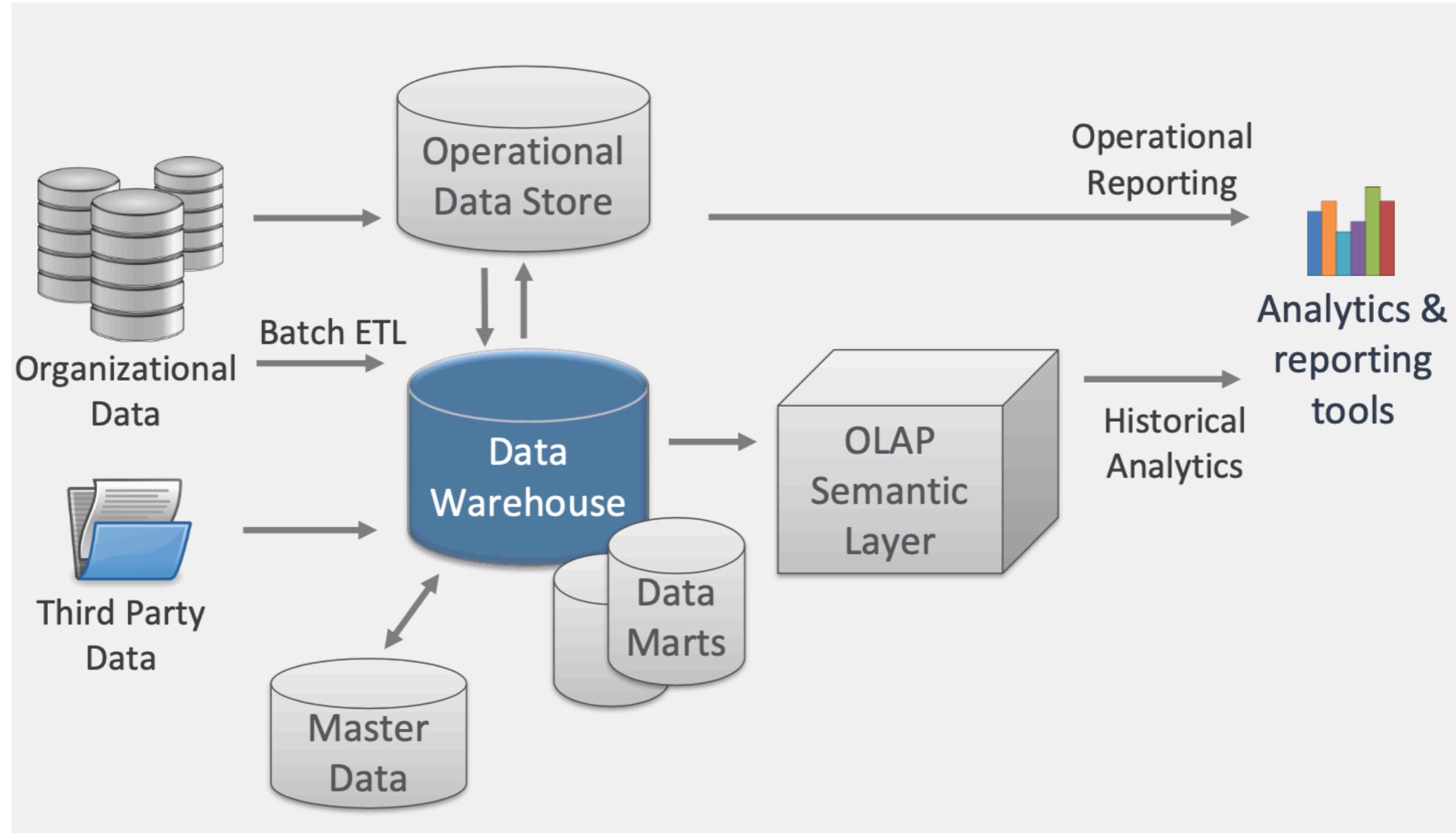
AWS Glue



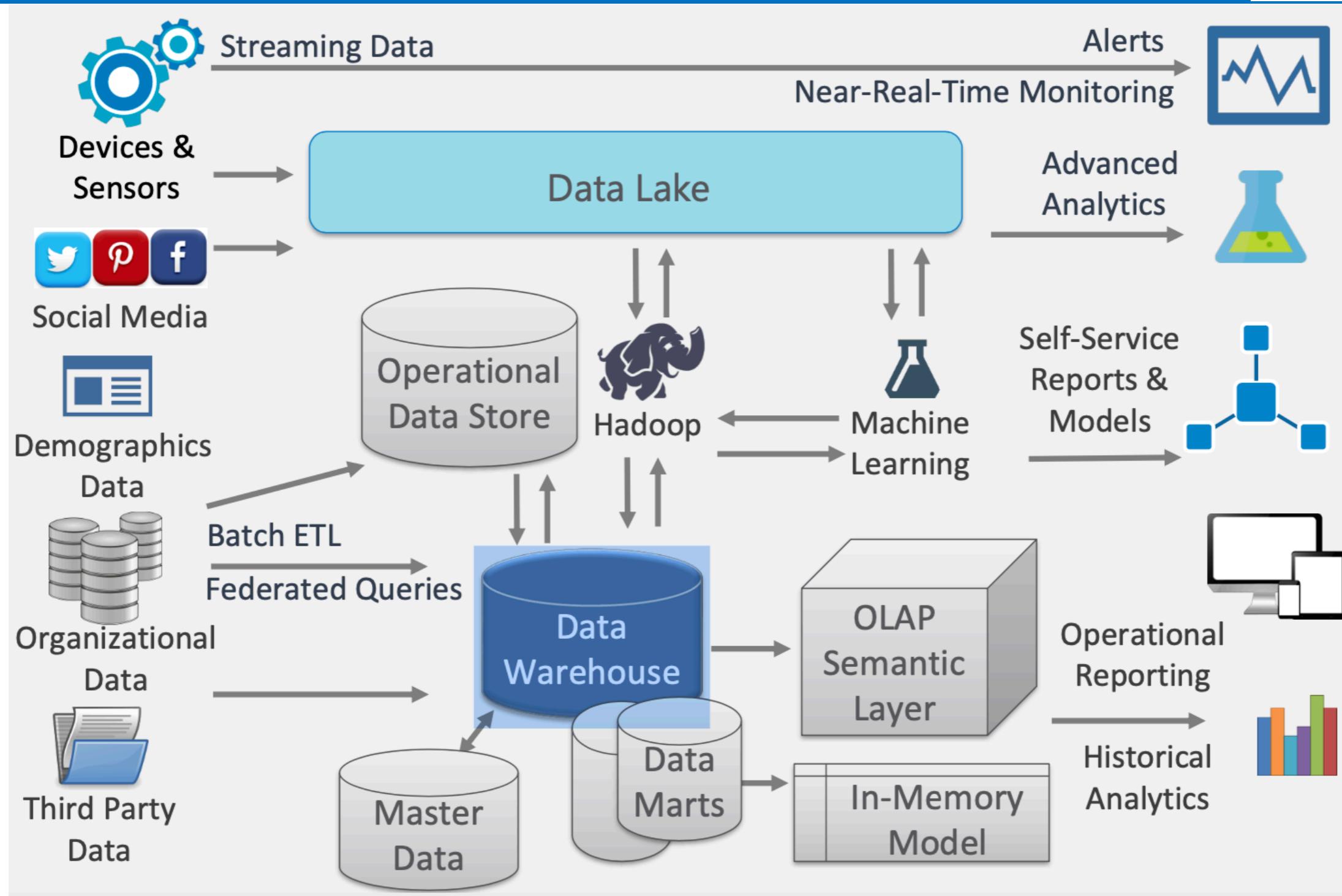
Cosa è un data warehouse?

Come deve essere fatto un data warehouse moderno?

