# Processamento de dados em Big Data

## Parte 2 - Frameworks

### Luiz Henrique Zambom Santana, D.Sc.
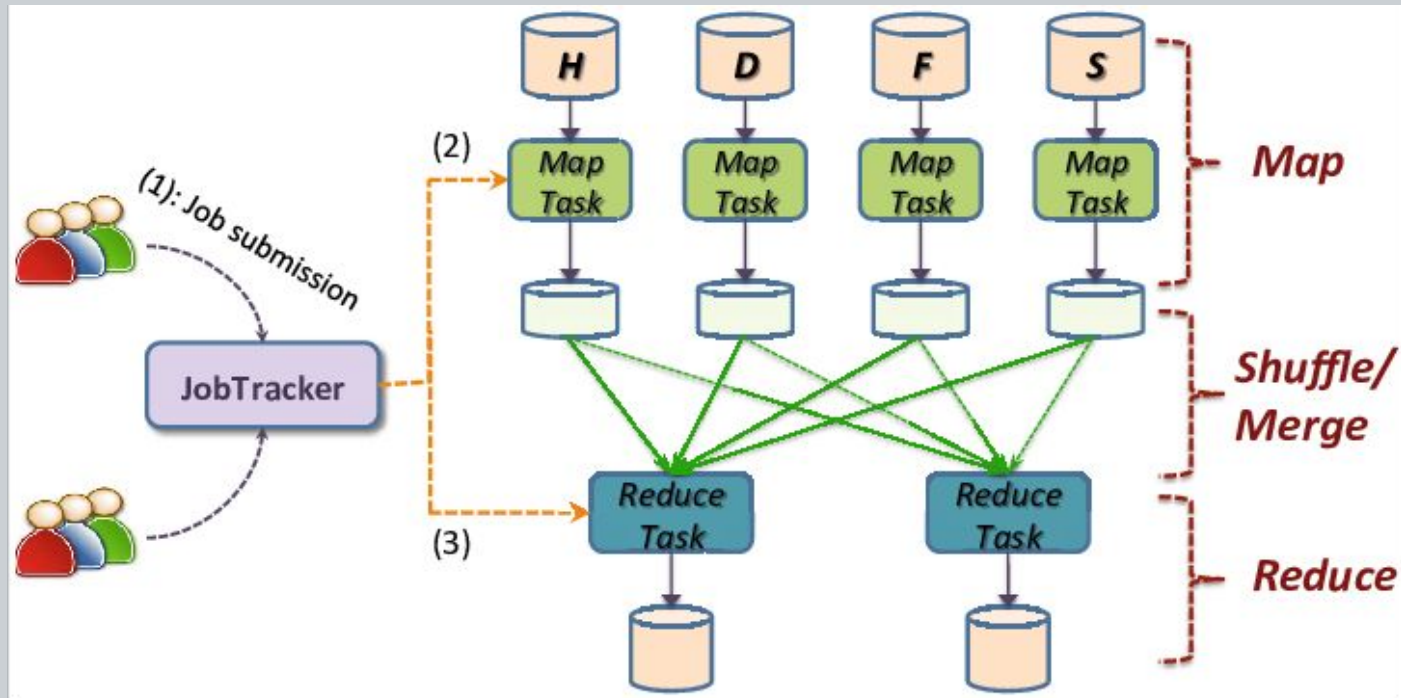
INE | CTC

UNIVERSIDADE FEDERAL
DE SANTA CATARINA

# Agenda

- Processamento:
  - Hadoop
  - Apache Spark
    - Core
    - Streaming
    - Machine learning
- Streaming: Apache Kafka e NiFi
- Orquestração: Apache Camel e Apache AirFlow
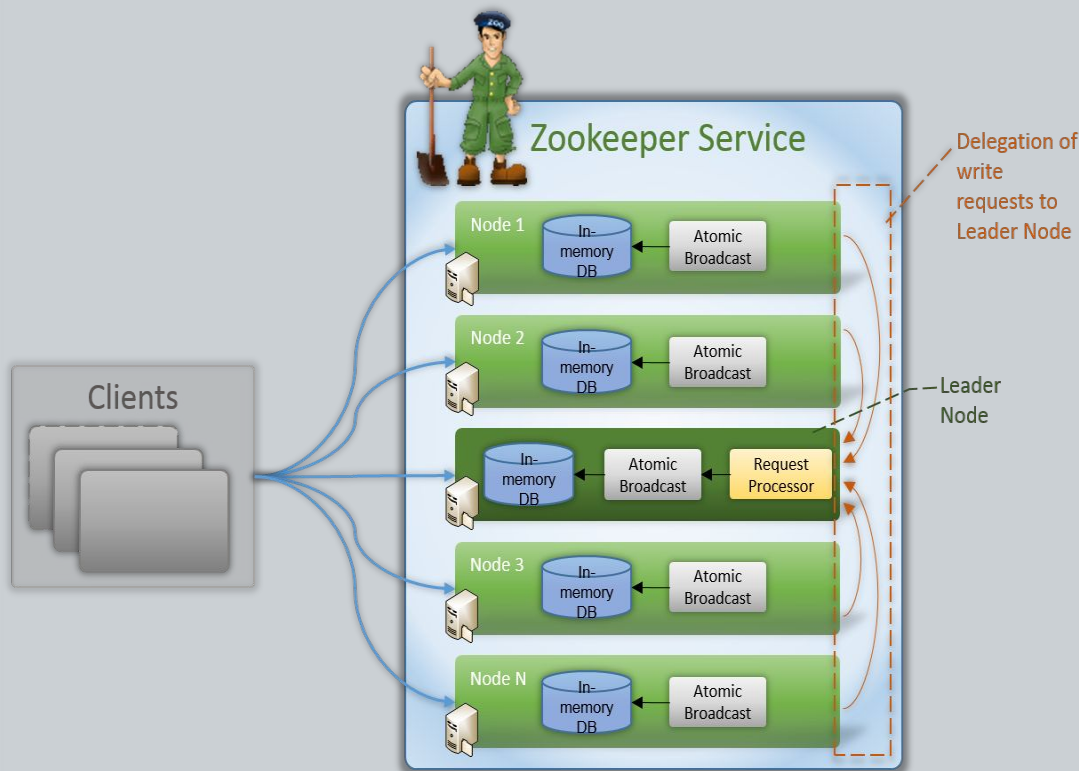- Data Lake: Apache Hudi, Delta Lake

