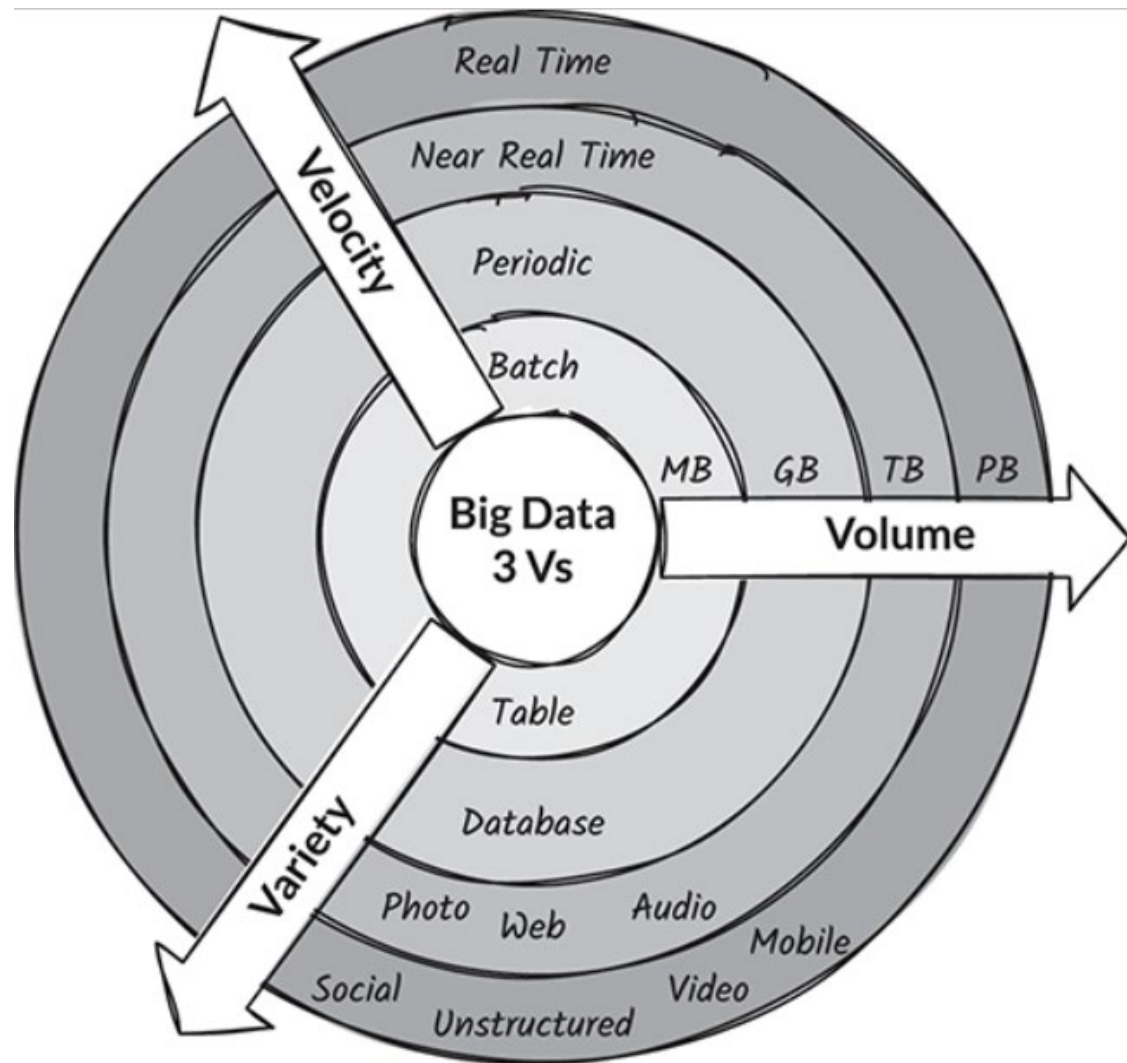


Big Data
a.a. 24/25

Prof. Alfredo Pulvirenti

Obiettivi del corso

- Comprendere i concetti fondamentali di Big Data
- Esplorare le tecnologie e gli strumenti principali
- Acquisire competenze pratiche per l'analisi e la gestione di grandi volumi di dati