Traditional Data Warehousing



Modernizing an Existing DW



Features for a Modern Data Warehouse



Coexists
with Data
lake

Variety of
data
sources

Multi
Platform
architecture

Flexible
deployment

Support all
user types
& levels

Near real-
time data

Advanced
analytics

Cloud
integration

APIs

Data
catalog

Scalable
architecture

Governanc
e model &
MDM

Data Lake



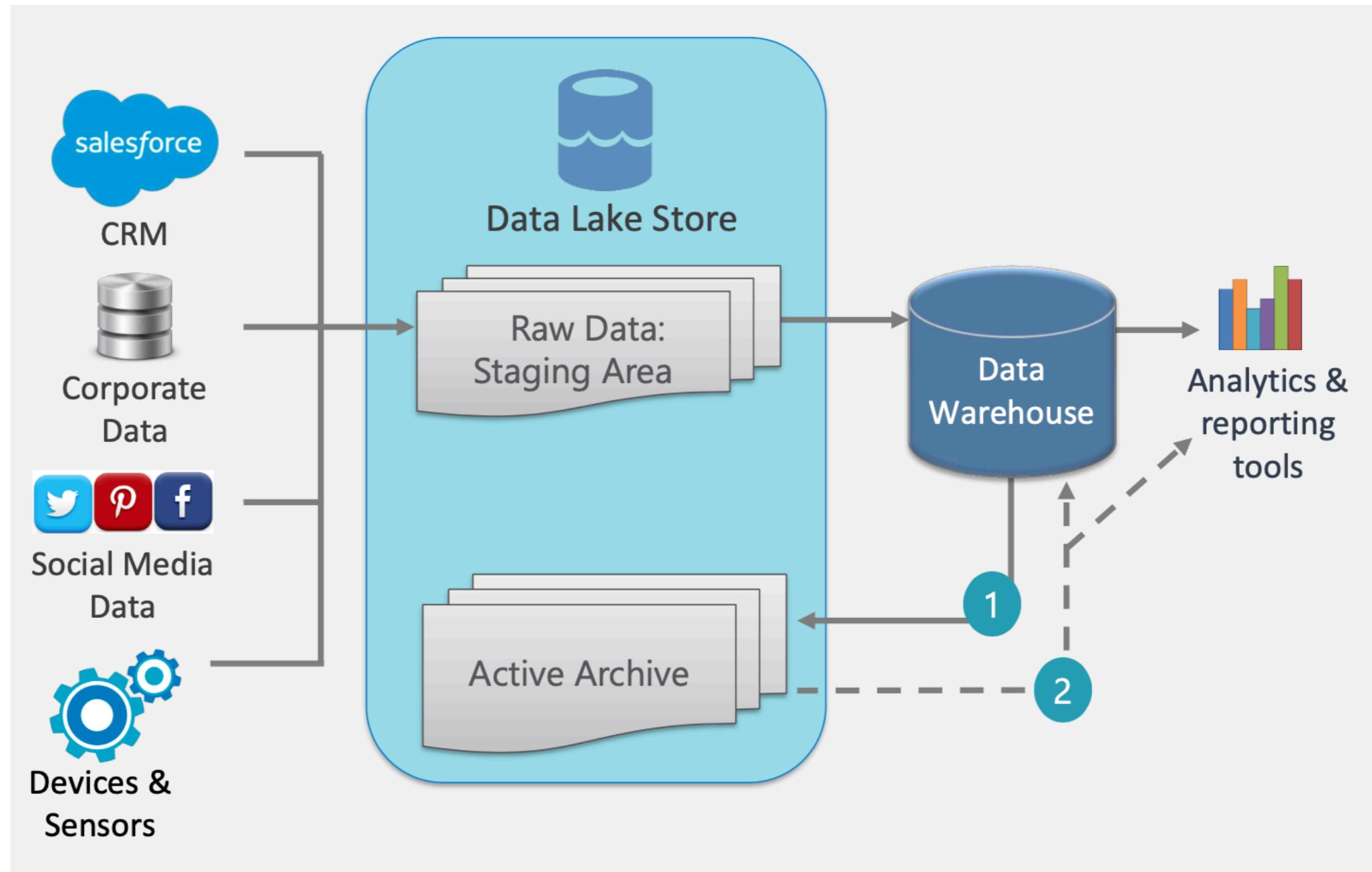
- A repository for analyzing large quantities of disparate sources of data in its native format
- One architectural platform to house all types of data:
 - Machine-generated data (ex: IoT, logs)
 - Human-generated data (ex: tweets, e-mail, video, images, ...)
 - Traditional operational data (ex: sales, inventory)



Data Lake

- Reduce up-front effort by ingesting data in any format without requiring a schema initially
- Make acquiring new data easy, so it can be available for data science & analysis quickly
- Store large volume of multi-structured data in its native format
- Defer work to schematize after value & requirements are known
- Achieve agility faster than a traditional data warehouse can
- Speed up decision-making ability
- Storage for additional types of data which were historically difficult to obtain

Data Lake



Data Lake Implementation



- A data lake is a conceptual idea. It can be implemented with one or more technologies.
 - HDFS (Hadoop Distributed File Storage) is a very common option for data lake storage.
 - NoSQL databases are also very common.
 - **MongoDB!!!**
 - Object stores (like Amazon S3 or Azure Blob Storage)

Data Lake Challenge



Complex,
multi-
layered
architecture

Unknown
storage &
scalability

Working
with un-
curated
data

Performance

Data
retrieval

Data quality

Governance

Redundant
Effort

Supponete di dover raccogliere ed elaborare dati per un grande sito di ecommerce...

Dovrete gestire dati transazionali, testi, immagini, video, log,

...

Da database tradizionali e db NoSQL

Per costruire un DWH per supportare le analisi delle vendite da una App

Quali sfide vi aspettano?

Challenges



- Ingest data from various data sources and join them together
- Enrich raw data
- Clean & Merge raw data
- Convert data to a new format to ensure efficient querying
- Grant access to roles based on the data classification
- SQL Access for Data Scientists
- Data Visualization with charts and graphs

Data Pipeline Challenges



Short timeline to produce a data layer.

Resources are not data scientists or ETL developers.

Data is in different spreadsheets or existing databases

Catalog vs. warehousing

What is AWS Glue?



Fully-managed,
server-less extract-transform-load (ETL) service

AWS Glue



- Serverless ETL
- Universal Data Catalog
- Open source Apache Spark environment
- DynamicFrame – Built in functions
- Seamless integration with AWS services
- Support for on-premises data stores

AWS Glue Components



Data Catalog

Discover

Automatic crawling
Apache Hive Metastore compatible
Integrated with AWS analytic services