# Apache Hadoop

# Apache Hadoop

- O Hadoop oferece um ecosistema completo:
  - Zookeeper
  - Hive
  - Pig
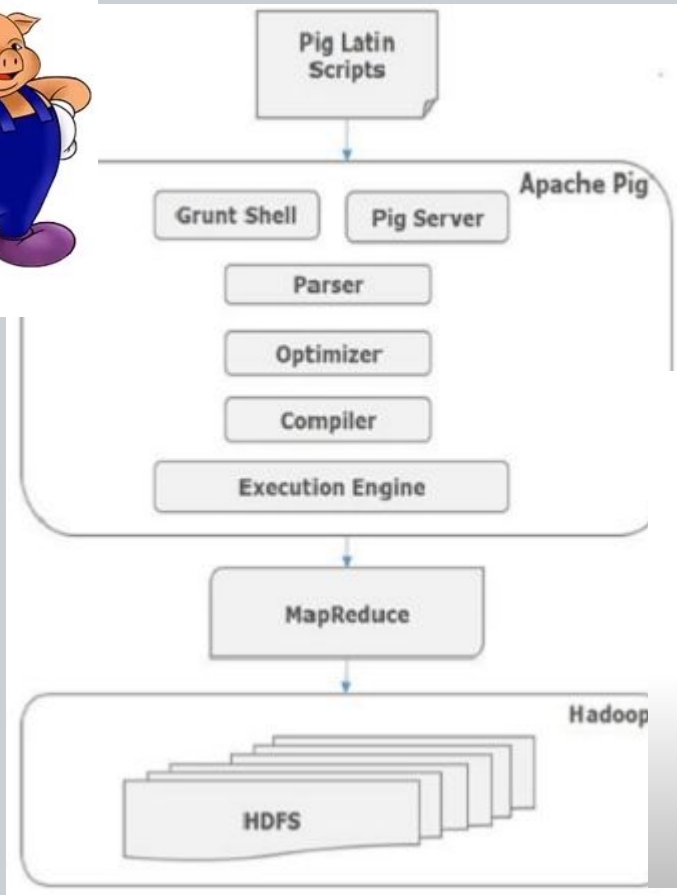  - HDFS

# Apache Hadoop - Zookeeper



- Serviço de nomeação
- Gerenciamento de configuração:
- Gerenciamento de cluster
- Eleição do líder
- Serviço de bloqueio e sincronização
- Registro de dados altamente confiável
- Kafka, Hadoop, Spark, NiFi…

# Apache Hadoop - Hive

- HQL: Hive Query Language
- Traduzido para jobs MapReduce no Hadoop
- OLAP, mas não suporta OLTP
- Funcionalidades
  - Suporta vários formatos, incluindo ORC
  - Usa vários tipos de compressão
  - Joins especiais para aumentar o desempenho
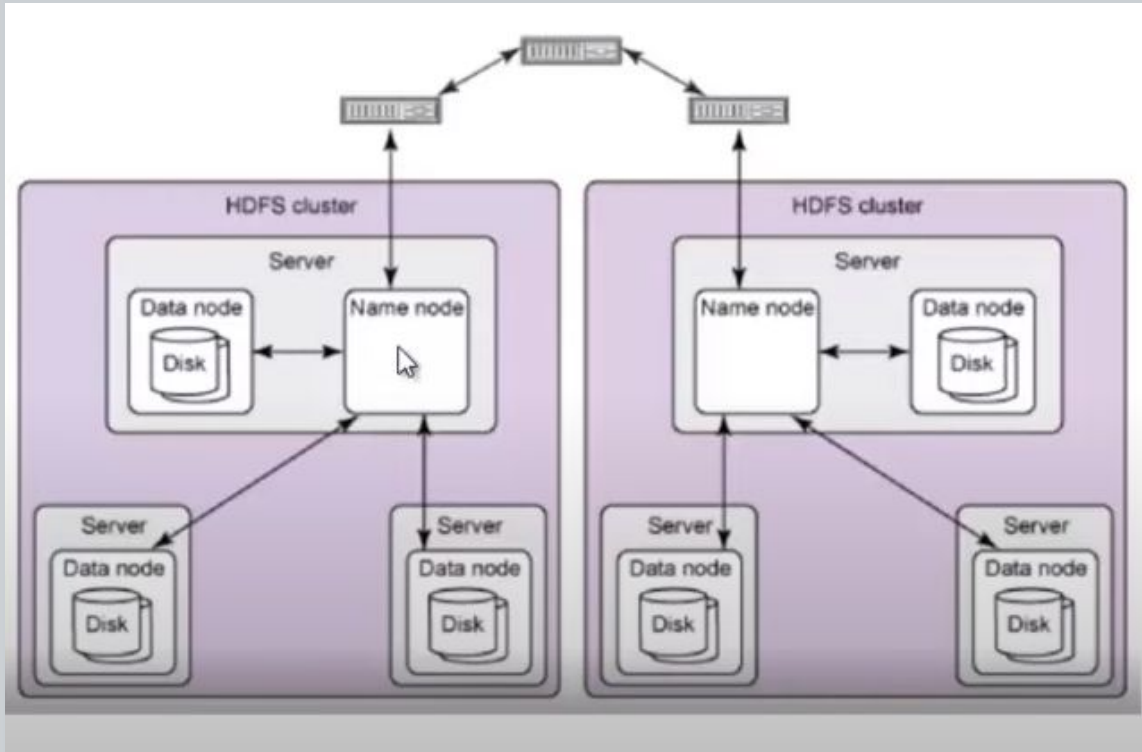  - Schema on Read

# Apache Hadoop - Pig



- Script para acessar Hadoop e fazer análise de dados
- Bom para dados desestruturados