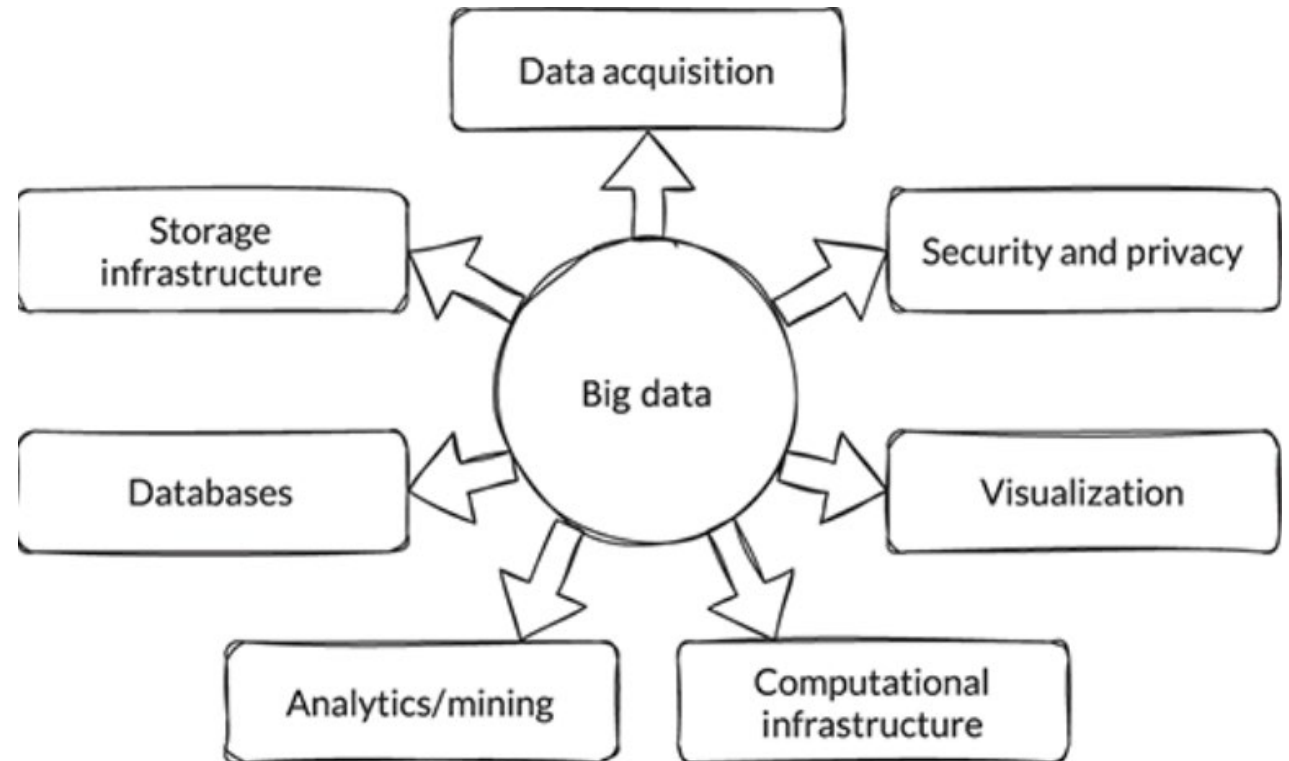
L'era dei dati

- Prodotti velocemente
- Eterogenei
- Potenzialmente ricchi di informazioni
- Volume, Velocità, Varietà, Veridicità, Valore

Big data: Qualsiasi processo di dati in cui la dimensione dei dati stessi è un problema: conservare, trasmettere, elaborare su scala.

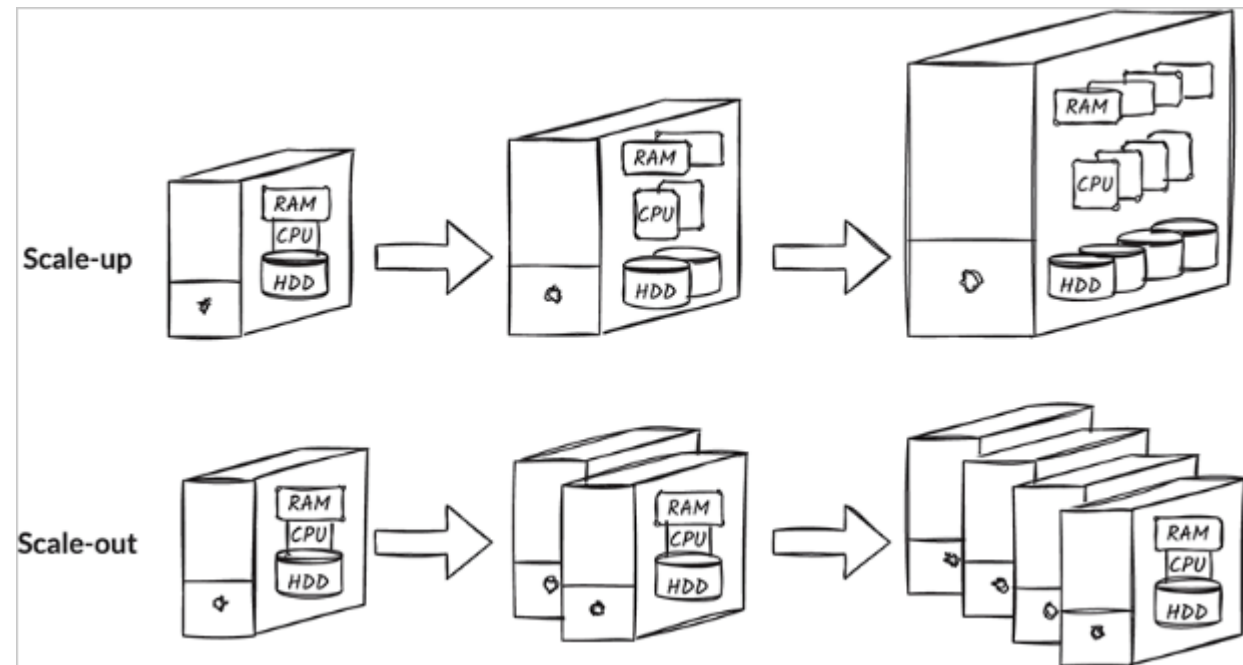
Diversi aspetti dei Big Data

- Concetti spesso confusi perché coinvolgono diversi aspetti



Come interagire con grandi quantità di dati

- Supponiamo di dover analizzare 100 terabyte di dati.
 - Non siamo in grado di memorizzare tutto in una sola macchina
 - Non abbiamo la possibilità di processare questi dati in un tempo ragionevole