Job Authoring

Develop

Auto-generates ETL code
Python and Apache Spark
Edit, Debug, and Explore

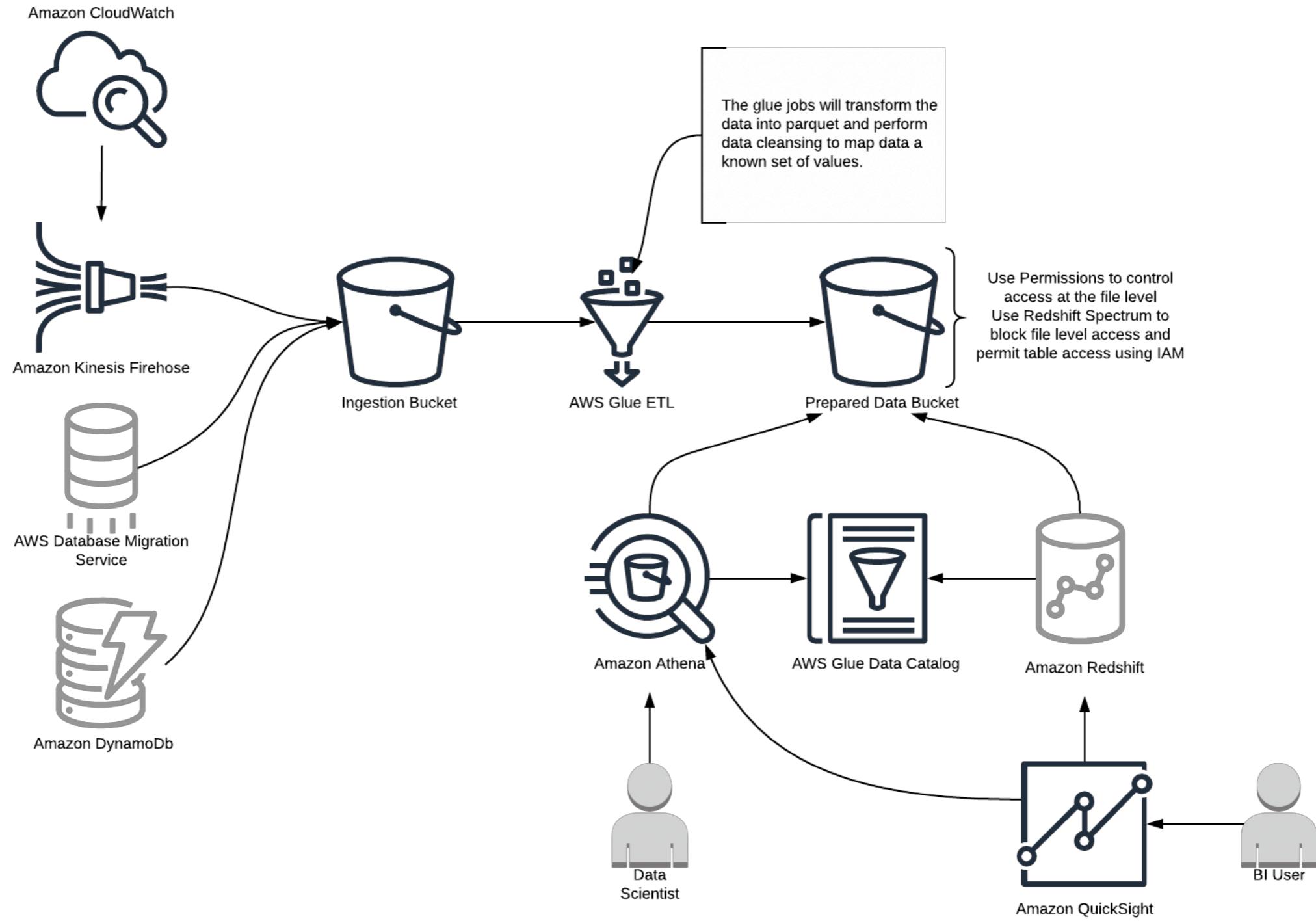


Job Execution

Deploy

Serverless execution
Flexible scheduling
Monitoring and alerting

High-level architecture





4 Steps

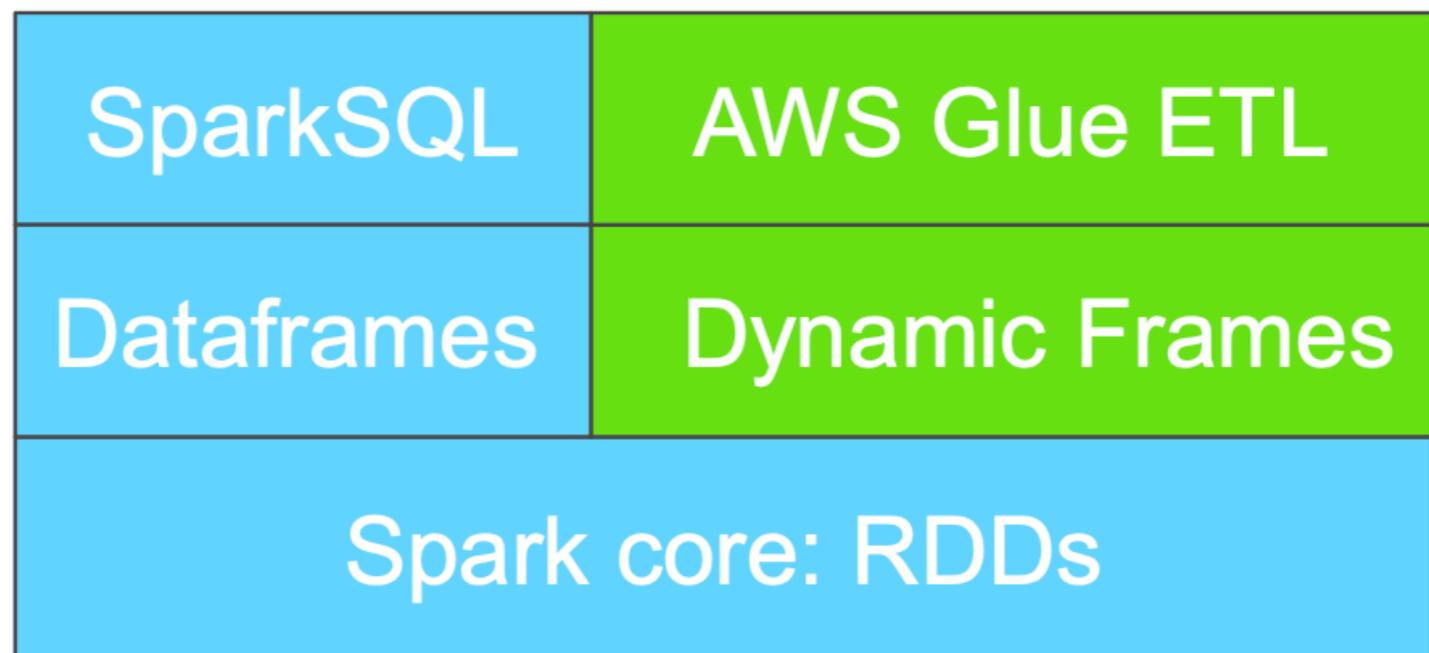
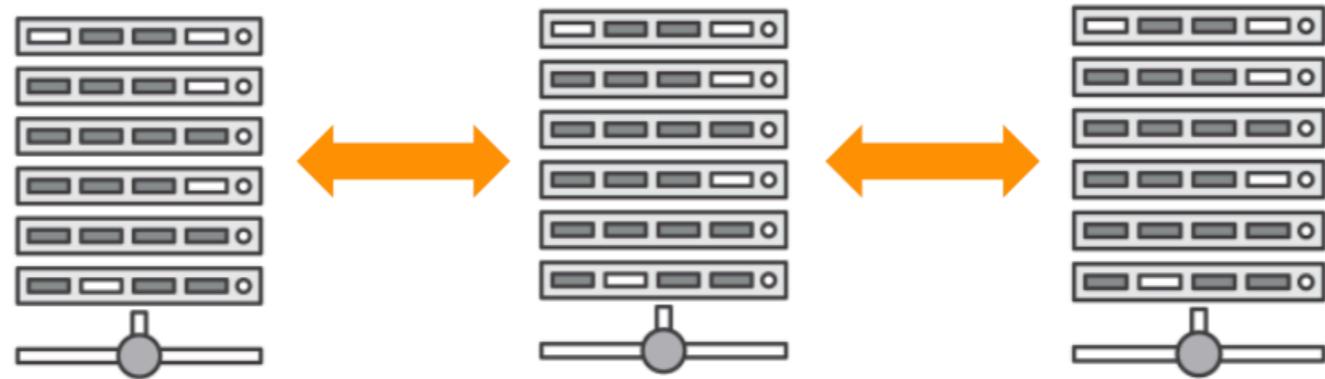
- Crawl and catalogue your data
- Specify mappings to generate scripts
- Interactively edit and explore with dev-endpoints
- Schedule a job for running in production

Server-less job execution



- No need to provision, configure, or manage servers
- Auto-configure VPC & role-based access security & isolation preserved
- Customers can specify job capacity (DPUs)
 - 1 DPU = 16GB x 4 vCPU
- Automatically scale resources
- Only pay for the resources you consume per-second billing (10-minute min)

Under the hood



MyTEDx



Iniziamo!



<https://console.aws.amazon.com/console/home?region=us-east-1>

AWS Management Console



AWS Management Console

AWS services

Find Services
You can enter names, keywords or acronyms.

Example: Relational Database Service, database, RDS

▼ Recently visited services

IAM S3 AWS Glue
 Step Functions CloudWatch

► All services

Build a solution
Get started with simple wizards and automated workflows.

Launch a virtual machine
With EC2
2-3 minutes

Build a web app
With Elastic Beanstalk
6 minutes

Build using virtual servers
With Lightsail
1-2 minutes

Access resources on the go

Access the Management Console using the AWS Console Mobile App. [Learn more](#)

Explore AWS

Amazon Redshift
Fast, simple, cost-effective data warehouse that can extend queries to your data lake. [Learn more](#)

Run Serverless Containers with AWS Fargate
AWS Fargate runs and scales your containers without having to manage servers or clusters. [Learn more](#)

Scalable, Durable, Secure Backup & Restore with Amazon S3
Discover how customers are building backup & restore solutions on AWS that save money.

Check your zone



Services ▾ Resource Groups ▾ ⚙

Pelucchi Mauro ▾ N. Virginia ▾

AWS Management Console

AWS services

Find Services
You can enter names, keywords or acronyms.

▼ Recently visited services
 AWS Glue Athena EC2
 S3 EMR

► All services

Build a solution
Get started with simple wizards and automated workflows.