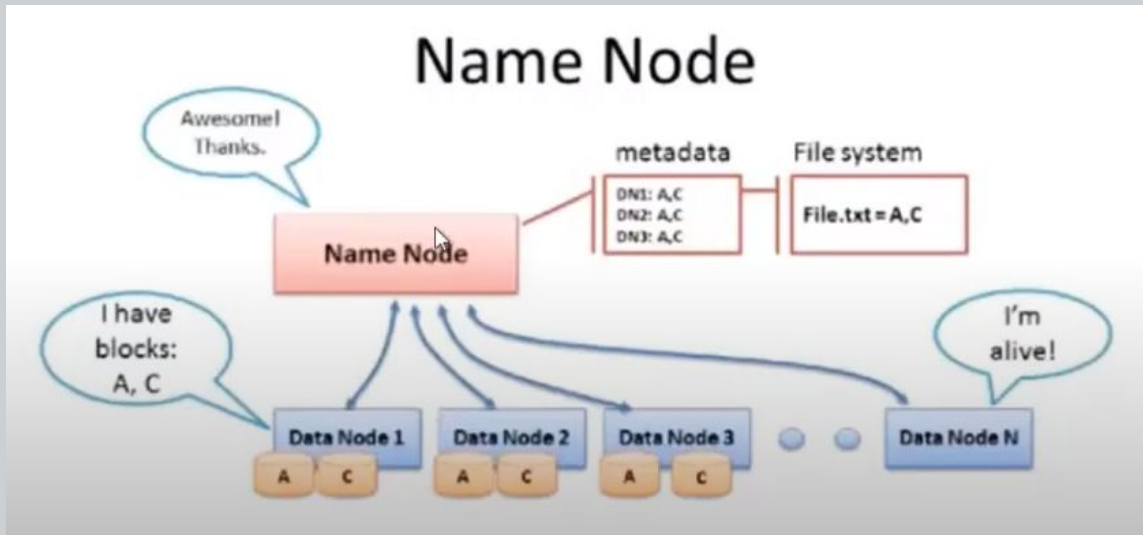| SQL | Pig |
|-----|-----|
| SELECT c_id , **SUM**(amount) AS CTotal | customer = **LOAD** '/data/customer.dat' **AS** (c_id,name,city); |
| **FROM** customers c | sales = **LOAD** '/data/sales.dat' **AS** (s_id,c_id,date,amount); |
| **JOIN** sales s ON c.c_id = s.c_id | salesBLR = **FILTER** customer **BY** city == 'Texas'; |
| **WHERE** c.city = 'Texas' | joined= **JOIN** customer **BY** c_id, salesTX **BY** c_id; |
| **GROUP BY** c_id | grouped = **GROUP** joined **BY** c_id; |
| **HAVING SUM**(amount) > 2000 | summed= **FOREACH** grouped **GENERATE GROUP**, **SUM**(joined.salesTX::amount); |
| **ORDER BY** CTotal **DESC** | spenders= **FILTER** summed **BY** $1 > 2000; |
| | sorted = **ORDER** spenders **BY** $1 **DESC**; |
| | **DUMP** sorted; |

# Apache Hadoop - HDFS



- Hadoop File System
- Sistema de arquivos distribuídos
- Forma uma abstração
- Baseado no Google File System

https://www.youtube.com/watch?v=Z4htZMwIfDs

# Apache Hadoop - HDFS



- Um grande índice de metadados

https://www.youtube.com/watch?v=
Z4htZMwlfDs
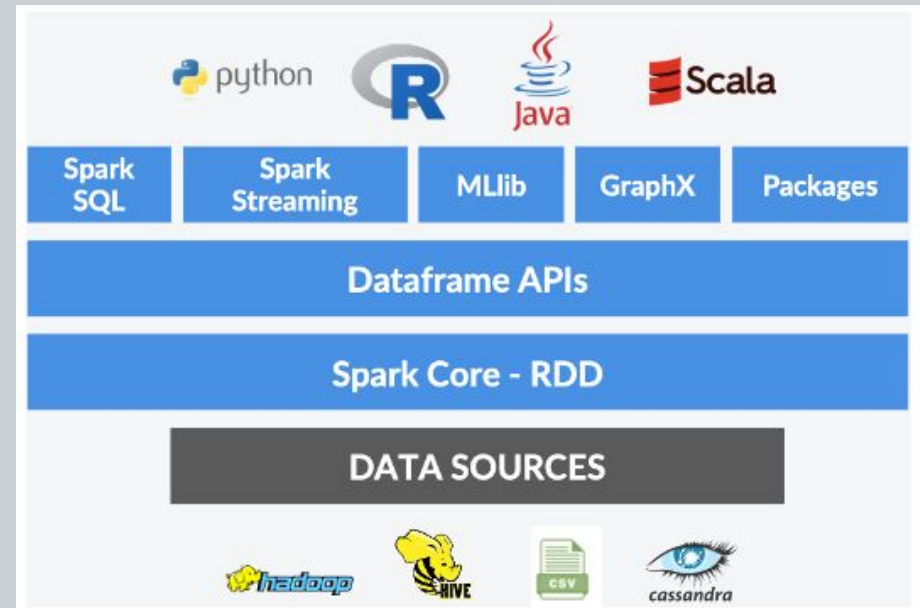
# Apache Hadoop - HDFS