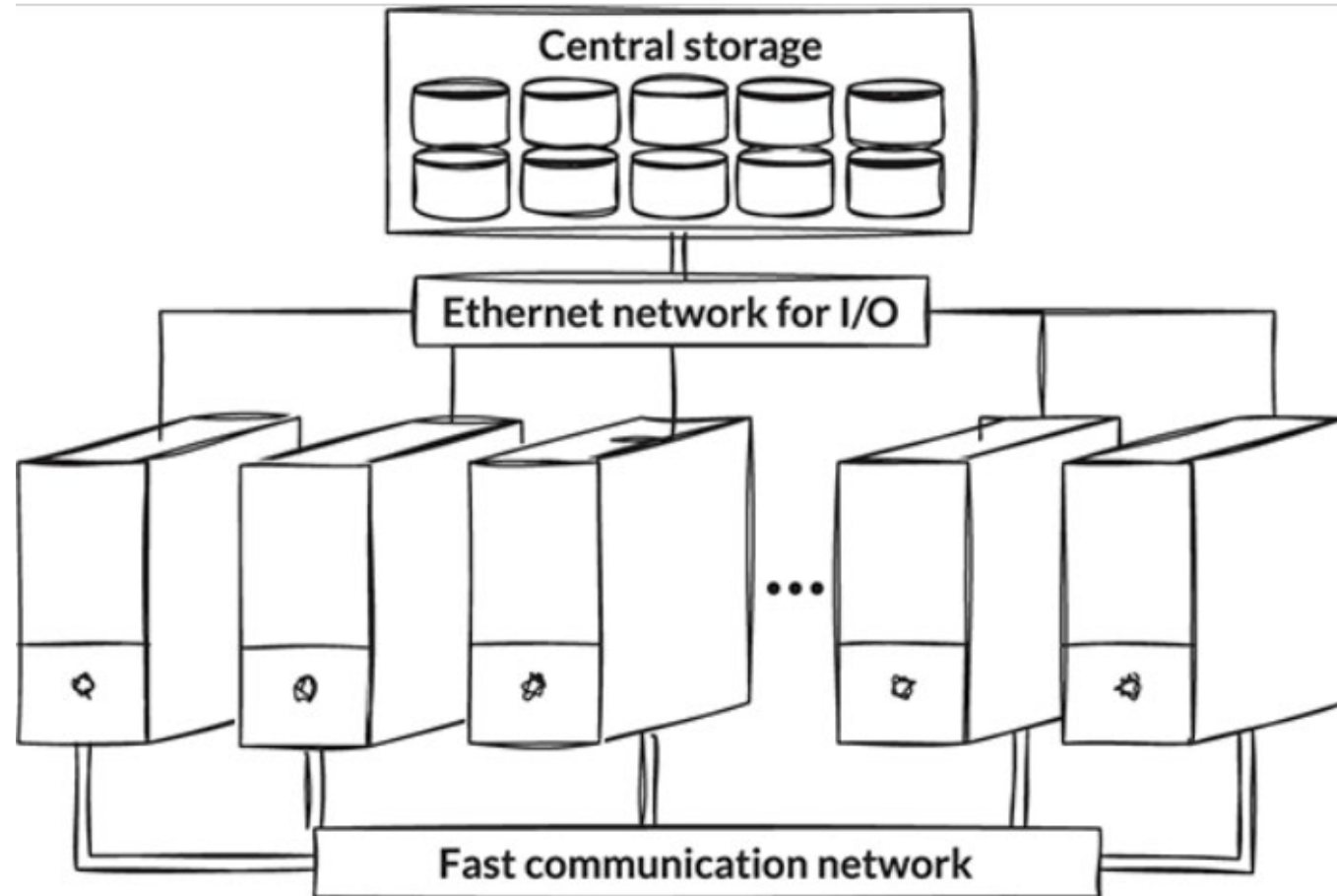


Scale-up	Scale-out
<i>Key idea</i>	<i>Key idea</i>
Add more memory, processors	Add more (cheaper) nodes
<i>Advantages</i>	<i>Advantages</i>
✓ Less energy consumption	✓ Cheaper
✓ Less expense on cooling systems	✓ Fault tolerance possible
✓ Easier to implement solutions	✓ Easy to grow
<i>Disadvantages</i>	<i>Disadvantages</i>
× Price	× More physical space
× No fault tolerance	× Energy costs (electricity and cooling)
× Limited hardware upgrades	× Network equipment required