Launch a virtual machine With EC2 2-3 minutes
Build a web app With Elastic Beanstalk 6 minutes
Build using virtual servers With Lightsail 1-2 minutes
Register a domain With Route 53 3 minutes

Connect an IoT device **Start migrating to AWS** **Start a development project** **Deploy a serverless microservice**

Access resources
 Access the Management Console Mobile

Explore AWS

EMR Migration Guide
Move your on-premises Amazon EMR. [Learn more](#)

AWS IQ
Connect with AWS Certified demand consultants

Amazon DocumentDB
New role-based access least privilege access and build multi-tenant applications. [Learn more](#)

AMD Powered EC2 Instances
Featuring AMD EPYC processors provide up to 10% lower cost than comparable instances. [Learn more](#)

US East (N. Virginia) us-east-1
US East (Ohio) us-east-2
US West (N. California) us-west-1
US West (Oregon) us-west-2

Asia Pacific (Hong Kong) ap-east-1
Asia Pacific (Mumbai) ap-south-1
Asia Pacific (Seoul) ap-northeast-2
Asia Pacific (Singapore) ap-southeast-1
Asia Pacific (Sydney) ap-southeast-2
Asia Pacific (Tokyo) ap-northeast-1

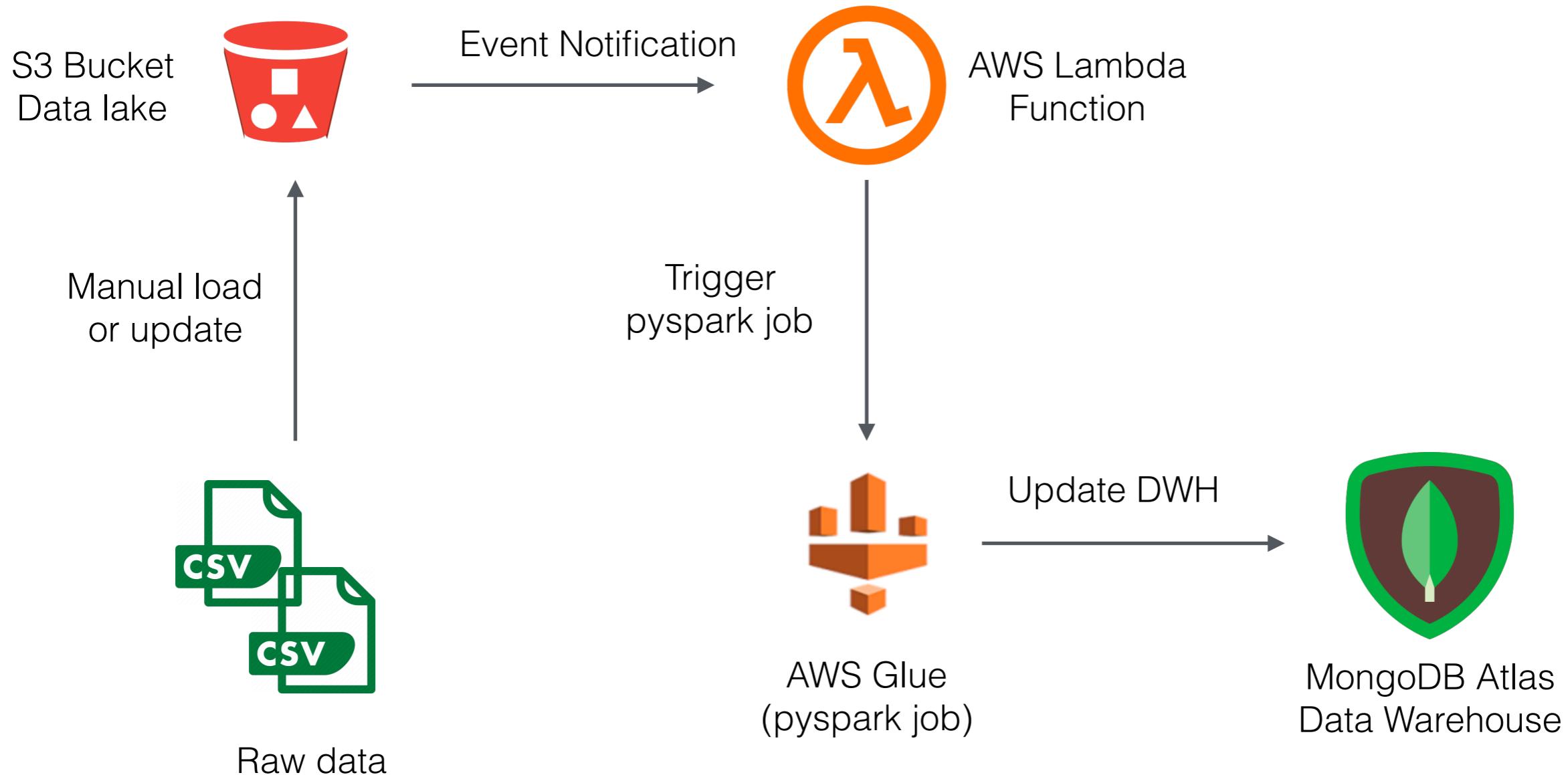
Canada (Central) ca-central-1

Europe (Frankfurt) eu-central-1
Europe (Ireland) eu-west-1
Europe (London) eu-west-2
Europe (Paris) eu-west-3
Europe (Stockholm) eu-north-1

Middle East (Bahrain) me-south-1

South America (São Paulo) sa-east-1

Main Architecture



Data lake



Create bucket

General configuration

Bucket name
 Bucket name must be unique and must not contain spaces or uppercase letters. [See rules for bucket naming](#)

Region

Bucket settings for Block Public Access

Public access is granted to buckets and objects through access control lists (ACLs), bucket policies, access point policies, or all. In order to ensure that public access to this bucket and its objects is blocked, turn on Block all public access. These settings apply only to this bucket and its access points. AWS recommends that you turn on Block all public access, but before applying any of these settings, ensure that your applications will work correctly without public access. If you require some level of public access to this bucket or objects within, you can customize the individual settings below to suit your specific storage use cases. [Learn more](#)

Block all public access
Turning this setting on is the same as turning on all four settings below. Each of the following settings are independent of one another.

Block public access to buckets and objects granted through new access control lists (ACLs)
S3 will block public access permissions applied to newly added buckets or objects, and prevent the creation of new public access ACLs for existing buckets and objects. This setting doesn't change any existing permissions that allow public access to S3 resources using ACLs.

Block public access to buckets and objects granted through any access control lists (ACLs)
S3 will ignore all ACLs that grant public access to buckets and objects.

Block public access to buckets and objects granted through new public bucket or access point policies
S3 will block new bucket and access point policies that grant public access to buckets and objects. This setting doesn't change any existing policies that allow public access to S3 resources.

Block public and cross-account access to buckets and objects through any public bucket or access point policies
S3 will ignore public and cross-account access for buckets or access points with policies that grant public access to buckets and objects.

Advanced settings

[Cancel](#) [Create bucket](#)

Create a bucket for our data lake, in the bucket we will put our raw data



Download raw data from GitHub



https://github.com/mauropelucchi/tedx_dataset

mauropelucchi / **tedx_dataset**

Watch 0 Star 0 Fork 0

Code Issues 0 Pull requests 0 Actions Projects 0 Wiki Security Insights Settings

TEDx Dataset Edit

tedx scraper dataset python selenium Manage topics

7 commits 1 branch 0 packages 0 releases 1 contributor MIT

Branch: master New pull request Create new file Upload files Find file Clone or download ▾

mauropelucchi add final dataset	Latest commit 553ee95 3 days ago
LICENSE	Initial commit 4 days ago
README.md	Update README.md 4 days ago
TEDx_Scraper.ipynb	add final dataset 3 days ago
tags_dataset.csv	add final dataset 3 days ago
tedx_dataset.csv	add final dataset 3 days ago
watch_next_dataset.csv	add final dataset 3 days ago
README.md	edit



Data lake



Amazon S3 > unibg-tedx-data

unibg-tedx-data

Overview Properties Permissions Map

Upload Create folder Download Actions

Upload an object

Buckets are globally unique containers for everything that you store.

Learn more

Get started

Next

Upload

① Select files ② Set permissions ③ Set properties

3 Files Size: 14.3 MB Target path: unibg-tedx-data

To upload a file larger than 160 GB, use the AWS CLI, AWS SDK, or Amazon S3 REST API. [Learn more](#)

+ Add more files

	tags_dataset.csv - 1.7 MB	
	tedx_dataset.csv - 2.3 MB	
	watch_next_dataset.csv - 10.3 MB	

tags_dataset.csv
tedx_dataset.csv
watch_next_dataset.csv

A red arrow points from the top right towards the 'Set properties' button in the upload dialog. A second red arrow points from the bottom left towards the 'Upload' button in the dialog.



Data lake on S3



unibg-tedx-data

Overview Properties Permissions Management Access points

Type a prefix and press Enter to search. Press ESC to clear.

Upload Create folder Download Actions

Name
<input type="checkbox"/> tags_dataset.csv
<input type="checkbox"/> tedx_dataset.csv
<input type="checkbox"/> watch_next_dataset.csv

S3 Bucket for Scripts and logs



Let's create a bucket to store scripts, logs, ...
unibg-tedx-script

Create bucket

General configuration

Bucket name
 Bucket name must be unique and must not contain spaces or uppercase letters. [See rules for bucket naming](#)

Region

Bucket settings for Block Public Access

Public access is granted to buckets and objects through access control lists (ACLs), bucket policies, access point policies, or all. In order to ensure that public access to this bucket and its objects is blocked, turn on Block all public access. These settings apply only to this bucket and its access points. AWS recommends that you turn on Block all public access, but before applying any of these settings, ensure that your applications will work correctly without public access. If you require some level of public access to this bucket or objects within, you can customize the individual settings below to suit your specific storage use cases. [Learn more](#)

Block all public access
Turning this setting on is the same as turning on all four settings below. Each of the following settings are independent of one another.