- Utiliza sempre blocos sequenciais
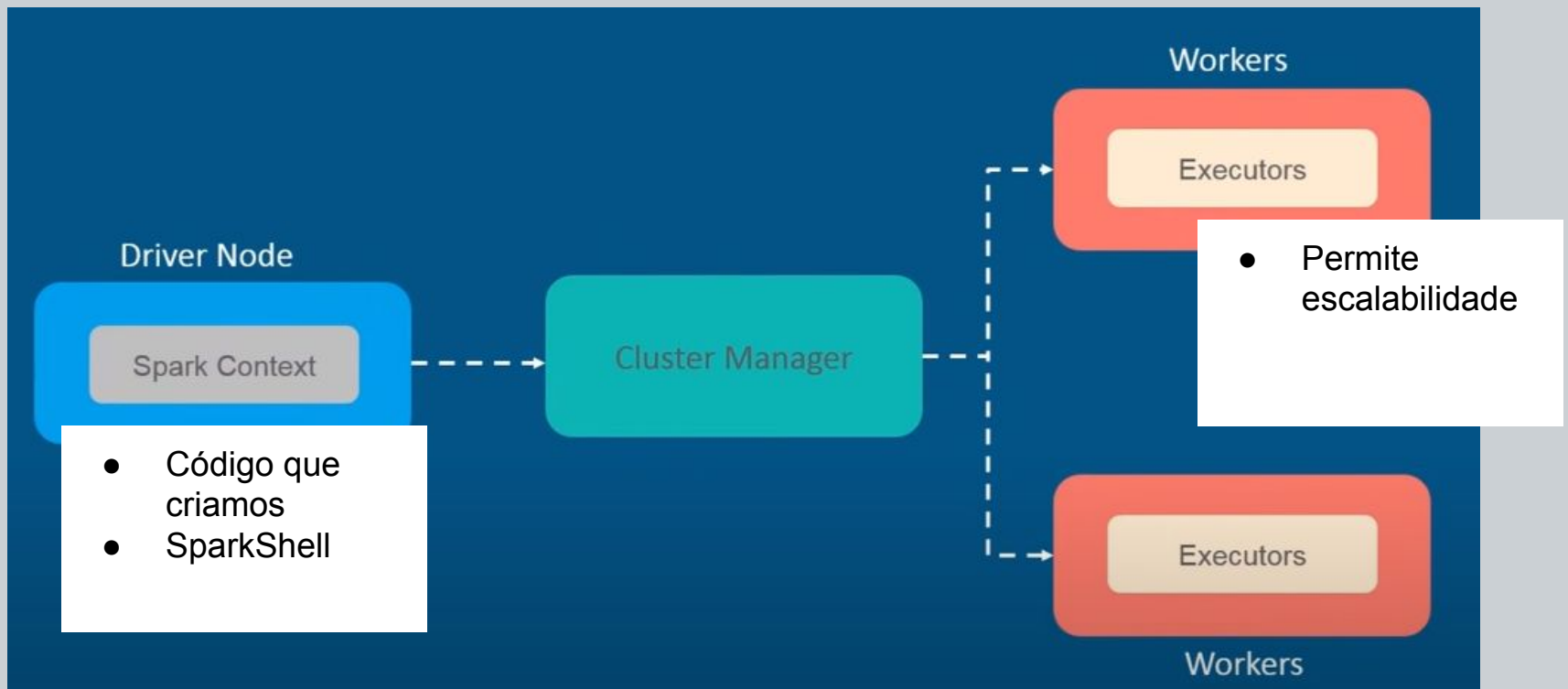- Replicação de blocos

https://bigishere.wordpress.com/2016/08/03/configuring-replication-factor-and-block-size-in-hdfs/

# Apache Spark

- Apache Spark
  - Core
    - RDD
      - Transformation
      - Actions
    - DAG
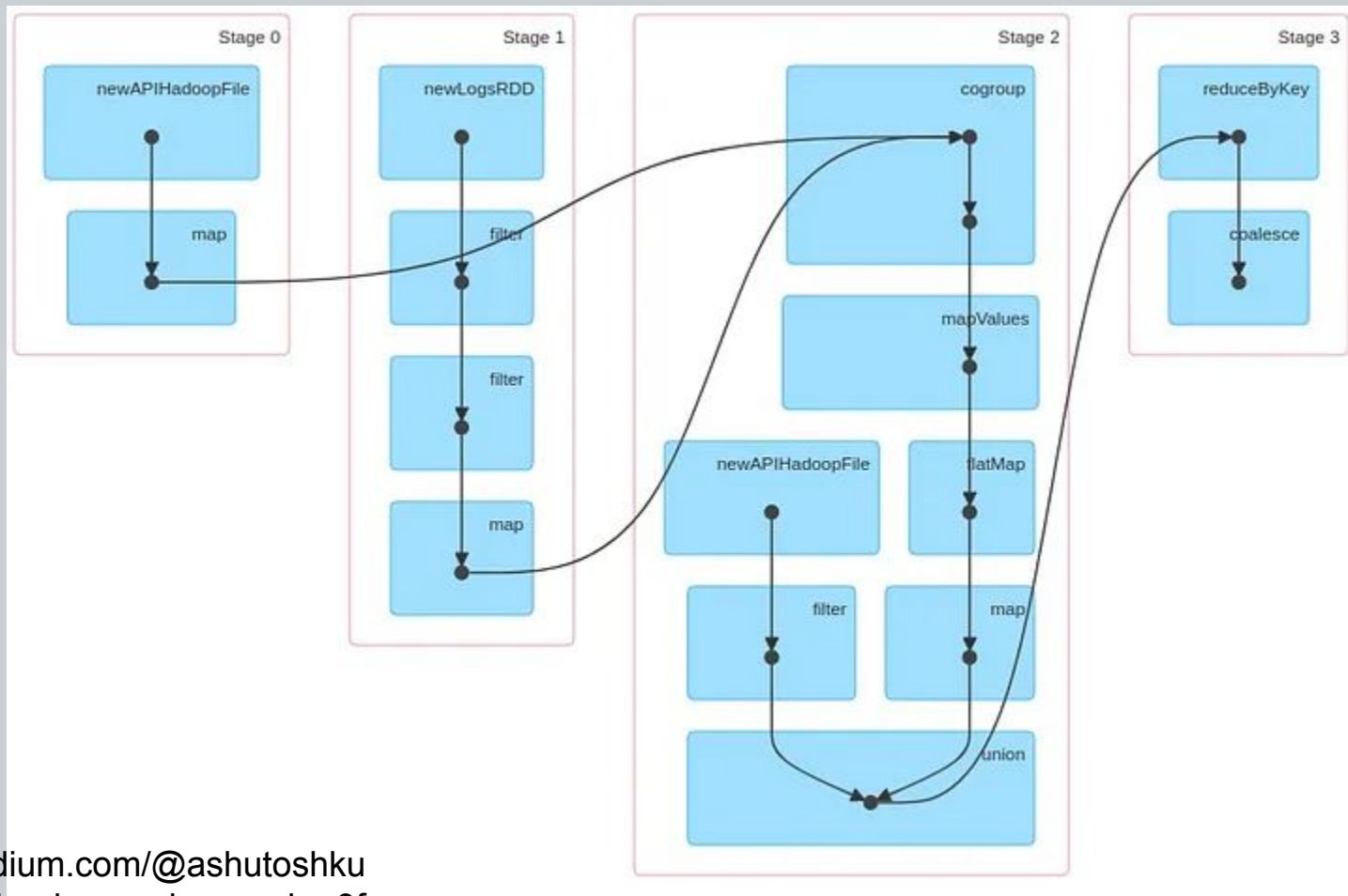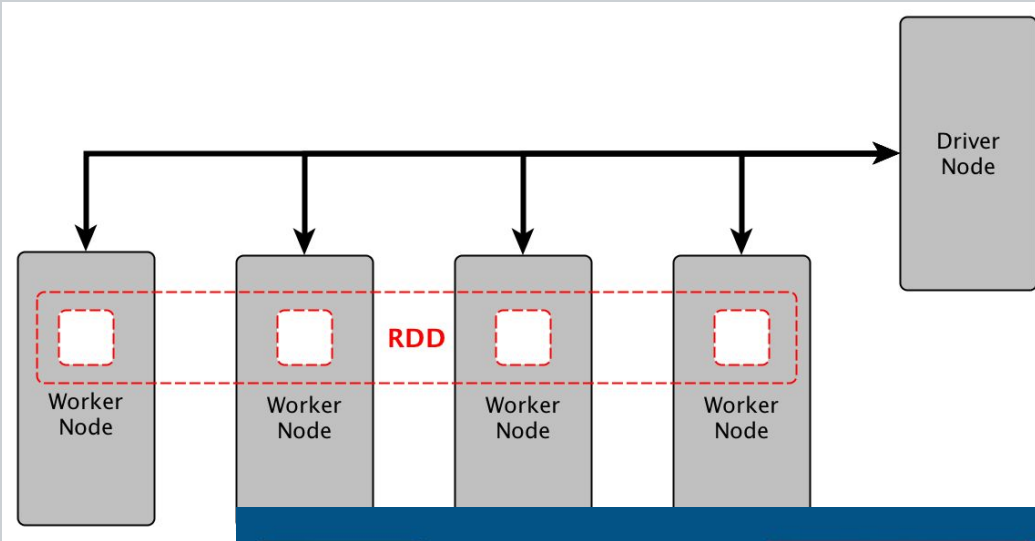  - Streaming
  - Machine learning

