

# Introduzione al calcolo parallelo con Apache Spark



Vincenzo Manzoni, Andrea Rota  
Data Science

Leader mondiale nella produzione e fornitura di prodotti tubolari e servizi per:

- Trivellazioni, estrazione e produzione di petrolio e gas
- Trasporto di petrolio e gas
- Impianti di trasformazione e centrali elettriche
- Applicazioni specialistiche industriali e automotive



**OCTG**



**Giunti  
Premium**



**Linepipe per  
applicazioni  
onshore e  
offshore**



**Trasformazione  
di idrocarburi**



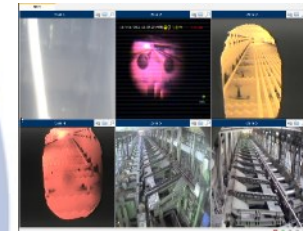
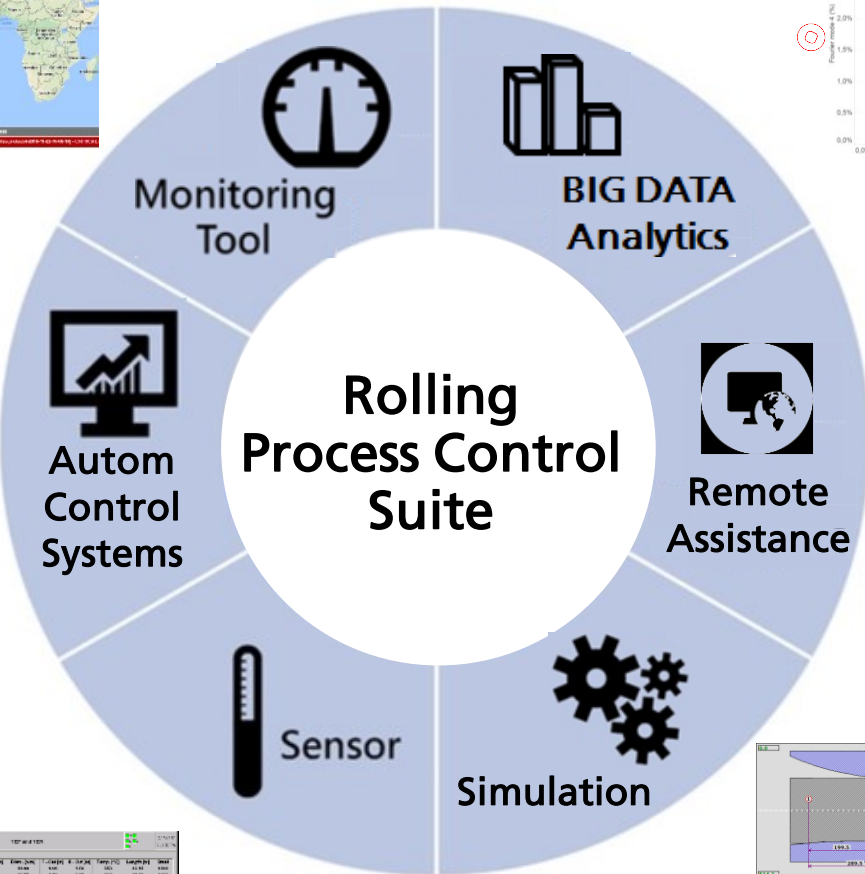
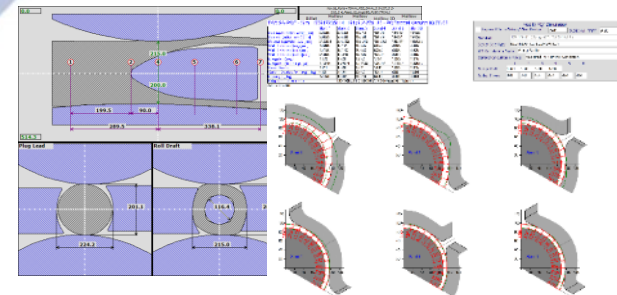
**Generazione di  
energia**



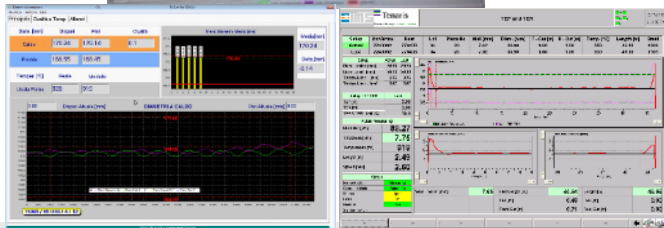
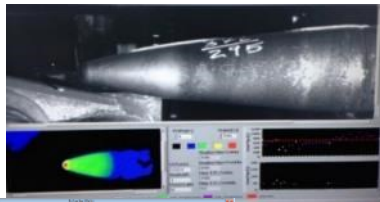
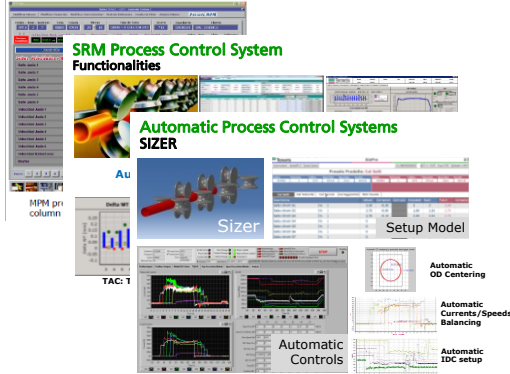
**Applicazioni  
industriali e  
automotive**