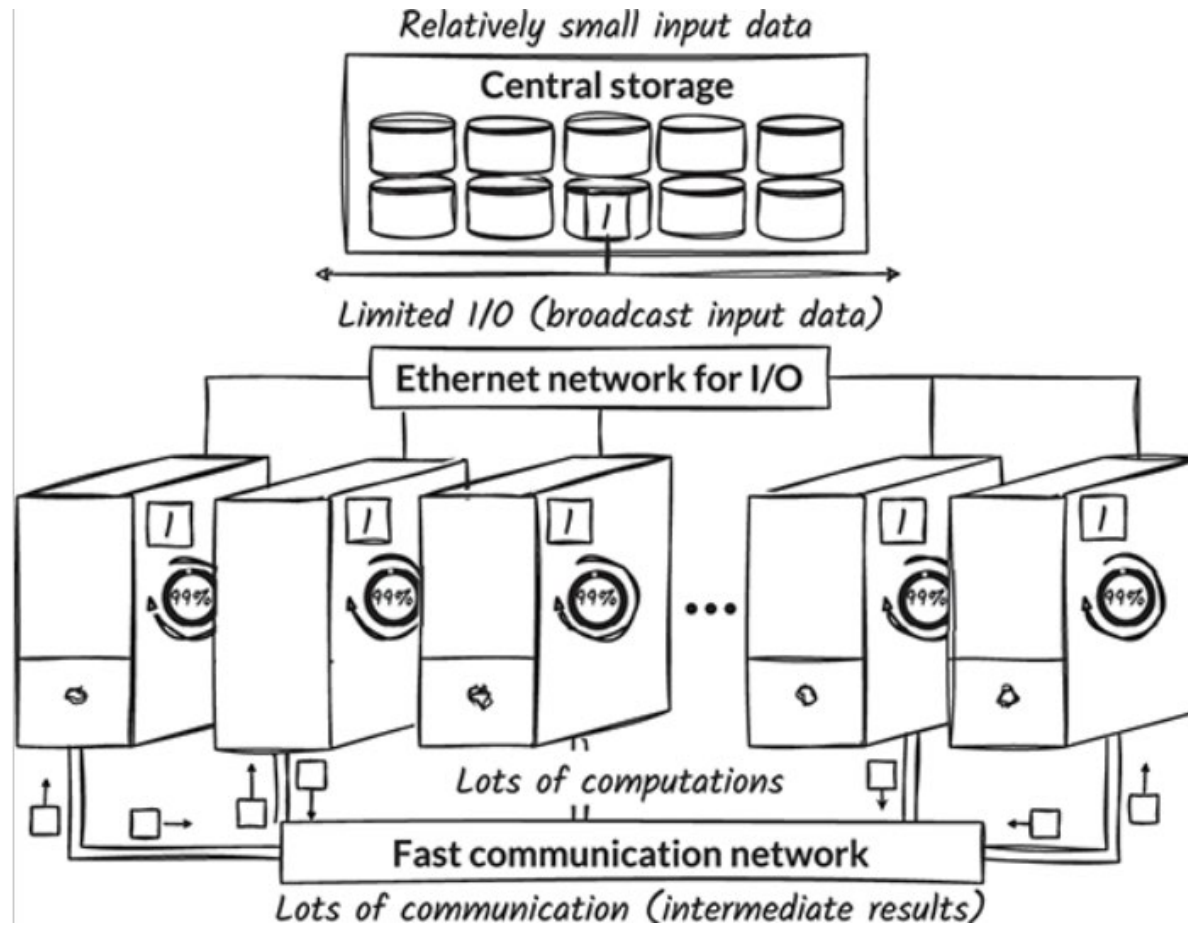
High Performance Computing Vs Big Data Computing

- HPC possiamo immaginare ad un cluster di macchine indipendenti (nodi) connessi tramite una rete di comunicazione ad alta velocità che hanno accesso ad uno storage centralizzato.
- Ogni nodo ha una CPU multicore e la sua RAM e HD.
- Ogni applicazione che viene eseguita su un nodo accede alla sua RAM (non quella degli altri nodi), parliamo di architettura a memoria distribuita.