Block public access to buckets and objects granted through new access control lists (ACLs)
S3 will block public access permissions applied to newly added buckets or objects, and prevent the creation of new public access ACLs for existing buckets and objects. This setting doesn't change any existing permissions that allow public access to S3 resources using ACLs.

Block public access to buckets and objects granted through any access control lists (ACLs)
S3 will ignore all ACLs that grant public access to buckets and objects.

Block public access to buckets and objects granted through new public bucket or access point policies
S3 will block new bucket and access point policies that grant public access to buckets and objects. This setting doesn't change any existing policies that allow public access to S3 resources.

Block public and cross-account access to buckets and objects through any public bucket or access point policies
S3 will ignore public and cross-account access for buckets or access points with policies that grant public access to buckets and objects.

► Advanced settings

[Cancel](#) [Create bucket](#)



S3 Bucket for Scripts and logs



Create a “logs” folder...

Amazon S3 > unibg-tedx-script

unibg-tedx-script

Overview Properties Permissions Management Access points

Type a prefix and press Enter to search. Press ESC to clear.

Upload Create folder Download Actions ▾

Name ▾

logs

When you create a new object in the S3 console creates an object with the above name appended by suffix "/" and that object is displayed in the S3 console. Choose the encryption setting for the object:

None (Use bucket's encryption setting)

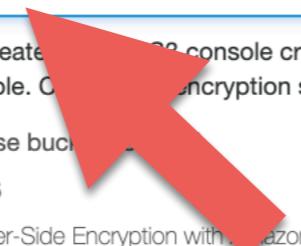
AES-256

Use Server-Side Encryption with Amazon S3-Managed Keys (SSE-S3)

AWS-KMS

Use Server-Side Encryption with AWS KMS-Managed Keys (SSE-KMS)

Save Cancel



Manage Access on AWS resources



Screenshot of the AWS Management Console showing the IAM service selected in the search bar.

The search bar at the top contains the text "iam". A dropdown menu is open, showing the "IAM" service with the description "Manage access to AWS resources". A large red arrow points from the text "iam" in the search bar to the "IAM" item in the dropdown menu.

The left sidebar shows the following navigation items:

- History
- AWS Glue
- Console Home
- S3
- Athena
- EMR
- EC2

The main content area displays several AWS services:

- IAM** Manage access to AWS resources
- EC2
- Lightsail
- Lambda
- Batch
- Elastic Beanstalk
- Serverless Application Repository
- AWS Outposts
- EC2 Image Builder
- Amazon Managed Blockchain
- Satellite**
Ground Station
- Quantum Technologies**
Amazon Braket

A search bar at the bottom right contains the placeholder text "Search IAM".

The right sidebar includes the following sections:

- Dashboard**
- Access management**
 - Groups
 - Users
 - Roles**
 - Policies
 - Identity providers
 - Account settings
- Access reports**
 - Access analyzer
 - Archive rules
 - Analyzers
 - Settings
- Credential report
- Organization activity
- Service control policies (SCPs)

Manage Access on AWS resources —> Roles



Identity and Access Management (IAM)

Dashboard

Access management

Groups

Users

Roles

Policies

Identity pro...

Account setti...

Access reports

Access analyzer

Archive rules

Analyzers

Settings

Credential report

Organization activity

Service control policies (SCPs)

Search IAM

AWS account ID:

386545469924

Roles

What are IAM roles?

IAM roles are a secure way to grant permissions to entities that you trust. Examples of entities include the following:

- IAM user in another account
- Application code running on an EC2 instance that needs to perform actions on AWS resources
- An AWS service that needs to act on resources in your account to provide its features
- Users from a corporate directory who use identity federation with SAML

IAM roles issue keys that are valid for short durations, making them a more secure way to grant access.

Additional resources:

- [IAM Roles FAQ](#)
- [IAM Roles Documentation](#)
- [Tutorial: Setting Up Cross Account Access](#)
- [Common Scenarios for Roles](#)

Create role

Delete role

Search

Role name ▾	Trusted entities	Last activity ▾
<input type="checkbox"/> AWSServiceRoleForSupport	AWS service: support (Service-Linked role)	None
<input type="checkbox"/> AWSServiceRoleForTrustedAdvisor	AWS service: trustedadvisor (Service-Linked ...)	None



Manage Access on AWS resources —> Roles



- Now we create a role to execute AWS Glue Job

AWS Glue Role



 AWS service EC2, Lambda and others	 Another AWS account Belonging to you or 3rd party	 Web identity Cognito or any OpenID provider	 SAML 2.0 federation Your corporate directory
---	--	--	---

Allows AWS services to perform actions on your behalf. [Learn more](#)

Choose a use case

Common use cases

[EC2](#)

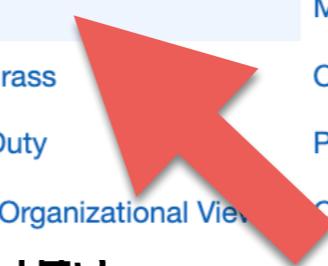
Allows EC2 instances to call AWS services on your behalf.

[Lambda](#)

Allows Lambda functions to call AWS services on your behalf.

Or select a service to view its use cases

API Gateway	CodeDeploy	EMR	KMS	RoboMaker
AWS Backup	CodeGuru	ElastiCache	Kinesis	S3
AWS Chatbot	CodeStar Notifications	Elastic Beanstalk	Lambda	SMS
AWS Support	Comprehend	Elastic Container Service	Lex	SNS
Amplify	Config	Elastic Transcoder	License Manager	SWF
AppStream 2.0	Connect	ElasticLoadBalancing	Machine Learning	SageMaker
AppSync	DMS	Forecast	Macie	Security Hub
Application Auto Scaling	Data Lifecycle Manager	Global Accelerator	MediaConvert	Service Catalog
Application Discovery Service	Data Pipeline	Glue	Migration Hub	Step Functions
Batch	DataSync	Greengrass	OpsWorks	Storage Gateway
Chime	DeepLens	GuardDuty	Personalize	Textract
CloudFormation	Directory Service	Health Organizational View	QLDB	Transfer



AWS Glue Role

TED

Create role

1 2 3 4

▼ Attach permissions policies

Choose one or more policies to attach to your new role.

Create policy



Filter policies		Showing 4 results
	Policy name	Used as
<input checked="" type="checkbox"/>	▶ AWSGlueConsoleFullAccess	None
<input checked="" type="checkbox"/>	▶ AWSGlueConsoleSageMakerNotebookFullAccess	None
<input checked="" type="checkbox"/>	▶ AWSGlueServiceNotebookRole	None
<input checked="" type="checkbox"/>	▶ AWSGlueServiceRole	None

Create role

1 2 3 4

Add tags (optional)

IAM tags are key-value pairs you can add to your role. Tags can include user information, such as an email address, or can be descriptive, such as a job title. You can use the tags to organize, track, or control access for this role. [Learn more](#)

Key	Value (optional)	Remove
Project	MyTEDx	×
Add new key		

You can add 49 more tags.

AWS Glue Role



Create role

Review

Provide the required information below and review this role before you create it.

Role name* Unibg-TEDxGlueRole

Use alphanumeric and '+=,.@-_ characters. Maximum 64 characters.

Role description Allows Glue to call AWS services on your behalf.

Maximum 1000 characters. Use alphanumeric and '+=,.@-_ characters.

Trusted entities AWS service: glue.amazonaws.com

Policies

- AWSGlueConsoleFullAccess
- AWSGlueServiceRole
- AWSGlueServiceNotebookRole
- AWSGlueConsoleSageMakerNotebookFullAccess

Permissions boundary Permissions boundary is not set

The new role will receive the following tag

Key	Value
Project	MyTEDx

* Required

Cancel Previous Create role



AWS Glue Role

Edit the Unibg-TEDx-Role and add
AmazonS3FullAccess

Add permissions to Unibg-TEDxGlueRole

Attach Permissions

Create policy

Filter policies ▾ Showing 4 results

	Policy name ▾	Type	Used as
<input type="checkbox"/>	▶ AmazonDMSRedshiftS3Role	AWS managed	None
<input checked="" type="checkbox"/>	▶ AmazonS3FullAccess	AWS managed	None
<input type="checkbox"/>	▶ AmazonS3ReadOnlyAccess	AWS managed	None
<input type="checkbox"/>	▶ QuickSightAccessForS3StorageManagementAnalyticsReadOnly	AWS managed	None

AWS Glue

TED

AWS services

Find Services
You can enter names, keywords or acronyms.