[illegible]

# Industry 4.0 topics in Tenaris

Diagnostic System: SisTec MPM



Vincenzo Manzoni

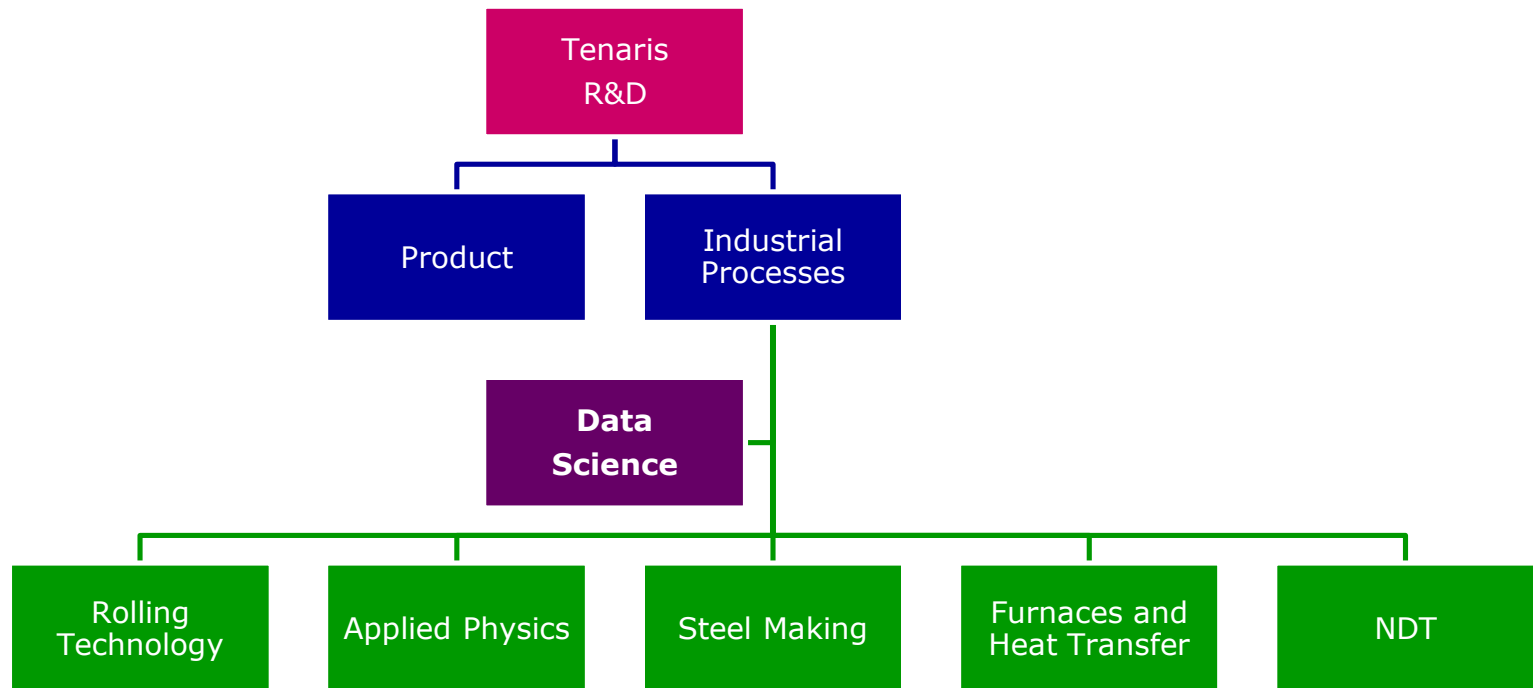
Tenaris



11 December 2017

4

# Chi siamo



# Big Data and AI examples



## Massive Data Analysis

**Offline** applications for **batch** decision making.