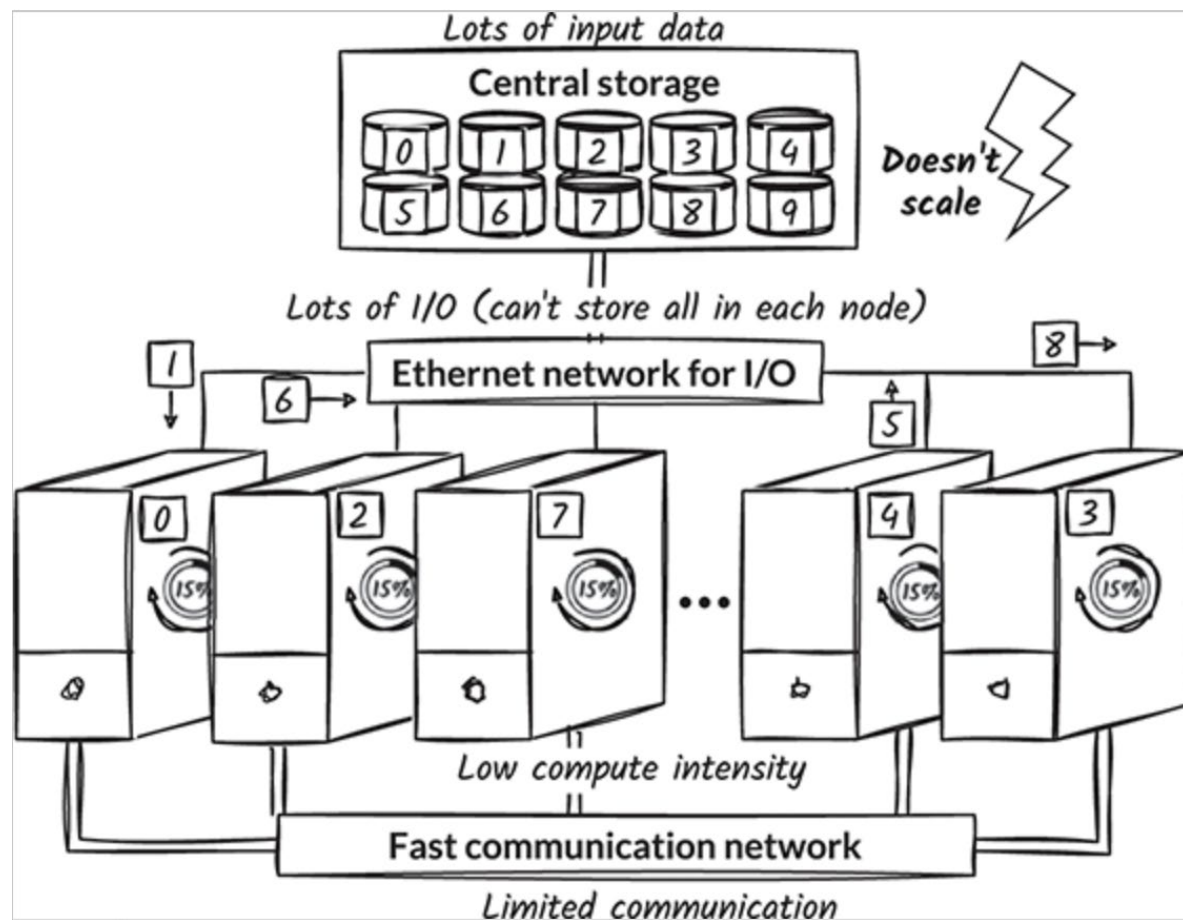
Architettura HPC semplificata



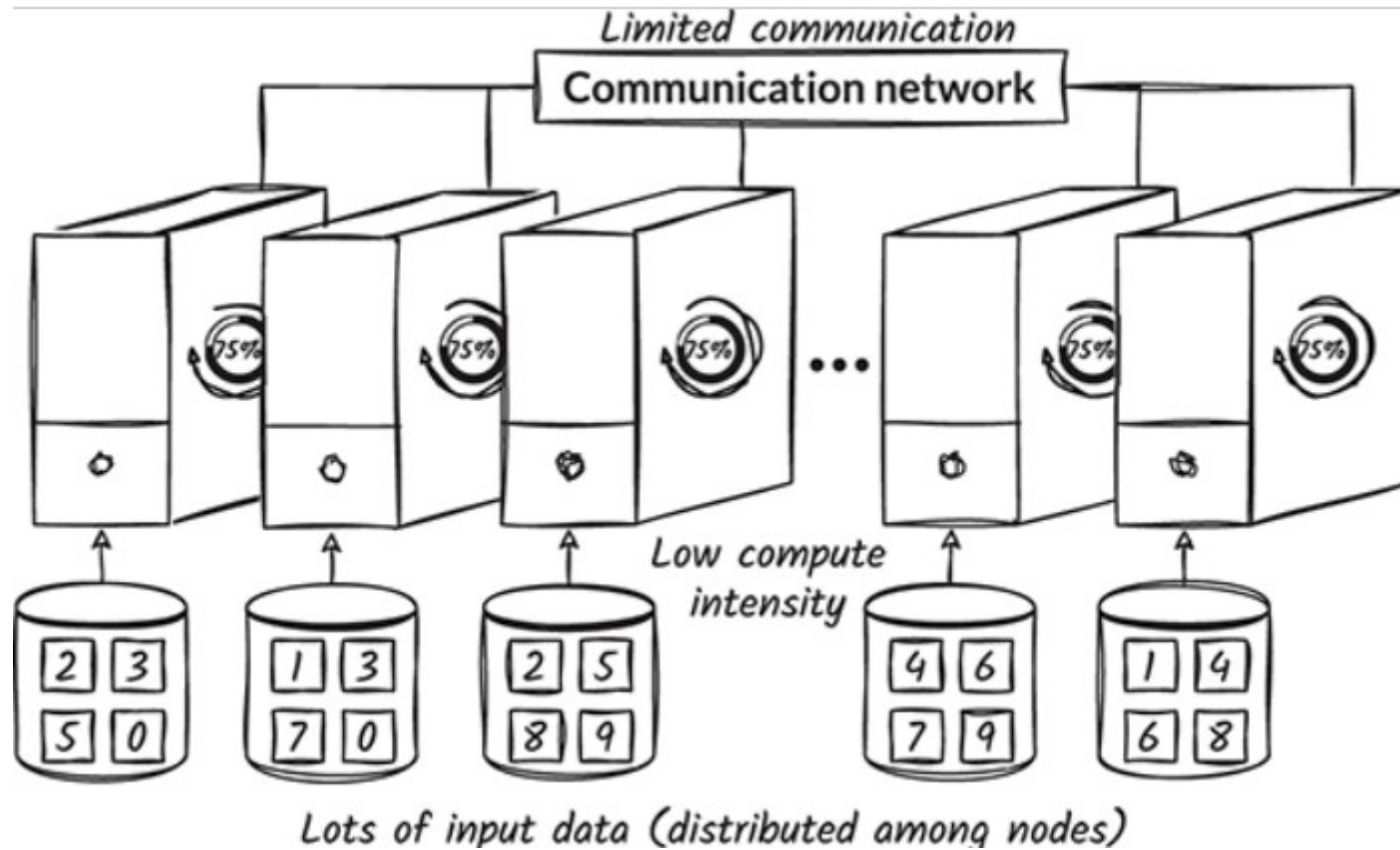
HPC



- Se abbiamo diversi computer connessi dobbiamo implementare programmi distribuiti e paralleli.
- I nodi:
 - Nodo master si occupa dell'orchestrazione e comunicazione
 - Nodi worker si occupano del calcolo
- Un modo classico per interagire con distributed computing è MPI (Message Passing Interface).



Architettura dei Sistemi Big Data



Big data

Focus on data-intensive jobs
Hardware failure common
Code: data science, graphs
Usually mix CPU/GPU and data
Job moved to where the data is located
SIMD model:^a data parallelism
Commodity hardware acceptable

HPC

Focus on computation-intensive jobs
Surprised by hardware failure
Code: simulation, optimization
Mix CPU/GPU
Data moved to where it will be processed
SIMD/MIMD^b model (more general parallelism)
Needs specialized hardware