X

AWS Glue
AWS Glue is a fully managed ETL (extract, transform, and load) service

AWS Lake Formation
AWS Lake Formation makes it easy to set up a secure data lake

 **S3**  **EMR**

▶ All services

AWS Glue

AWS Glue is a fully managed ETL (extract, transform, and load) service that makes it simple and cost-effective to categorize your data, clean it, enrich it, and move it reliably between various data stores.

Get started

[Getting started guide](#)

OMENSI.

AWS Glue

Add a new JOB



AWS Glue

Data catalog
Databases
Tables
Connections
Crawlers
Classifiers
Settings

ETL
Workflows
Jobs (highlighted)
ML Trans
Triggers
Dev endpoints
Notebooks

Security
Security configurations
Tutorials
Add crawler

Jobs A job is your business logic required to perform extract, transform and load (ETL) work. Job runs are initiated by triggers which can be sch

New in AWS Glue
Streaming ETL in AWS Glue (preview): Process streaming data and make it available for analysis in seconds. [Learn More](#)
Reduced start times for AWS Glue Spark jobs (preview): Glue Spark jobs will start in under a minute. [Learn More](#)

Add job Action Filter by tags and attributes

Name Type ETL language Script

You don't have any jobs defined yet.

Add job

We have to create a new job to read the raw data and populate our NoSql DWH

Configure the job properties



Configure the job properties

Name
TEDx-Load-Aggregate-Model

IAM role i
Unibg-TEDxGlueRole

Ensure that this role has permissions to access Amazon S3 sources, targets, temporary directory, scripts, and any libraries used by the job. [Create IAM role](#).

Type
Spark

Glue version
Spark 2.4, Python 3 (Glue Version 1.0)

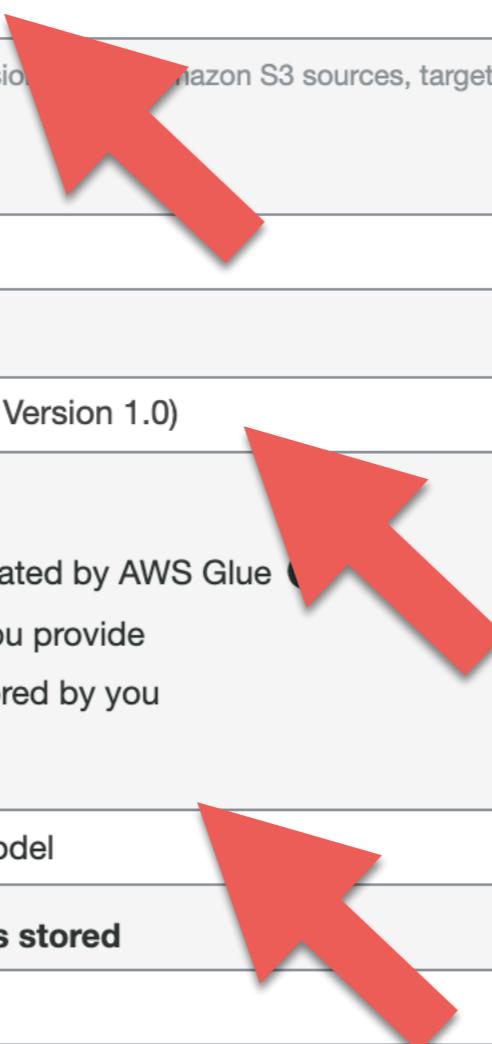
This job runs

- A proposed script generated by AWS Glue
- An existing script that you provide
- A new script to be authored by you

Script file name
TEDx-Load-Aggregate-Model

S3 path where the script is stored
s3://unibg-tedx-script

Temporary directory i
s3://unibg-tedx-script/logs



Save and edit script



Connections

Choose connections required by this job. These connections are used to set up access to your data and must match connections referenced in the