# Apache Spark



**Driver Node**
- Spark Context

**Cluster Manager**

**Workers**
- Executors

- Código que criamos
- SparkShell

- Permite escalabilidade

- Executors

**Workers**

https://www.youtube.com/watch?v=jffQhcweGwY

# Apache Spark - Directed Acyclic Graph (DAG)

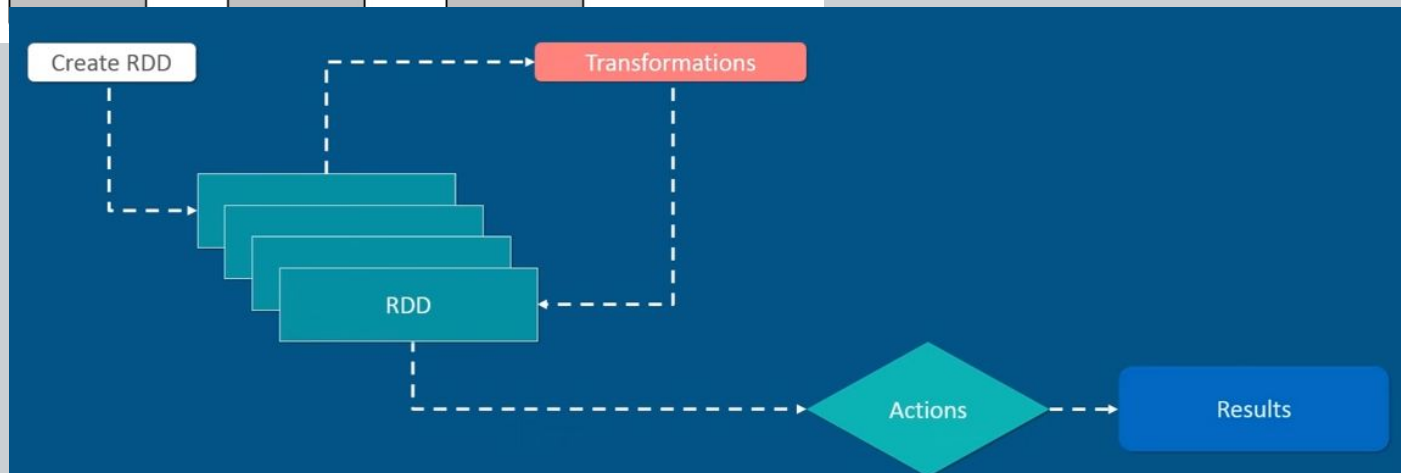UFSC UNIVERSIDADE FEDERAL DE SANTA CATARINA

# Apache Spark



- Resilient Distributed Datasets
- Memory organization
- Immutable objects in JVM

# Apache Spark

- Transformations e Actions no RDD



| Transformations | Actions |
| --- | --- |
| map (func) | reduce(func) |
| flatMap(func) | collect() |
| filter(func) | count() |
| groupByKey() | first() |
| reduceByKey(func) | take(n) |
| mapValues(func) | saveAsTextFile(path) |
| sample(...) | countByKey() |
| union(other) | foreach(func) |
| distinct() | ... |
| sortByKey() | |
| ... | |

# Exercício

- Instalar Spark localmente
- Usar transformations e actions

## DataFrames

| dept | age | name |
|------|-----|------|
| Bio | 48 | H Smith |
| CS | 54 | A Turing |
| Bio | 43 | B Jones |
| Chem | 61 | M Kennedy |