## Process Monitoring

**Streaming** application for **real-time** decision making.



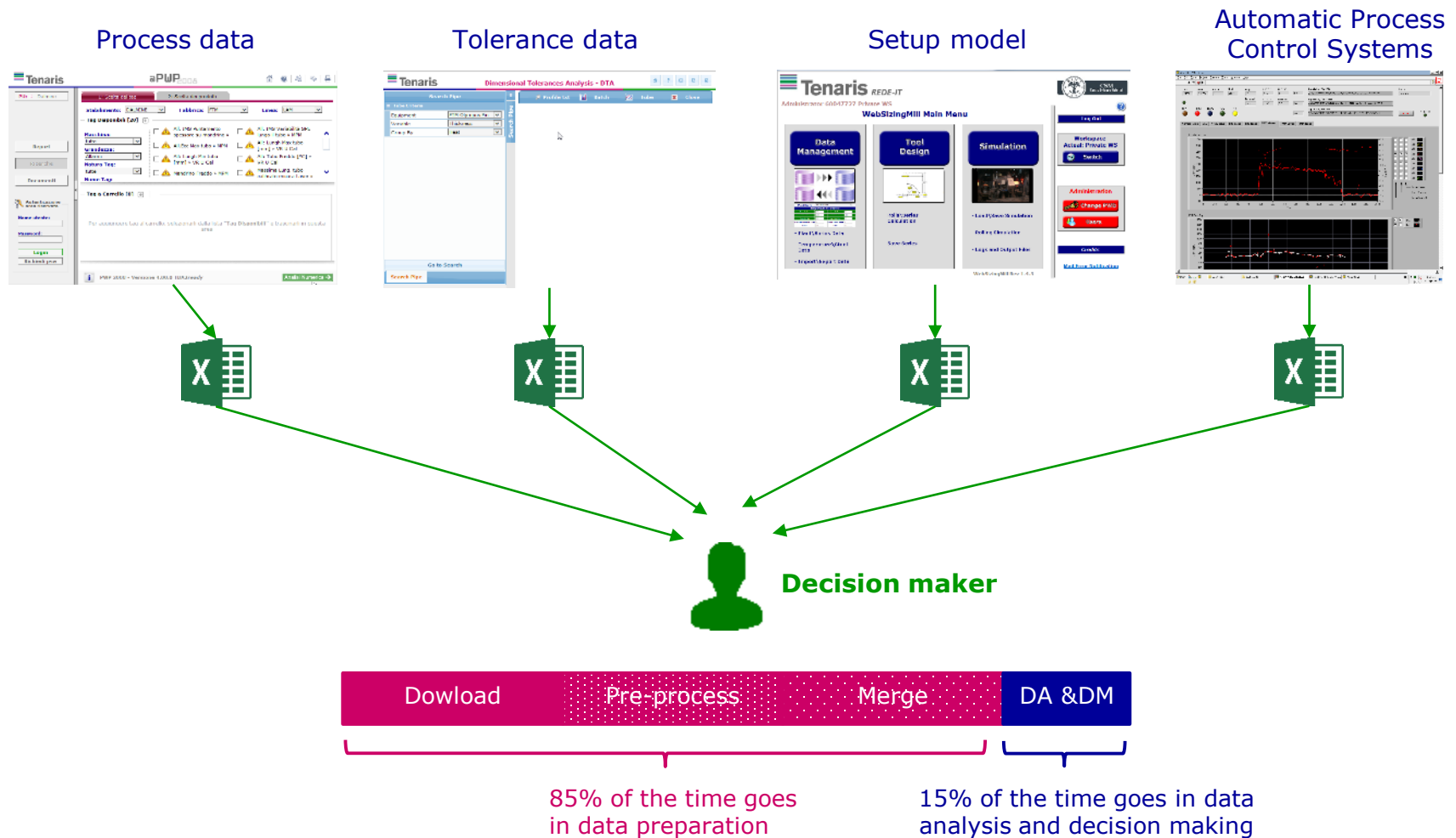
## Artificial Intelligence and Machine Learning Modeling

Model Industrial Processes behavior from data through Artificial Intelligence and Machine Learning algorithms.