Showing: 0 - 0 < >

All connections

No items available

Required connections

No

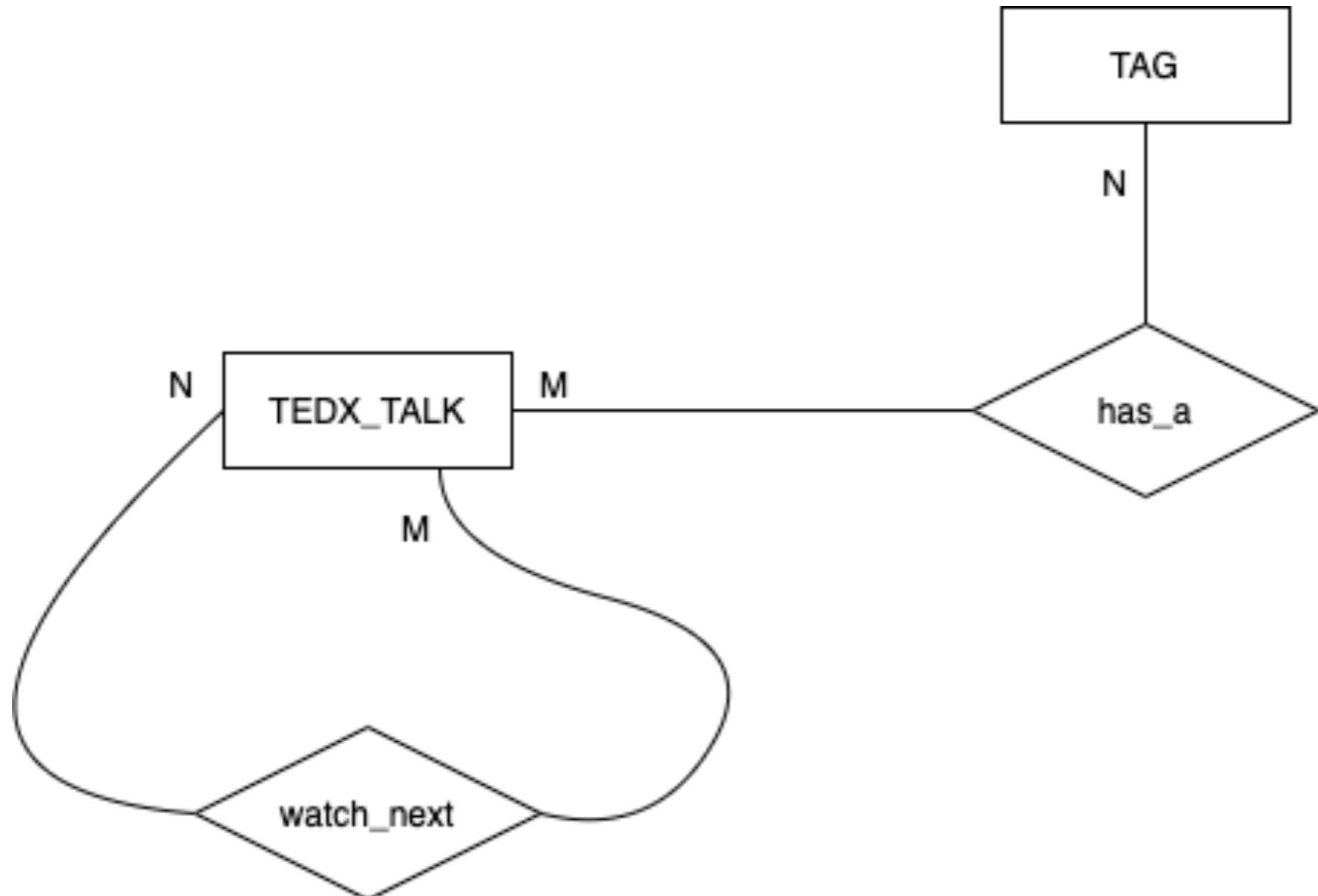
Add connection

Back

Save job and edit script

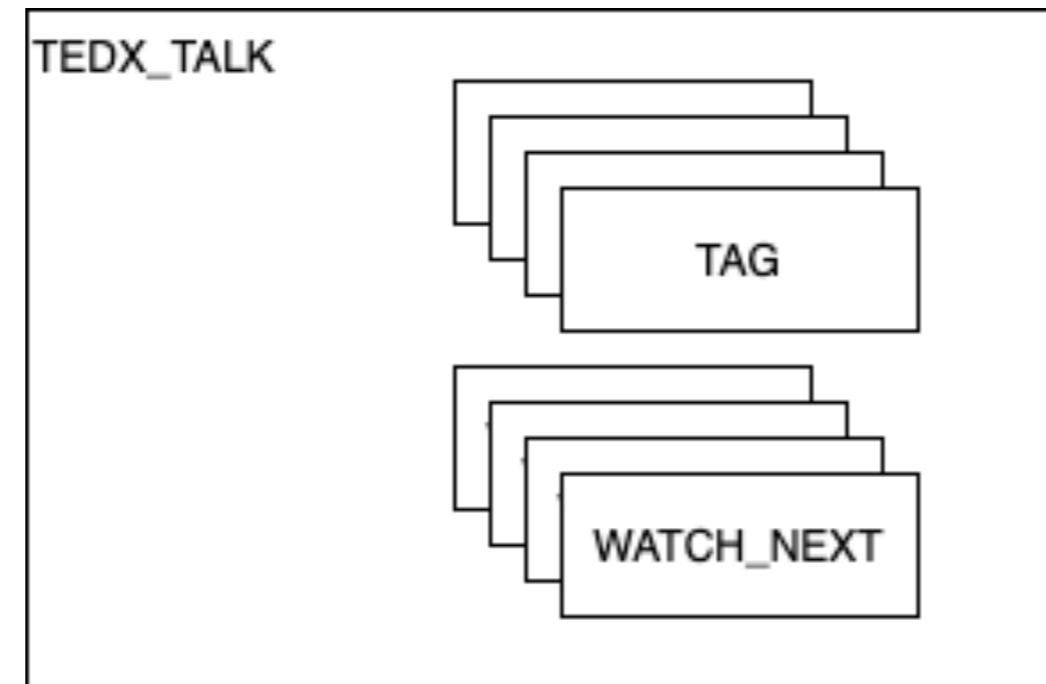


ER Model



Come può essere secondo voi un modello aggregato?

Aggregate model



Pyspark Job



- Import modules
- Read talks dataset —> tedx_dataset
- Init job
- Read tags dataset —> tags_dataset
- Group by tags by talk —> tedx_dataset_agg
- Join tedx_dataset with tedx_dataset_agg
- Store data on MongoDB

Import modules

```
##### TEDx-Load-Aggregate-Model
#####

import sys
import json
import pyspark
from pyspark.sql.functions import col, collect_list, array_join

from awsglue.transforms import *
from awsglue.utils import getResolvedOptions
from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsglue.job import Job
```

Read talks dataset tedx_dataset



```
##### FROM FILES
tedx_dataset_path = "s3://unibg-tedx-data/tedx_dataset.csv"

##### READ PARAMETERS
args = getResolvedOptions(sys.argv, ['JOB_NAME'])

##### START JOB CONTEXT AND JOB
sc = SparkContext()

glueContext = GlueContext(sc)
spark = glueContext.spark_session

job = Job(glueContext)
job.init(args['JOB_NAME'], args)
```

Init job

Read talks dataset tedx_dataset



```
#### READ INPUT FILES TO CREATE AN INPUT DATASET
tedx_dataset = spark.read \
    .option("header", "true") \
    .option("quote", "\") \
    .option("escape", "\\") \
    .csv(tedx_dataset_path)

tedx_dataset.printSchema()
```

Read the CSV file

```
#### FILTER ITEMS WITH NULL POSTING KEY
count_items = tedx_dataset.count()
count_items_null = tedx_dataset.filter("idx is not null").count()

print(f"Number of items from RAW DATA {count_items}")
print(f"Number of items from RAW DATA with NOT NULL KEY {count_items_null}")
```

Check & filter record with
idx is null

Read tags dataset

tags_dataset



```
## READ TAGS DATASET
tags_dataset_path = "s3://unibg-tedx-data/tags_dataset.csv"
tags_dataset = spark.read.option("header","true").csv(tags_dataset_path)
```

Group by tags by talk

tedx_dataset_agg



```
# CREATE THE AGGREGATE MODEL, ADD TAGS TO TEDX_DATASET
tags_dataset_agg = tags_dataset.groupBy(col("idx").alias("idx_ref")).agg(collect_list("tag").alias("tags"))
tags_dataset_agg.printSchema()
tedx_dataset_agg = tedx_dataset.join(tags_dataset_agg, tedx_dataset.idx == tags_dataset_agg.idx_ref, "left") \
    .drop("idx_ref") \
    .select(col("idx").alias("_id"), col("*")) \
    .drop("idx") \
tedx_dataset_agg.printSchema()
```

Join the 2 datasets

Store data on MongoDB



```
mongo_uri = "mongodb://mycluster-shard-00-00-250eh.mongodb.net:27017,mycluster-shard-00-01-250eh.mongodb.net:27017,mycluster-shard-00-02-250eh.mongodb.net:27017"

write_mongo_options = {
    "uri": mongo_uri,
    "database": "unibg_tedx",
    "collection": "tedz_data",
    "username": "admin",
    "password": "VFT0TMyEbfu2CJh0",
    "ssl": "true",
    "ssl.domain_match": "false"}
from awsglue.dynamicframe import DynamicFrame
tedx_dataset_dynamic_frame = DynamicFrame.fromDF(tedx_dataset_agg, glueContext, "nested")

glueContext.write_dynamic_frame.from_options(tedx_dataset_dynamic_frame, connection_type="mongodb", connection_options=write_mongo_options)
```

MongoDB Atlas

How Connect to MyCluster?



X

Connect to MyCluster

✓ Setup connection security

Choose a connection method

Connect

Choose a connection method [View documentation ↗](#)

Get your pre-formatted connection string by selecting your tool below.



Connect with the mongo shell

Interact with your cluster using MongoDB's interactive Javascript interface



Connect your application

Connect your application to your cluster using MongoDB's native drivers



Connect using MongoDB Compass

Explore, modify, and visualize your data with MongoDB's GUI



Go Back

Close



MongoDB Atlas

How Connect to MyCluster?



Connect to MyCluster

✓ Setup connection security

✓ Choose a connection method

Connect

1 Select your driver and version

DRIVER

Scala

VERSION

2.2 or later

2 Add your connection string into your application code

Connection String Only

Full Driver Example

```
mongodb+srv://admin:<password>@mycluster-250eh.mongodb.net/test?r
```

Copy

Replace **<password>** with the password for the user, **admin**, and ensure all special characters are URL encoded.

Having trouble connecting? [View our troubleshooting documentation](#)

Go Back

Close



Run your job



Action ▾ Save Run job Generate diagram ⓘ Insert template at cursor ⓘ Source Target Target

el

```
    42 COUNT_ELEMS_1
    43     print(f"Number of elements: {len(tags_dataset_1)}")
    44     print(f"Number of elements: {len(tags_dataset_2)}")
    45
    46
    47
    48
    49     ## READ TAGS
    50     tags_dataset_1 = pd.read_csv('tags_dataset_1.csv')
    51     tags_dataset_2 = pd.read_csv('tags_dataset_2.csv')
    52
    53
    54
    55     # CREATE THE DATASETS
    56     tags_dataset_1['category'] = 'A'
    57     tags_dataset_2['category'] = 'B'
    58     tedx_dataset = pd.concat([tags_dataset_1, tags_dataset_2])
    59
    60     tedx_dataset.to_csv('tedx_dataset.csv', index=False)
    61
    62
    63
```

Logs

Parameters (optional)

Review and override parameter values, as needed, before running this job. Changes affect this run only. Edit a job to change default parameter values.

- ▶ Advanced properties
- ▶ Monitoring options
- ▶ Tags
- ▶ Security configuration, script libraries, and job parameters

Only job **TEDx-Load-Aggregate-Model** is run. Jobs dependent on the completion of job **TEDx-Load-Aggregate-Model** will not be run. To run a job and trigger dependent jobs, define an on-demand trigger.

Run job

RUM BERGOMENSIS

Job History



Add job Action ▾ Filter by tags and attributes

Name

TEDx-Load-Aggregate-Model

History Details Script Metrics

View run metrics Rewind job bookmark

Run ID	Retry attempt	Run status	Logs	Error logs
jr_deb6e9f845dd...	-	Succeeded	Logs	
jr_cbc1d4a1a07c...	-	Failed	! Sy... Logs	Error logs



Check result on your DWH



MyCluster

Overview Real Time Metrics **Collections** Profiler Performance Advisor Command Line Tools

DATABASES: 1 COLLECTIONS: 1

+ Create Database

NAMESPACES

unibg_tedx.tedz_data

COLLECTION SIZE: 3.3MB TOTAL DOCUMENTS: 4494 INDEXES TOTAL SIZE: 208KB

Find Indexes Aggregation Search^{BETA}

FILTER {"filter": "example"}

QUERY RESULTS 1-20 OF MANY

```
_id: "8d2005ec35280deb6a438dc87b225f89"
main Speaker: "Alexandra Auer"
title: "The intangible effects of walls"
details: "More barriers exist now than at the end of World War II, says designer..."
posted: "Posted Apr 2020"
url: "https://www.ted.com/talks/alexandra_auer_the_intangible_effects_of_wal..."
tags: Array
```

```
_id: "b3072cd11f40eb57fd259555264476c6"
main Speaker: "Elizabeth Gilbert"
title: "It's OK to feel overwhelmed. Here's what to do next"
details: "If you're feeling anxious or fearful during the coronavirus pandemic, ..."
posted: "Posted Apr 2020"
url: "https://www.ted.com/talks/elizabeth_gilbert_it_s_ok_to_feel_overwhelme..."
tags: Array
```

AWS Lambda Functions



Create a new Lambda function



Screenshot of the AWS Lambda console showing the 'Functions' page. The sidebar on the left is titled 'AWS Lambda' and includes links for 'Dashboard', 'Applications', 'Functions' (which is highlighted in orange), and 'Layers'. The main content area shows a table with the following columns: 'Function name', 'Description', 'Runtime', 'Code size', and 'Last modified'. A search bar at the top says 'Filter by tags and attributes or search by keyword'. A red arrow points from the bottom right towards the 'Create function' button, which is located at the top right of the table header. The status message 'There is no data to display.' is centered below the table.