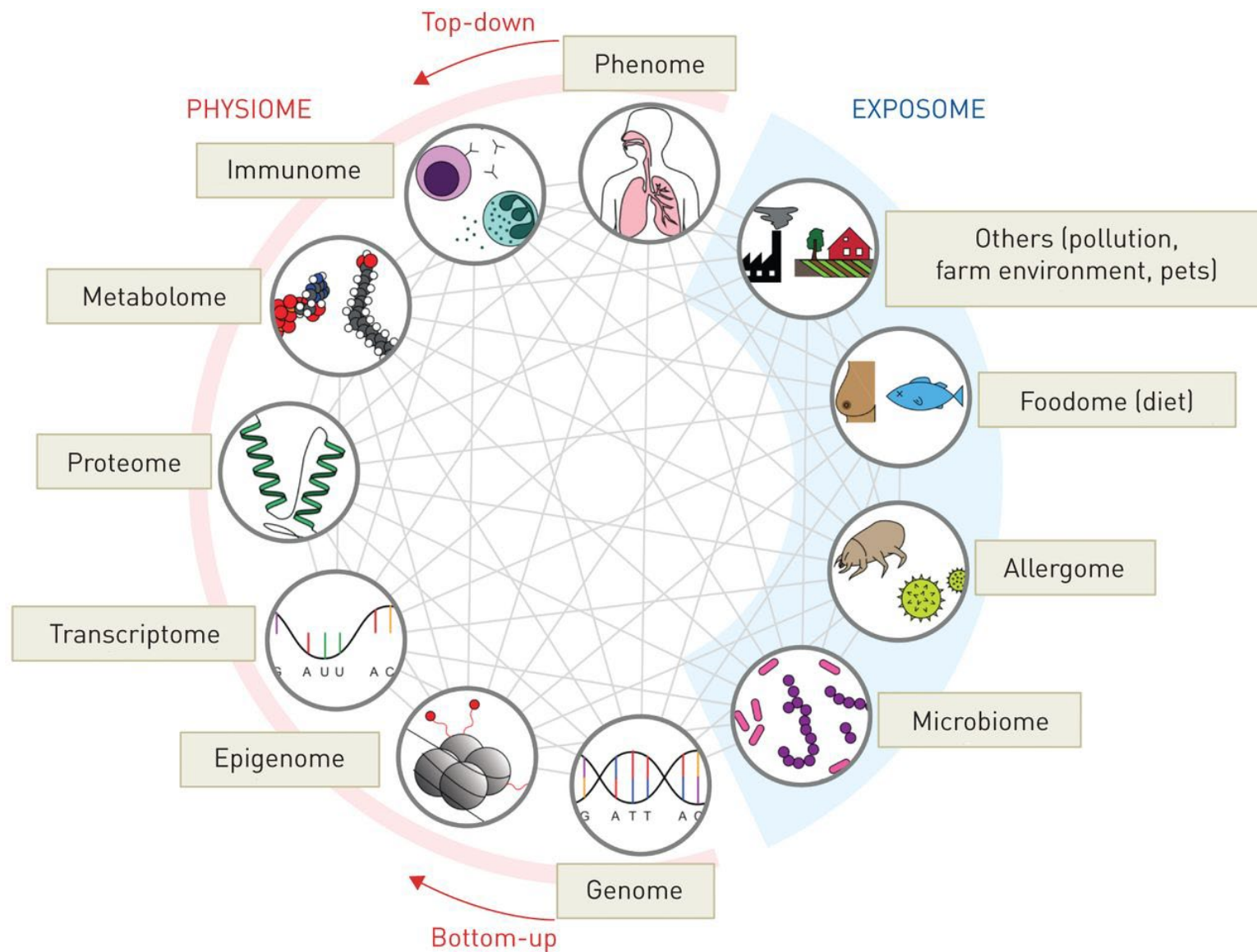
Applicazioni dei Big Data

- Settori che utilizzano Big Data:
 - Sanità,
 - Finanza,
 - Marketing,
 - Industria, etc.

Alcuni esempi

- **Sanità:** Analisi predittiva per diagnosi precoci, gestione delle pandemie, personalizzazione delle cure
- **Finanza:** Rilevamento delle frodi, algoritmi di trading ad alta frequenza, analisi del rischio creditizio
- **Marketing:** Pubblicità mirata basata su dati comportamentali, personalizzazione dell'esperienza utente, analisi delle tendenze di mercato
- **Industria:** Manutenzione predittiva, ottimizzazione delle catene di approvvigionamento, automazione basata su AI

Sanità



Teconologie

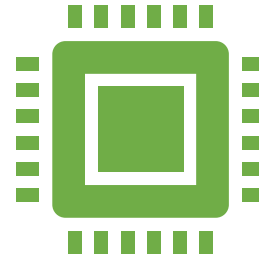
- Hadoop e l'ecosistema (HDFS, MapReduce, YARN)
- Apache Spark e il calcolo distribuito
- Strumenti NoSQL (MongoDB, Cassandra, HBase)

Cosa impareremo?



Impareremo a analizzare diversi tipi di dati:

high dimensional
graph
labeled
text



Impareremo a usare diversi modelli di computazione e tecnologie:

MapReduce
Single machine in-memory
Spark

Che problem risolveremo?



Probelmi della vita reale:

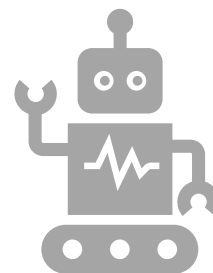
Recommender systems

Market Basket Analysis

Spam detection

Duplicate document detection

..