# MASSIVE DATA ANALYSIS

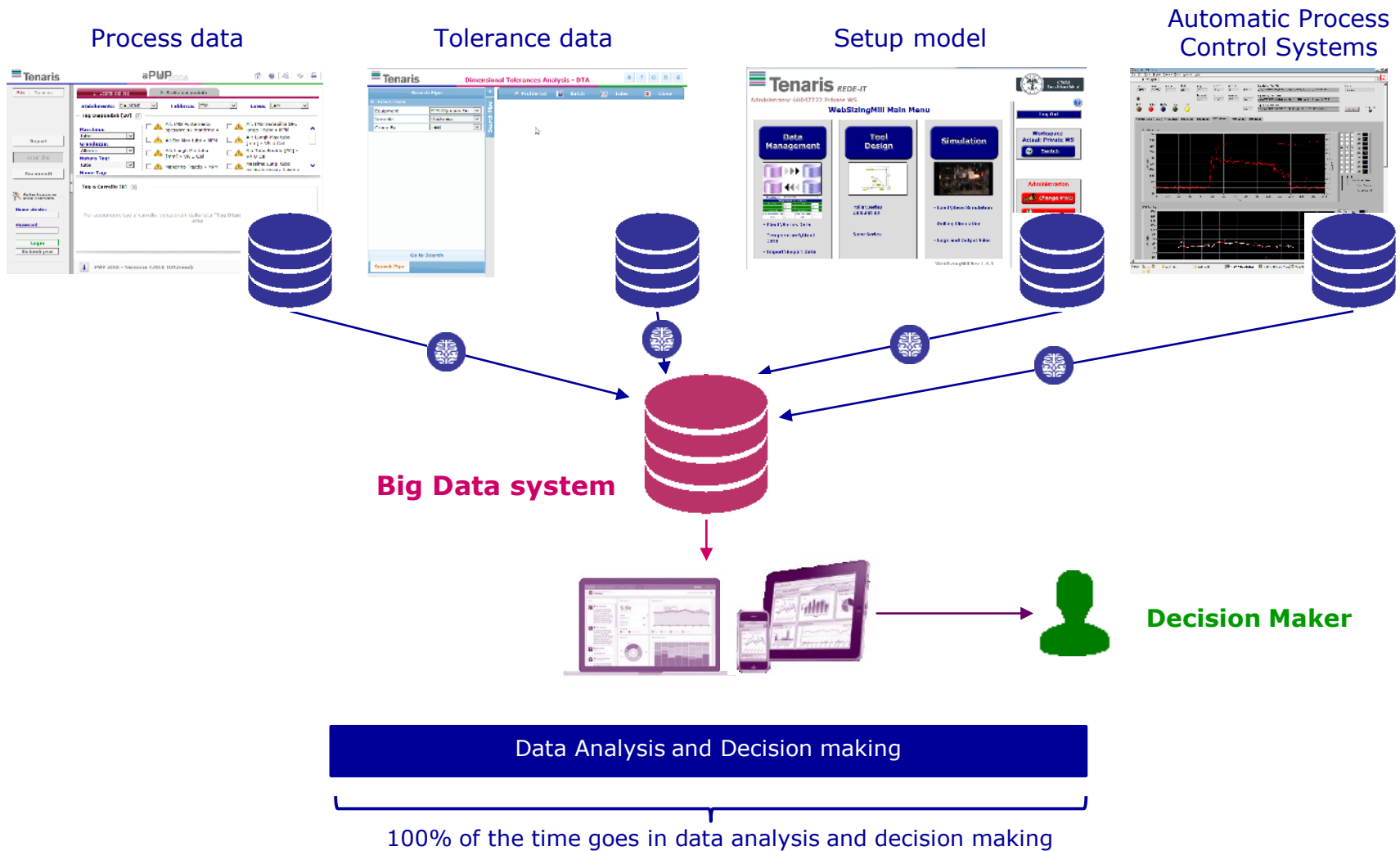
# Processo decisionale tradizionale



Munson, "A Study on the Importance of and Time Spent on Different Modeling Steps" ([link](#))



# Processo decisionale supportato da tecnologie Big Data



# Perché Big Data



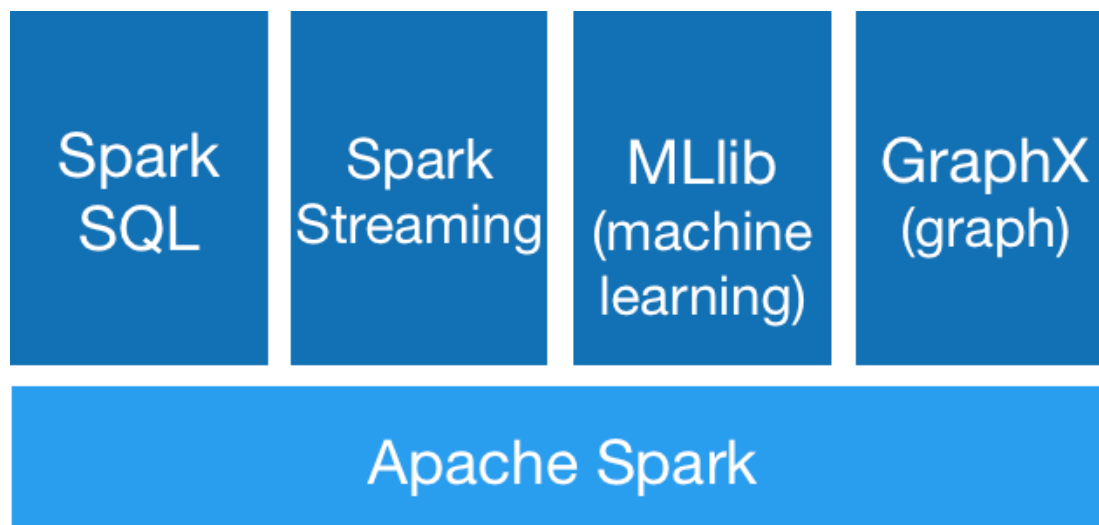
- La **velocità** con cui vengono prodotti i dati cresce più velocemente della capacità di una singola macchina per elaborarli.
- Il **volume** dei dati cresce più velocemente della capacità di storage di una singola macchina.