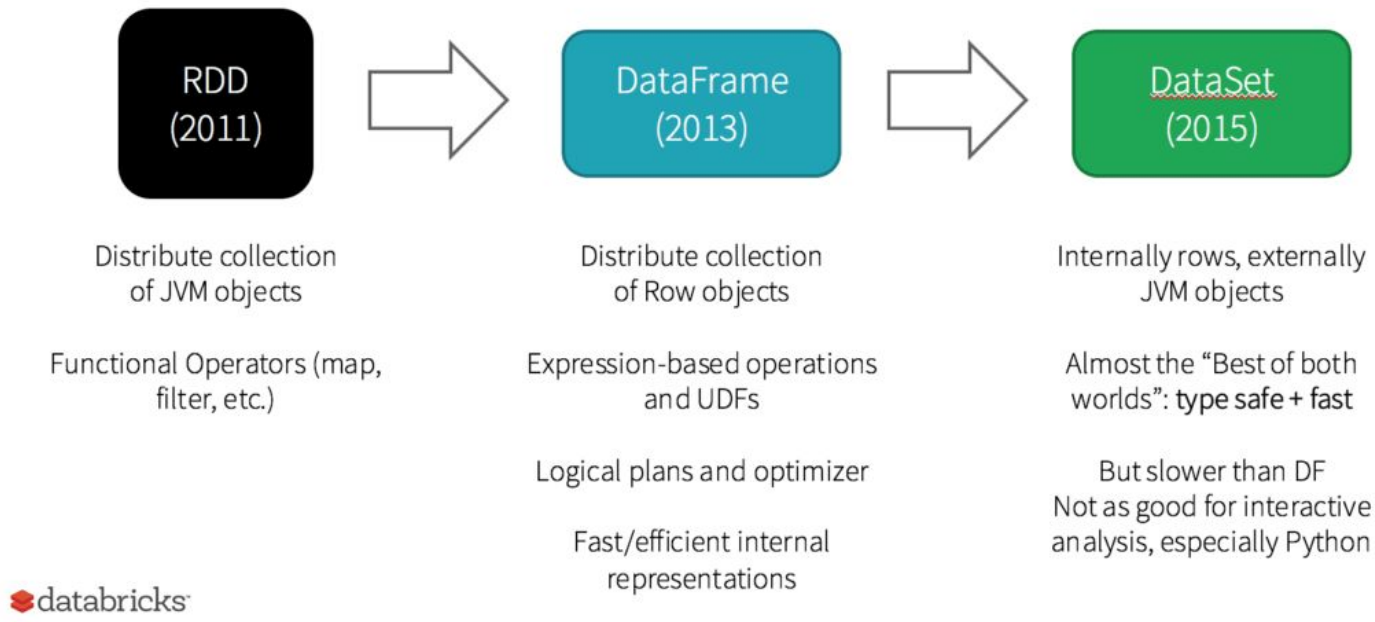Data grouped into named columns

*RDD API*

```
pdata.map(lambda x: (x.dept, [x.age, 1])) \
  .reduceByKey(lambda x, y: [x[0] + y[0], x[1] + y[1]]) \
  .map(lambda x: [x[0], x[1][0] / x[1][1]]) \
  .collect()
```

*DataFrame API*

```
data.groupBy("dept").avg("age")
```

UFSC UNIVERSIDADE FEDERAL DE SANTA CATARINA

# Apache Spark

## History of Spark APIs



| RDD (2011) | | DataFrame (2013) | | DataSet (2015) |
|---|---|---|---|---|
| Distribute collection of JVM objects | | Distribute collection of Row objects | | Internally rows, externally JVM objects |
| Functional Operators (map, filter, etc.) | | Expression-based operations and UDFs | | Almost the "Best of both worlds": type safe + fast |
| | | Logical plans and optimizer | | But slower than DF Not as good for interactive analysis, especially Python |
| | | Fast/efficient internal representations | | |

databricks

# Apache Spark

| | RDD | DATAFRAME | DATASET |
|---|---|---|---|
| What | Distributed collection of elements | Organized into Named Columns | Extension of Dataframe |
| When | 1.0 | 1.3 | 1.6 |
| Compile-time type safety | No | No | Yes |
| APIs | No | Yes | Yes |
| Spark SQL | No | Yes | Yes |
| Catalyst Optimizer | No | Yes | Yes |
| Tungsten component | No | Yes | Yes |
| Advanced Encoders | No | No | Yes |

https://www.youtube.com/watch?v=26zl50iNBp8

# Exercício

- Executar SQL sobre um DataSet

# Spark Streaming



https://spark.apache.org/docs/latest
/streaming-programming-guide.html
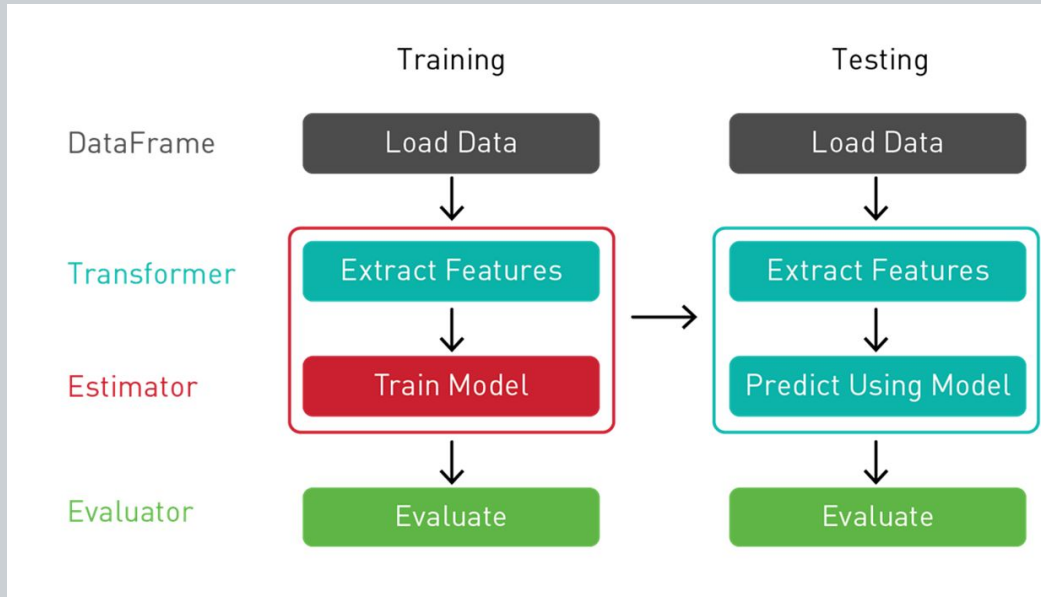#discretized-streams-dstreams