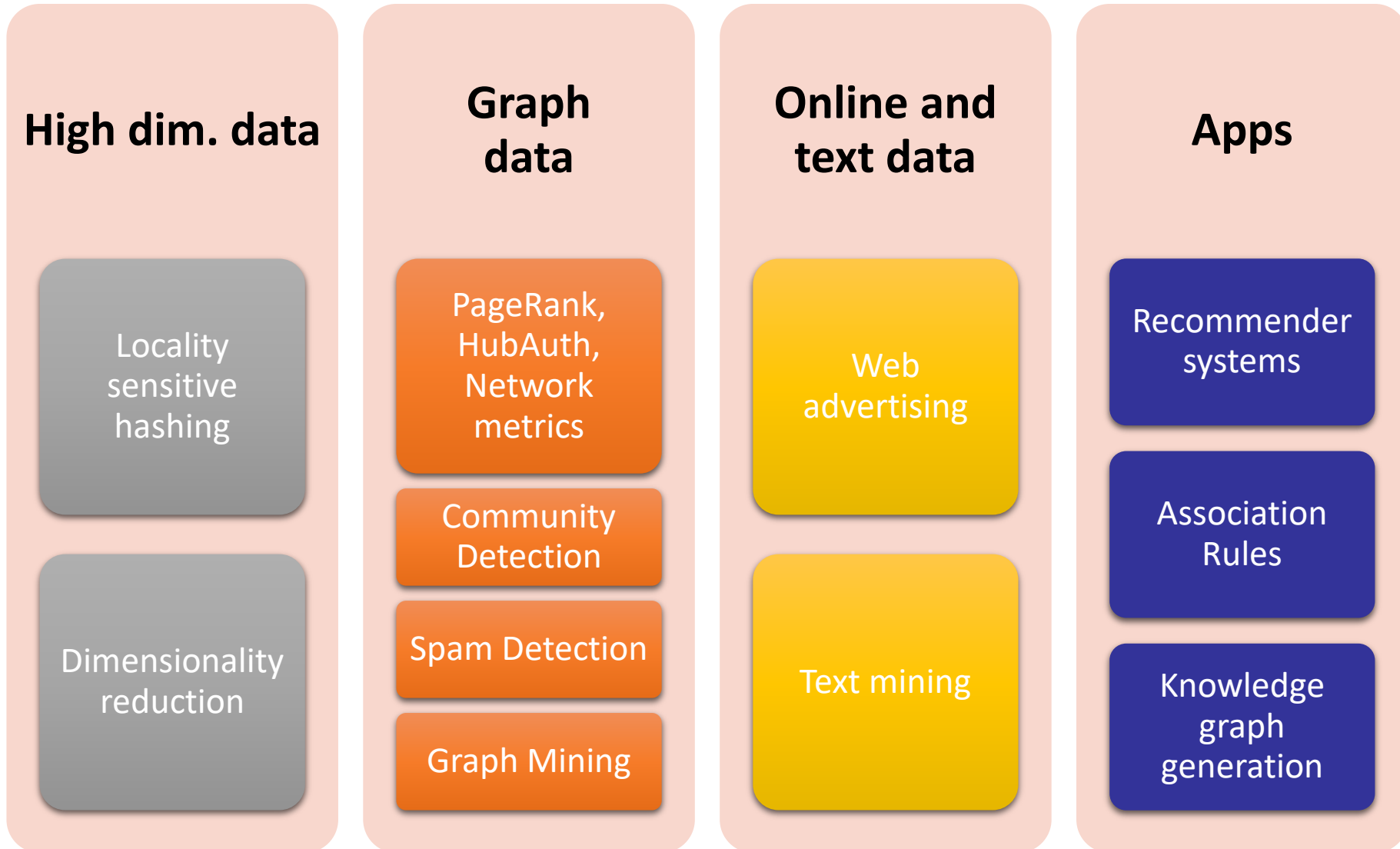
Impareremo diverse “metodologie”:

Linear algebra (SVD, Rec. Sys., Communities)

Optimization (stochastic gradient descent)

Hashing (LSH)

Tutto assieme



Programma /1

Sistemi di raccomandazione:

- Collaborative filtering,
- Modelli a semantica latente
- Network based inference
- Modelli ibridi

Map-Reduce:

- Concetti, motivazioni e algoritmi:
- conteggio parole documenti, prodotto vettore/matrice; Prodotto matrice/matrice; Join; Group By
- Beyond map-reduce: Spark

Programma /2

Ricerca di similarità su alte dimensioni:

- Shingling
- Min-Hashing
- Locality Sensitive Hashing (LSH)
- Min-LSH
- Applicazioni

Dimensionality reduction:

- PCA
- SVD
- CUR
- Proiezioni random e Teorema di Johnson-Lindenstrauss

Link Analysis:

- PageRank e sue estensioni
- Link spam
- Hub-Authorities
- Metriche per l'analisi di reti
- Applicazioni su Map-Reduce

Programma /3

Web Advertising:

- Algoritmi online
- Adword e sue implementazioni

Graph mining e network analysis:

- Network models
- community detection: overlapping communities
- Applicazioni
- Graph Neural Networks
- Knowledge graph generation

Text mining

- TF.IDF, Bag-Of-Word,
- Entity annotation based on AI
- Applicazioni

https://colab.research.google.com/



The screenshot displays the Google Colaboratory (Colab) interface. At the top, the Colab logo is followed by the text "Un benvenuto a Colaboratory". A navigation bar includes links for "File", "Modifica", "Visualizza", "Inserisci", "Runtime", "Strumenti", and "Guida". On the right, there are settings and a "Conc" button. The left sidebar, titled "Sommario", lists sections: "Introduzione", "Data science", "Machine learning", "Altre risorse", and "Esempi in primo piano", with a "+ Sezione" button at the bottom. The main content area features the heading "Ti diamo il benvenuto in Colab" and a sub-heading "(Novità) Prova l'API Gemini". Below these are six links: "Generate a Gemini API key", "Talk to Gemini with the Speech-to-Text API", "Gemini API: Quickstart with Python", "Gemini API code sample", "Compare Gemini with ChatGPT", and "More notebooks". A paragraph at the bottom suggests watching a video for more information on interactive tables, code history, and command palettes.