**Serve una tecnologia per esprimere  
ad alto livello analisi in parallelo su  
grandi quantità di dati.**

Framework open-source sviluppato dall'AMPLab dell'Università di Berkeley.

A differenza di paradigmi precedenti – tipo Map-Reduce – Spark usa primitive in memoria (fino a 100x in termini di prestazioni).

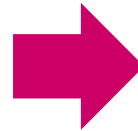
Supporta diversi linguaggi di programmazione, tra cui Java, Scala e Python.



# Contare le parole di un testo



«I am Sam  
I am Sam  
Sam I am  
Do you like  
Green eggs and ham  
...»



Parola	Occorrenze
I	3
Am	3
Sam	3
Do	1
You	1
Like	1
...	...

# Contare le parole di un testo



«I am Sam  
I am Sam  
Sam I am  
Do you like  
Green eggs and ham  
...»

Parola	Occorrenze
I	1

# Contare le parole di un testo



«I am Sam  
I am Sam  
Sam I am  
Do you like  
Green eggs and ham  
...»

Parola	Occorrenze
I	1
Am	1

# Contare le parole di un testo



«I am Sam  
I am Sam  
Sam I am  
Do you like  
Green eggs and ham  
...»

Parola	Occorrenze
I	1
Am	1
Sam	1

# Contare le parole di un testo



«I am Sam  
I am Sam  
Sam I am  
Do you like  
Green eggs and ham  
...»