# Spark Streaming



https://spark.apache.org/docs/latest
/streaming-programming-guide.html
#discretized-streams-dstreams

UFSC UNIVERSIDADE FEDERAL DE SANTA CATARINA

# Spark Streaming



- Transform (RDD -> RDD)
  - Map
  - FlatMap
  - Filter
  - Count
  - CountByValue
  - GroupByKey
  - Reduce
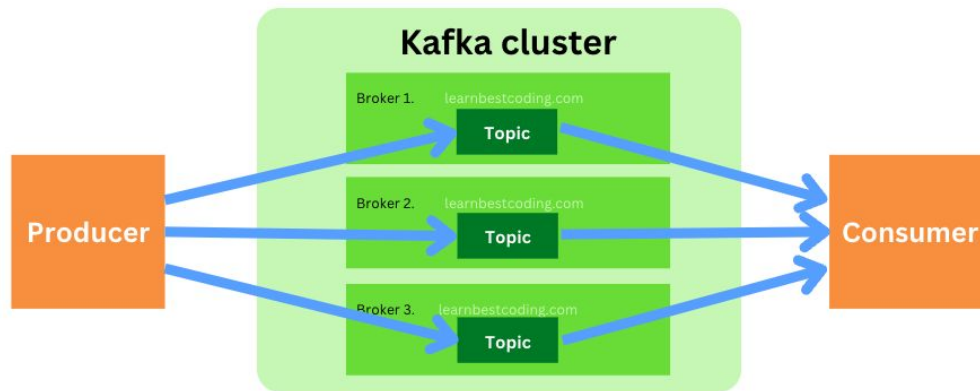  - ReduceByKey
  - Join
  - Cogroup
  - Transform
  - UpdateStateByKey

https://spark.apache.org/docs/latest
/streaming-programming-guide.html
#discretized-streams-dstreams

UFSC UNIVERSIDADE FEDERAL DE SANTA CATARINA

# SparkML

- Permite usar as estruturas do Spark para executar algoritmos de Machine Learning
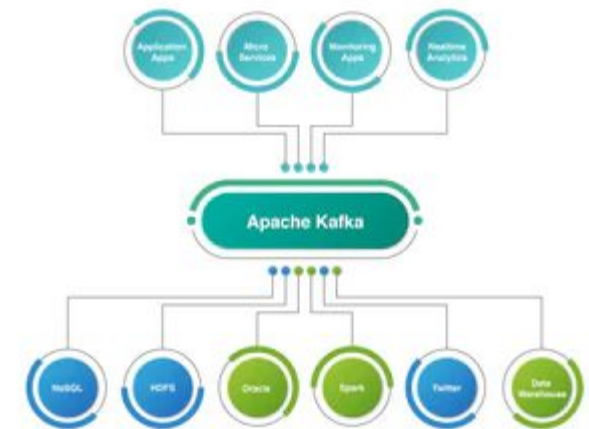- Tem ferramentas para treinamento e uso de modelos
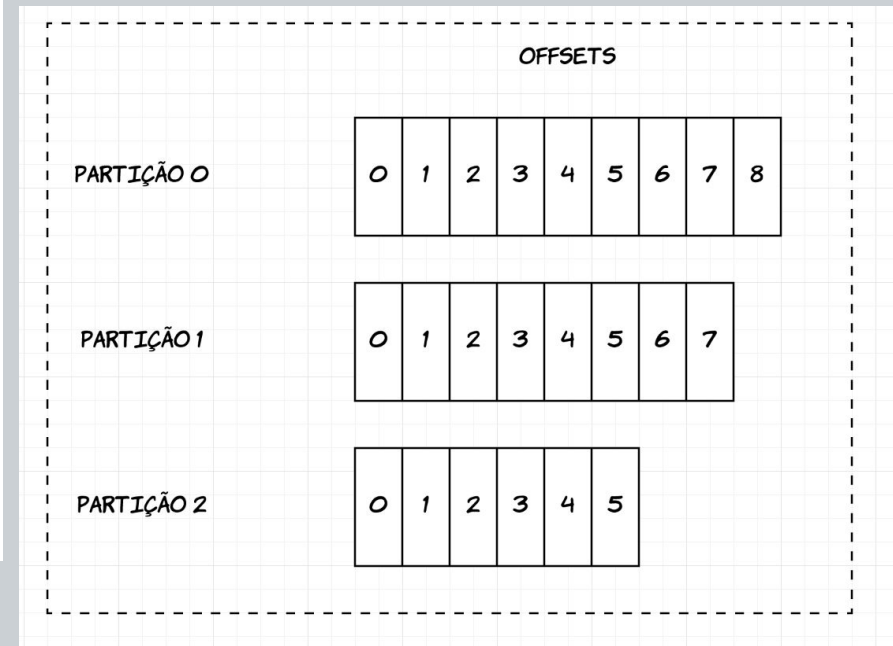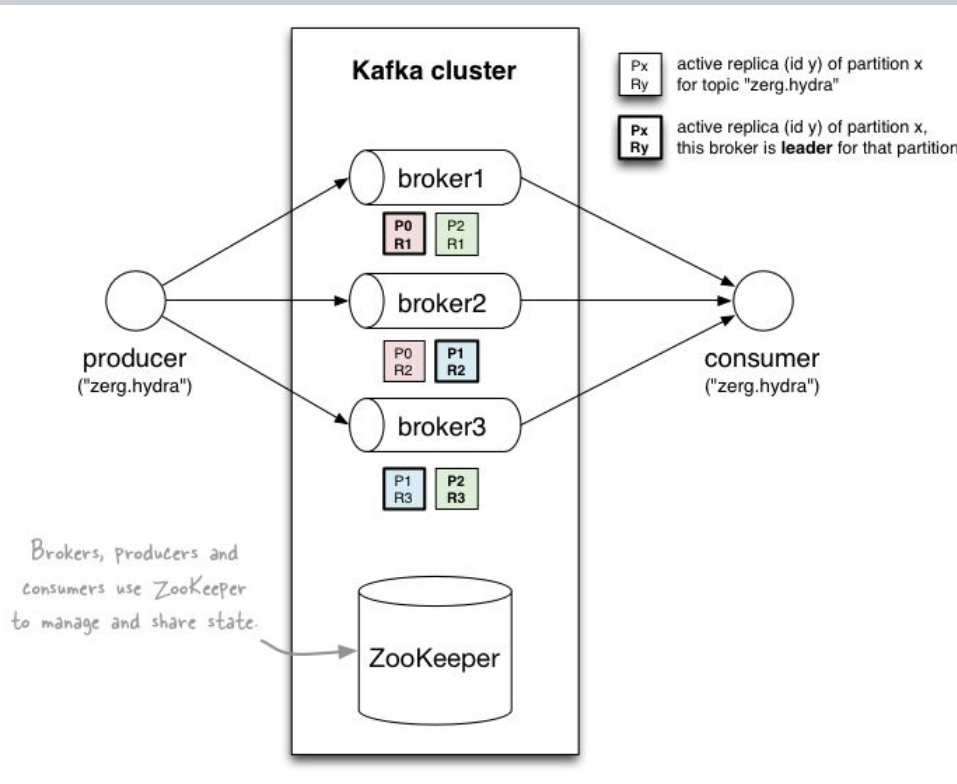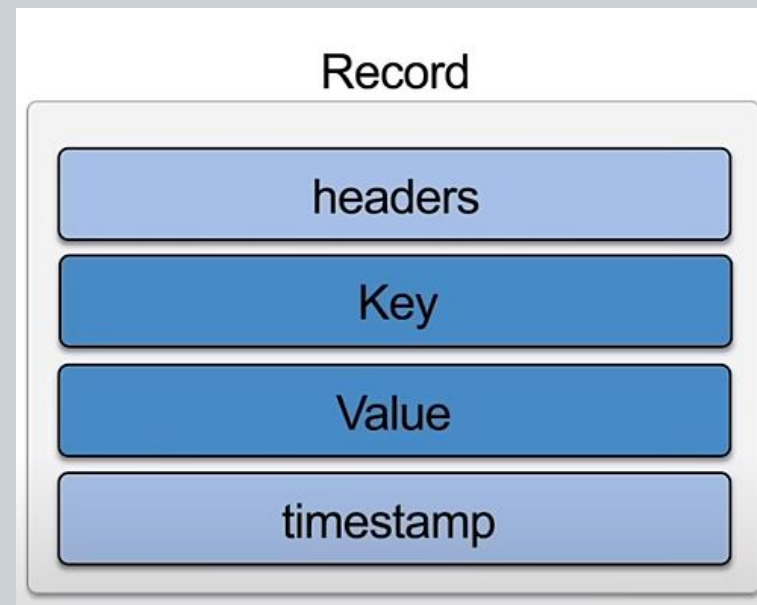
# Apache Kafka