Un benvenuto a Colaboratory

File Modifica Visualizza Inserisci Runtime Strumenti Guida

Sommario

- Introduzione
- Data science
- Machine learning
- Altre risorse
- Esempi in primo piano

+ Sezione

+ Codice + Testo Copia su Drive

Ti diamo il benvenuto in Colab

(Novità) Prova l'API Gemini

- [Generate a Gemini API key](#)
- [Talk to Gemini with the Speech-to-Text API](#)
- [Gemini API: Quickstart with Python](#)
- [Gemini API code sample](#)
- [Compare Gemini with ChatGPT](#)
- [More notebooks](#)

Se conosci già Colab, guarda questo video per avere informazioni sulle tabelle interattive, sulla visualizzazione della cronologia del codice eseguito e sulla tavolozza dei comandi.

Informazioni

Contatti

- Prof. Alfredo Pulvirenti
 - Stanza 35 terzo blocco, Dipartimento di Matematica e Informatica
 - Tel. 095-7383087
 - e-mail: apulvirenti@dmf.unict.it
 - Homepage: <http://www.dmf.unict.it/~apulvirenti/>
 - Materiale: <http://studium.unict.it>

A large orange circle occupies the left side of the slide, partially cut off by the edge.

Ricevimento

- **Mercoledì' 11-13 (in presenza)**



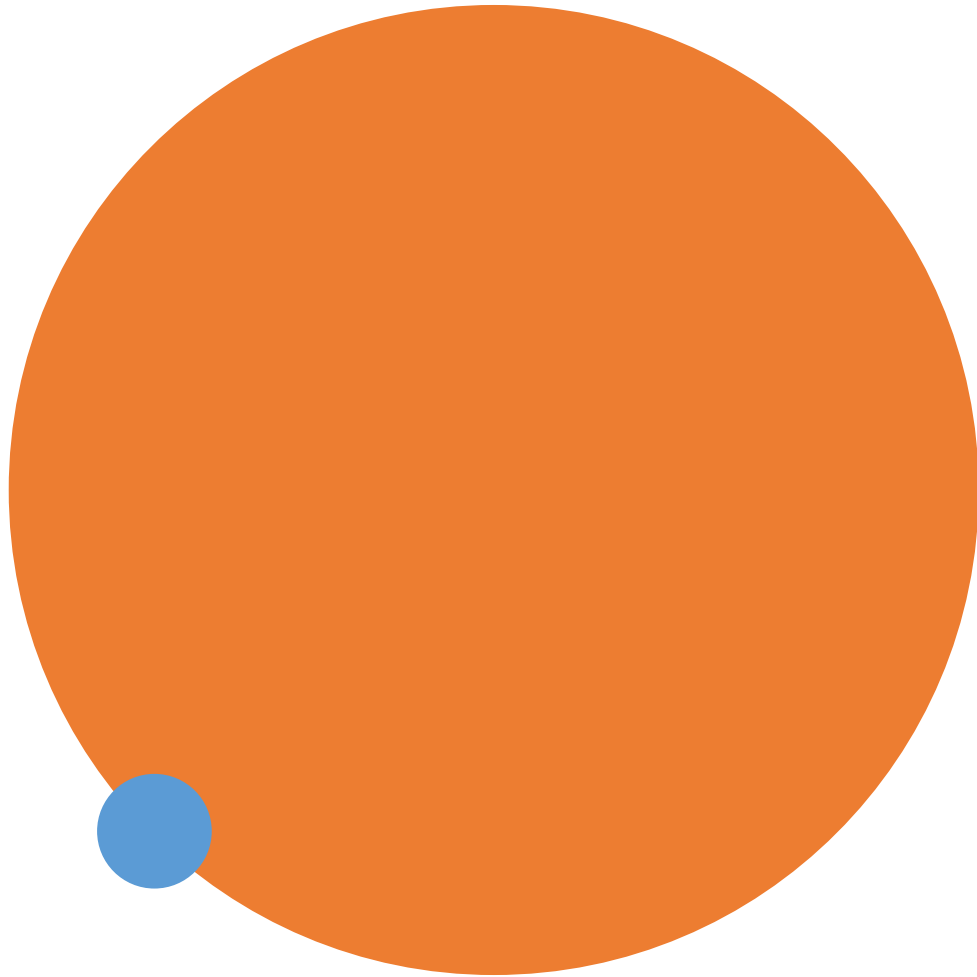
Logistica

- Course website:

<http://studium.unict.it>

- Slide
- Esercizi, soluzioni
- Letture





- Libri
 - **Mining of Massive Datasets**, Leskovec, Rajaraman, Ullman , Free online: <http://www.mmds.org>
 - **Large-Scale Data Analytics with Python and Spark** Isaac Triguero, Mikel Galar
 - **Introduction to Data Mining**, Tan, Steinbach, Kumar, Pearson Ed.
 - **Data Mining: Concepts and Techniques**, Han, Kamber, Morgan Kaufmann Ed.

Comunicazioni

- **Studium**
- **email:**
 - apulvirenti@dmf.unict.it
- **Messaggi sul corso saranno pubblicati su studium**

Prerequisiti

Algoritmi

- Programmazione dinamica, strutture dati

Basic probability

- Momenti, distribuzioni, MLE, ...

Programmazione

- C++/Java/R ecc. saranno utili

Esame

Scritto: 25%

**Progetto: 75% (da consegnare
entro 60 giorni dal superamento
dello scritto)**

**Attività laboratoriale svolta in aula
3 punti extra**