Parola	Occorrenze
I	2
Am	1
Sam	1

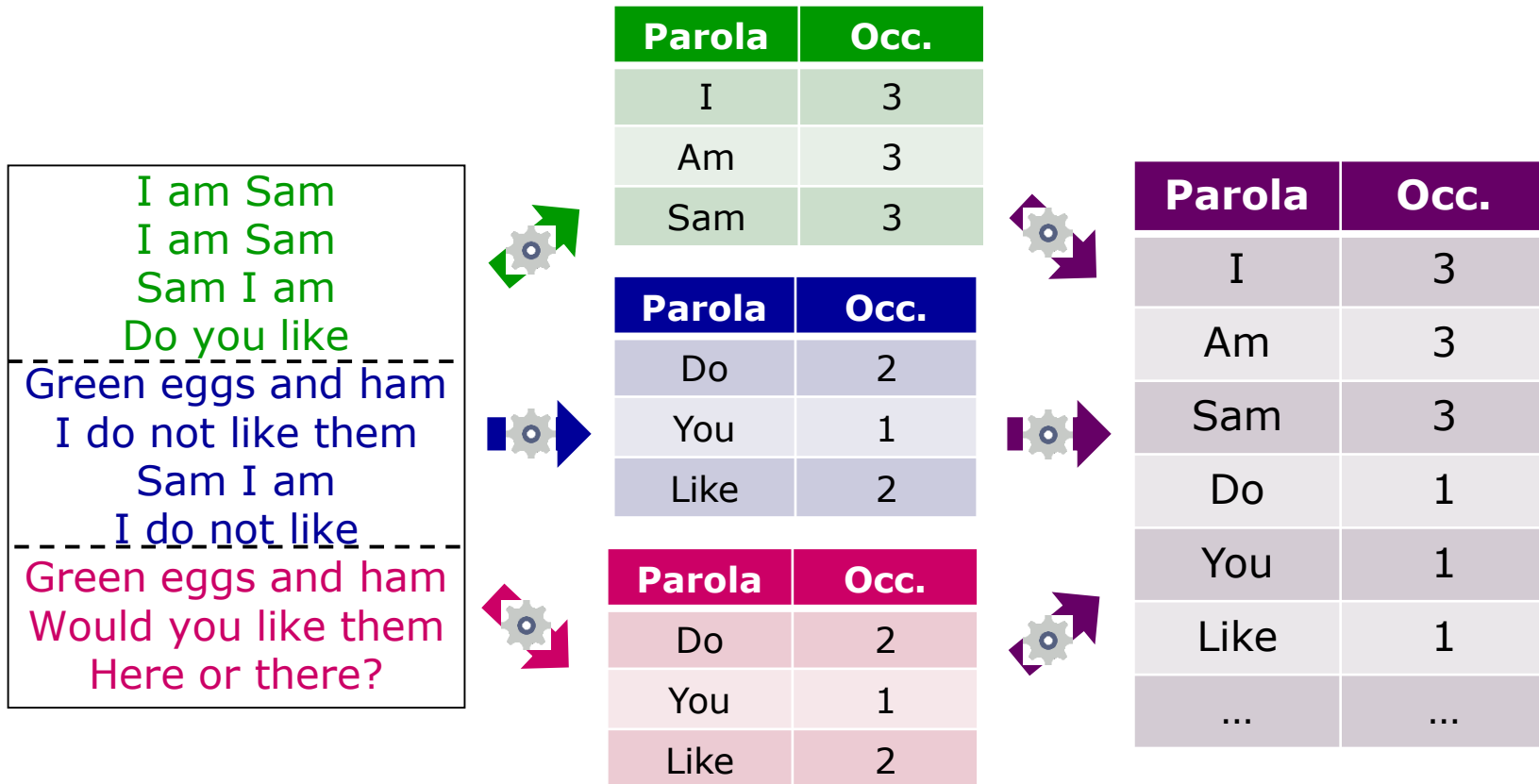


# Contare le parole di un testo (molto) lungo (1/3)



I am Sam  
I am Sam  
Sam I am  
Do you like  
Green eggs and ham  
I do not like them  
Sam I am  
I do not like  
Green eggs and ham  
Would you like them  
Here or there?

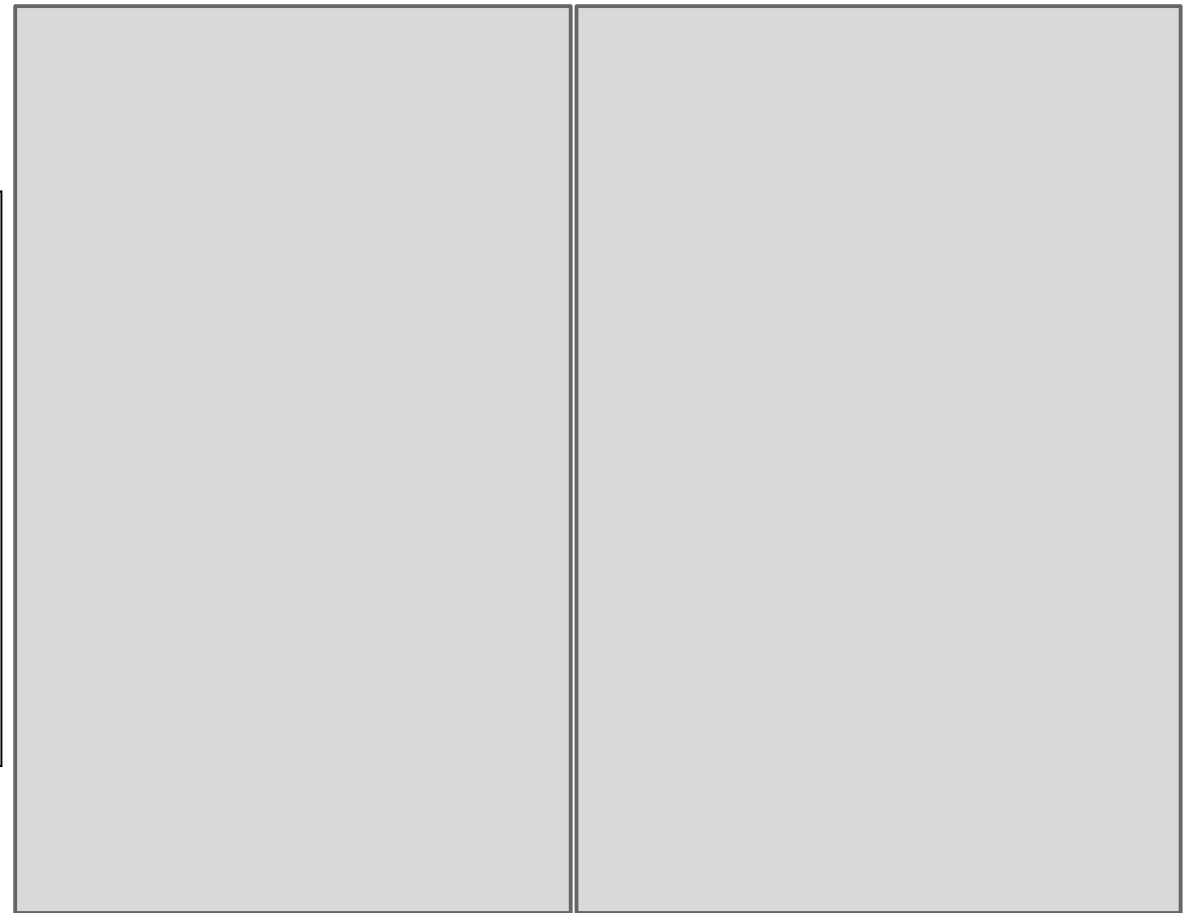
# Contare le parole di un testo (molto) lungo (2/3)



