

Introduction to Hadoop MapReduce and Spark

Corso di Sistemi e Architetture per Big Data A.A. 2021/22 Valeria Cardellini

Laurea Magistrale in Ingegneria Informatica

The ML, AI and Data (MAD) landscape

MACHINE LEARING, ARTIFICIAL INTELLIGENCE, AND DATA (MAD) LANDSCAFE 2021								
INFRASTRUCTURE	ANALYTICS	MACHINE LEARNING & ARTIFICIAL INTELLIGENCE	APPLICATIONS – ENTERPRISE					
Source DALLOUIN Source DALLOUIN Source person Source Dallouin Source person Source persource Source person Source person	Sec. Sec. <th< td=""><td></td><td></td></th<>							
- DATA RESOURCES 4.0% DATA RESOURCES -								
	ALTINA PODU/INIS P		DECEMBER Description Description					
Version 3.0 - November 2021	© Matt Turck (@mattturck), John Wu (@john_d_wu) & FirstMark (@first	tmarkcap) mattturck.com/da	ta2021 FIRSTMARK					

https://mattturck.com/data2021/

Zooming on open-source



Valeria Cardellini - SABD 2021/22

Zooming on open-source

ML OPS & INFRA	AI / MACHINE LEARNING / DEEP	LEARNING		
mlflow @Kubeflow	Transf	formers 원 OpenCV 🍱 🐚	tean K Keras BE	RT XGBoost Caffe
Pachyderm	Microsoft DM OpenAI OP	yTorch Lightning theano	de Apache SINGA DIMSU	M FeatureFu VELES
DC SELDON snorkel	mxnet neon™ chainer M.Mchalangelo	🖗 ONNX 🐗 🕮 🖽 🙀 🤝		🙀 DL4J 🛞 МАНОUТ
Polyanon 👔 BENTOML 🔠 MediaPipe	Aerosolve fast.ai Imir CopenML	mindsdb spaCy &Kubeflow Allen	NLP Catilioost	
SEARCH		VISUALIZATION	COLLABORATION 7	SECURITY ——
elasticsearch Solr [®]	elasticsearch 🔣 kibana 🔊 SENTRY	Superset matpl%tlib	Beake	Apache Ranger KNOX
<i>■Lucere Sphinx</i>	elogstash OPrometheus Stuentoit	Metabase redash TensorBoard		Sentry accumulo
🔟 🧊 Sonic	Sfluentd OGrafana		Zeppelin	44
Toshi Search	Den Telemetry	📾 seaborn bokeh 🔛	ANACONDA	iii snyk

The basic Hadoop ecosystem



See <u>https://hadoopecosystemtable.github.io</u> for more products

Valeria Cardellini - SABD 2021/22

The reference Big Data stack





Valeria Cardellini - SABD 2021/22

Parallel programming: background

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





Parallel programming: background

• Simplest environment for parallel programming

- Master/worker architecture

- Master
 - Gets data and splits it into chunks according to the number of workers
 - Sends each worker equal number of chunks
 - Receives results from each worker
- Workers:
 - Receive some chunks of data from master
 - Perform processing
 - Send back results to master

Valeria Cardellini - SABD 2021/22

Parallel programming: background

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



- Estimation algorithm for calculating π

 Relies on Monte Carlo method
- Let's first consider the sequential version of the algorithm
- Then, how to realize a parallel and faster version

Valeria Cardellini - SABD 2021/22

Example: Pi estimation

 By definition, π is the area of a circle with radius equal to 1



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

Valeria Cardellini - SABD 2021/22

Example: Pi estimation

• In formulas:



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

Example: Pi estimation



Total Number of points: 249 Points within circle: 185 Pi estimation: 2.97189

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

Valeria Cardellini - SABD 2021/22

14

Example: Pi estimation



Total Number of points: 4163 Points within circle: 3259 Pi estimation: 3.13140

Example: Pi estimation



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

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

Valeria Cardellini - SABD 2021/22

Example: Pi estimation

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

Key idea behind MapReduce and Spark: Divide and conquer

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



Valeria Cardellini - SABD 2021/22

Divide and conquer: how?

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

Key idea behind MapReduce and Spark: scale out, not up!

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

Valeria Cardellini - SABD 2021/22