Create a new Lambda function



Basic information

Function name

Enter a name that describes the purpose of your function.

Trigger_Refresh_DWH1

Use only letters, numbers, hyphens, or underscores with no spaces.

Runtime Info

Choose the language to use to write your function.

Python 3.8

Permissions Info

Lambda will create an execution role with permissions to upload logs to Amazon CloudWatch Logs. You can configure and modify permissions further when you add triggers.

▼ Choose or create an execution role

Execution role

Choose a role that defines the permissions of your function. To create a custom role, go to the [IAM console](#).

- Create a new role with basic Lambda permissions
- Use an existing role
- Create a new role from AWS policy templates

(i) Role creation might take a few minutes. Please do not delete the role or edit the trust or permissions policies in this role.

Lambda will create an execution role named Trigger_Refresh_DWH1-role-4ryv1as3, with permission to upload logs to Amazon CloudWatch Logs.

Cancel

Create function

Create a new Lambda function



▼ Designer



Layers Info

Add a layer

▲ Merge earlier

▼ Merge later

Remove

Call AWS Glue Job From lambda function



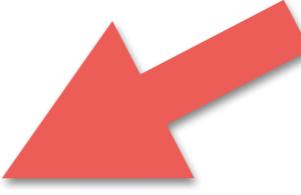
```
# Set up logging
import json
import os
import logging
logger = logging.getLogger()
logger.setLevel(logging.INFO)

# Import Boto 3 for AWS Glue
import boto3
client = boto3.client('glue')

# Variables for the job:
glueJobName = "TEDx-Load-Aggregate-Model"

# Define Lambda function
def lambda_handler(event, context):

    response = client.start_job_run(JobName = glueJobName)
    logger.info('## STARTED GLUE JOB: ' + glueJobName)
    logger.info('## GLUE JOB RUN ID: ' + response['JobRunId'])
    return response
```



Call AWS Glue Job From lambda function



Function code Info

Code entry type

Edit code inline

Runtime

Python 3.8

Handler Info

index.lambda_handler

```
1 # Set up logging
2 import json
3 import os
4 import logging
5 logger = logging.getLogger()
6 logger.setLevel(logging.INFO)
7
8 # Import Boto 3 for AWS Glue
9 import boto3
10 client = boto3.client('glue')
11
12 # Variables for the job:
13 glueJobName = "TEDx-Load-Aggregate-Model"
14
15 # Define Lambda function
16 def lambda_handler(event, context):
17
18     response = client.start_job_run(JobName = glueJobName)
19     logger.info('## STARTED JOB: ' + glueJobName)
20     logger.info('## GLUE JOB RUN ID: ' + response['JobRunId'])
21
22     return response
```

Enable Events on Our S3 Data Lake



unibg-tedx-data

Overview Properties Permissions Management Access points

Versioning

Keep multiple versions of an object in the same bucket.

[Learn more](#)

Disabled

Server access logging

Set up access log records that provide details about access requests.

[Learn more](#)

Disabled

Static website hosting

Host a static website, which does not require server-side technologies.

[Learn more](#)

Disabled

Object-level logging

Record object-level API activity using the CloudTrail data events feature (additional cost).

[Learn more](#)

Disabled

Default encryption

Automatically encrypt objects when stored in Amazon S3

[Learn more](#)

Disabled

Advanced settings

Object lock

Prevent objects from being deleted.

[Learn more](#)

Disabled

Tags

Use tags to track your cost against projects or other criteria.

[Learn more](#)

0 Tags

Transfer acceleration

Enable fast, easy and secure transfers of files to and from your bucket.

[Learn more](#)

Suspended

Events

Receive notifications when specific events occur in your bucket.

[Learn more](#)

0 Active notifications

Requester pays

The requester (instead of the bucket owner) will pay for requests and data transfer.

[Learn more](#)

Disabled

Enable Events on Our S3 Data Lake



Events

Add notification Delete Edit

Name	Events	Filter	Type
New event			

Name Notify_Refresh_Data_Lake

Events

PUT All object delete events
 POST Restore initiated
 COPY Restore completed
 Multipart upload completed Replication time missed threshold
 All object create events Replication time completed after threshold
 Object in RRS lost Replication time not tracked
 Permanently deleted Replication failed
 Delete marker created

Prefix e.g. images/

Suffix e.g. .jpg

Send to Lambda Function

Lambda Trigger_Refersh_DWH

0 Active notifications

Cancel Save

Try to add a new file to the bucket



Amazon S3 > unibg-tedx-data

unibg-tedx-data

Overview Properties Permissions More

Type a prefix and press Enter to search. Press ESC to clear.

Upload + Create folder Download Actions ▾

Name ▾

tags_dataset.csv

tedx_dataset.csv

watch_next_dataset.csv

Storage class ▾

Standard

Standard

Standard

Upload

① Select files ② Set permissions ③ Set properties ④ Review X

1 Files **Size: 2.3 MB** Target path: unibg-tedx-data

To upload a file larger than 160 GB, use the AWS CLI, AWS SDK, or Amazon S3 REST API. [Learn more ↗](#)

+ Add more files

tedx_dataset.csv - 2.3 MB X

Upload Next

Perché la lambda function fallisce?

Manage Access on AWS resources → Roles



Roles > Trigger_Refresh_DWH-role-y67bgjru

Summary

Delete role

Role ARN	arn:aws:iam::386545469924:role/service-role/Trigger_Refresh_DWH-role-y67bgjru
Role description	Edit
Instance Profile ARNs	
Path	/service-role/
Creation time	2020-04-10 13:53 UTC+0200
Last activity	Not accessed in the tracking period
Maximum CLI/API session duration	1 hour Edit

Permissions

Trust relationships

Tags

Access Advisor

Revoke sessions

▼ Permissions policies (1 policy applied)

Attach policies

Add inline policy

Policy name ▾

Policy type ▾

▶ [AWSLambdaBasicExecutionRole-cdca2e44-e2b6-4af8-9a61-8191e5ca4eaa](#)

Managed policy

X

▶ Permissions boundary (not set)

Manage Access on AWS resources —> Roles



Add permissions to Trigger_Refresh_DWH-role-y67bgjru

Attach Permissions

Create policy

Filter policies ▾

	Policy name ▾
<input checked="" type="checkbox"/>	AWSGlueConsoleFullAccess
<input type="checkbox"/>	AWSGlueConsoleSageMakerNotebookFullAccess
<input type="checkbox"/>	AWSGlueServiceNotebookRole
<input type="checkbox"/>	AWSGlueServiceRole

Try to re-add a new file to the bucket



Amazon S3 > unibg-tedx-data

unibg-tedx-data

Overview Properties Permissions Map

Type a prefix and press Enter to search. Press ESC to clear.

Upload Create folder Download Actions ▾

Name ▾

tags_dataset.csv

tedx_dataset.csv

watch_next_dataset.csv

Storage class ▾

Standard

Standard

Standard

Upload

Next

1 Files Size: 2.3 MB Target path: unibg-tedx-data

To upload a file larger than 160 GB, use the AWS CLI, AWS SDK, or Amazon S3 REST API. [Learn more](#)

+ Add more files

tedx_dataset.csv - 2.3 MB

The screenshot shows the AWS S3 console interface. On the left, there's a list of files in the 'unibg-tedx-data' bucket: 'tags_dataset.csv', 'tedx_dataset.csv', and 'watch_next_dataset.csv'. On the right, a modal window titled 'Upload' is open, showing a progress bar for a file named 'tedx_dataset.csv' which is 2.3 MB in size. The modal has four steps: 'Select files', 'Set permissions', 'Set properties', and 'Review'. The 'Select files' step is active. Below the steps, it says 'To upload a file larger than 160 GB, use the AWS CLI, AWS SDK, or Amazon S3 REST API.' and provides a link to learn more. At the bottom of the modal are 'Upload' and 'Next' buttons.

Ora tocca a voi...