# Contare le parole di un testo (molto) lungo (3/3)



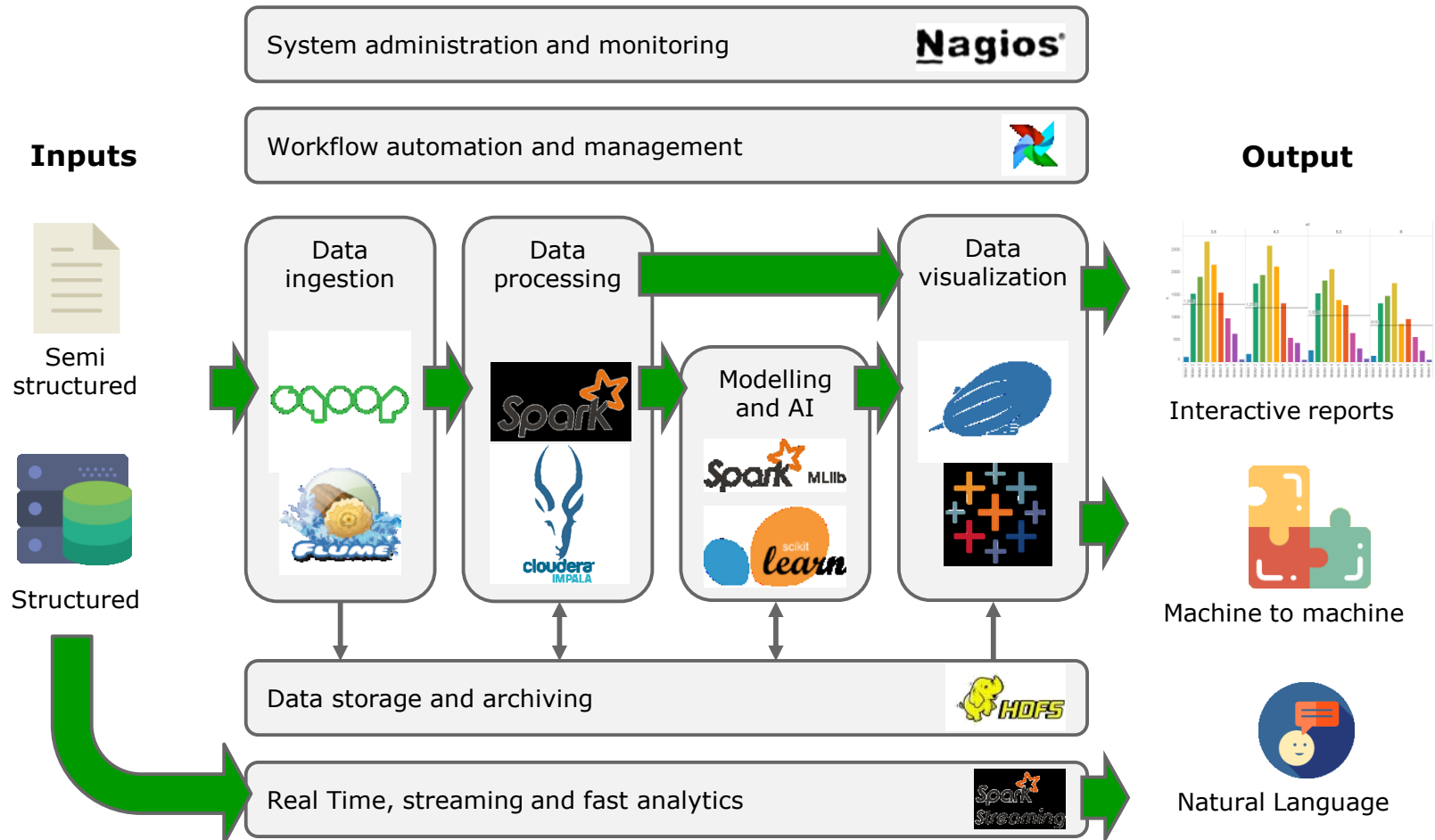
I am Sam  
I am Sam  
Sam I am  
Do you like  
-----  
Green eggs and ham  
I do not like them  
Sam I am  
-----  
I do not like  
-----  
Green eggs and ham  
Would you like them  
Here or there?



Map

Reduce

# Tenaris Big Data technology stack



# Tenaris Big Data technologies



## October 17, 2016 - Apache Flume 1.7.0 Released

The Apache Flume team is pleased to announce the release of Flume 1.7.0.

Flume is a distributed, reliable, and available service for efficiently collecting, aggregating, and moving large amount

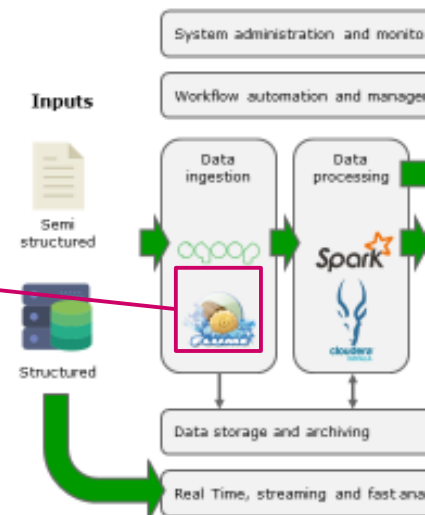
Version 1.7.0 is the tenth Flume release as an Apache top-level project. Flume 1.7.0 is stable, production-ready software versions of the Flume 1.x codeline.

Several months of active development went into this release: almost 100 patches were committed since 1.6.0, representing While the full change log can be found on the 1.7.0 release page (link below), here are a few new feature highlights:

- Taildir source
- Kafka integration improvements (eg. security)

Below is the list of people (from Git/SVN logs) who submitted and/or reviewed improvements to Flume during the 1.7.0 development cycle:

- Abraham Fine
- Alexandre Dutra
- **Andrea Rota**
- Ashish Paliwal
- Attila Simon
- Bessenyei Balázs Donát
- Daniel Templeton
- Deepesh Khandelwal





# DEMO TIME!

# Setup Databricks

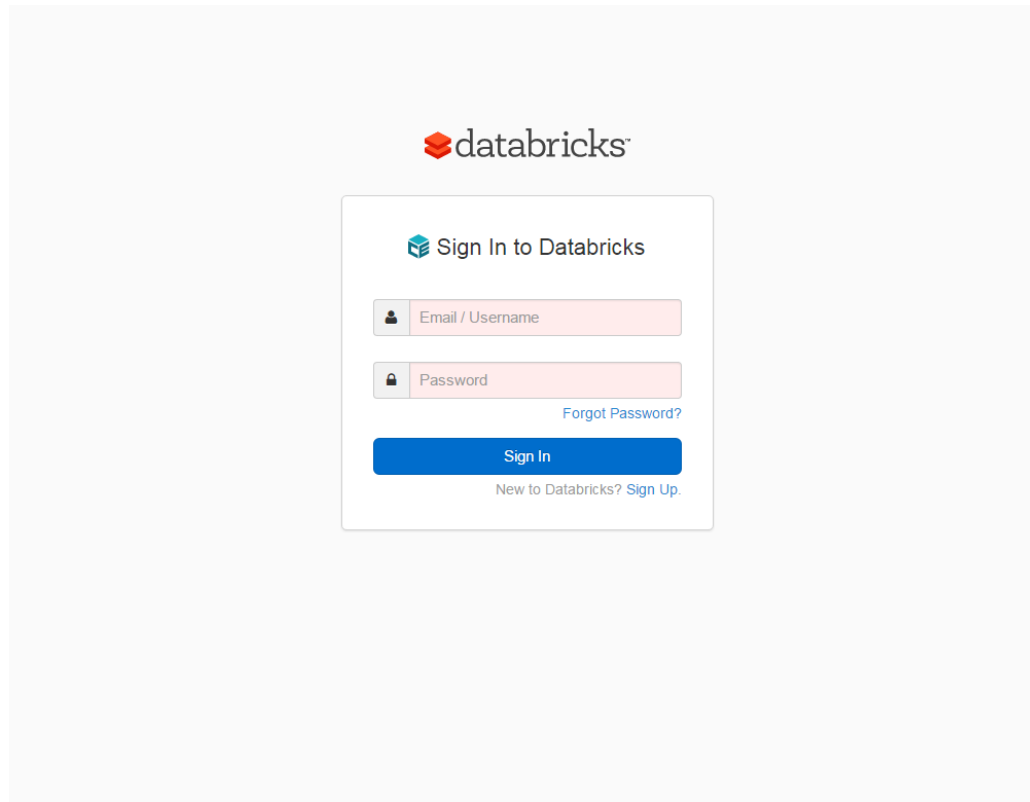


<https://databricks.com/try-databricks>

# Login databricks



<https://community.cloud.databricks.com/>



The screenshot shows the Databricks login interface. At the top is the Databricks logo. Below it is a 'Sign In to Databricks' section. This section contains two input fields: 'Email / Username' and 'Password'. Below the password field is a link for 'Forgot Password?'. A blue 'Sign In' button is positioned below the input fields. At the bottom of the sign-in section is a link for 'New to Databricks? Sign Up.'.



# GitHub repository



<https://github.com/tenaris/scala-spark-workshop>

# Join the dark side!



Siamo alla ricerca di studenti del 5 anno per tesi, tirocini e assegni di ricerca su tematiche **Big Data**, **Machine Learning** e **Artificial Intelligence**.

Per candidarvi, scrivete una mail a:

Vincenzo Manzoni:  
[vmanzoni@tenaris.com](mailto:vmanzoni@tenaris.com)

Andrea Rota  
[arota@tenaris.com](mailto:arota@tenaris.com)