# Apache Kafka

# Exercício

- Instalar Apache Kafka
- Criar um tópico
- Enviar uma mensagem neste tópico



Record
headers
Key
Value
timestamp

UFSC UNIVERSIDADE FEDERAL DE SANTA CATARINA

# Apache NiFi

- NiFi was built to automate the flow of data between systems





NiFi Data Provenance
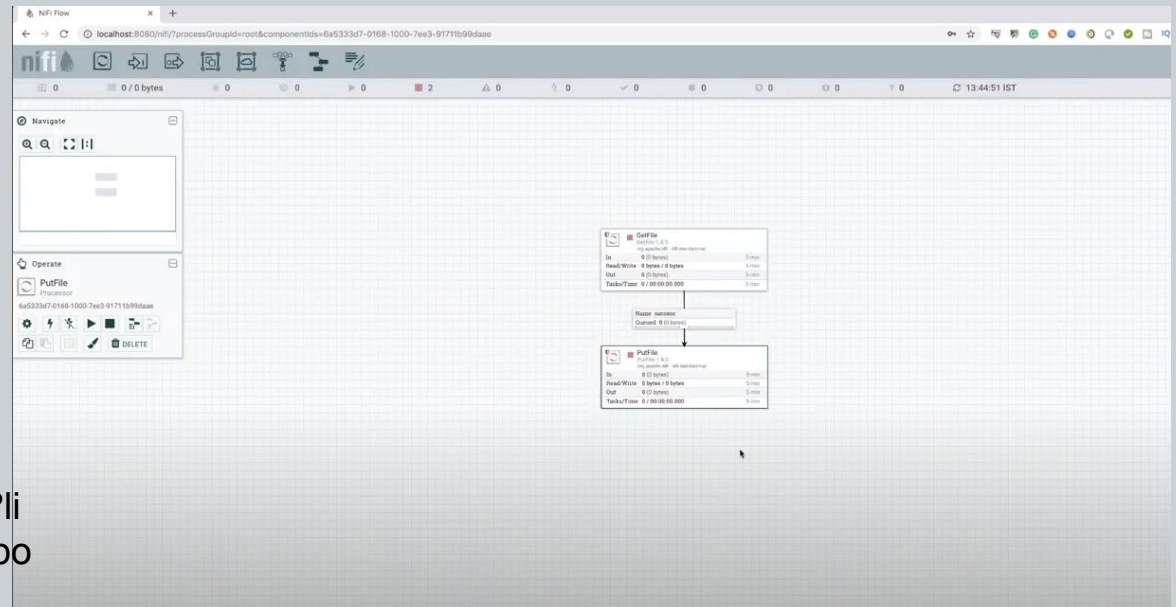
RECEIVE → ATTRIBUTES MODIFIED → FORK → DROP / ROUTE → DROP

# Apache NiFi

- Configuration over coding
- Flow based programming
  - Processors (GetFile)
  - Data Source e Data Sink (Kafka)



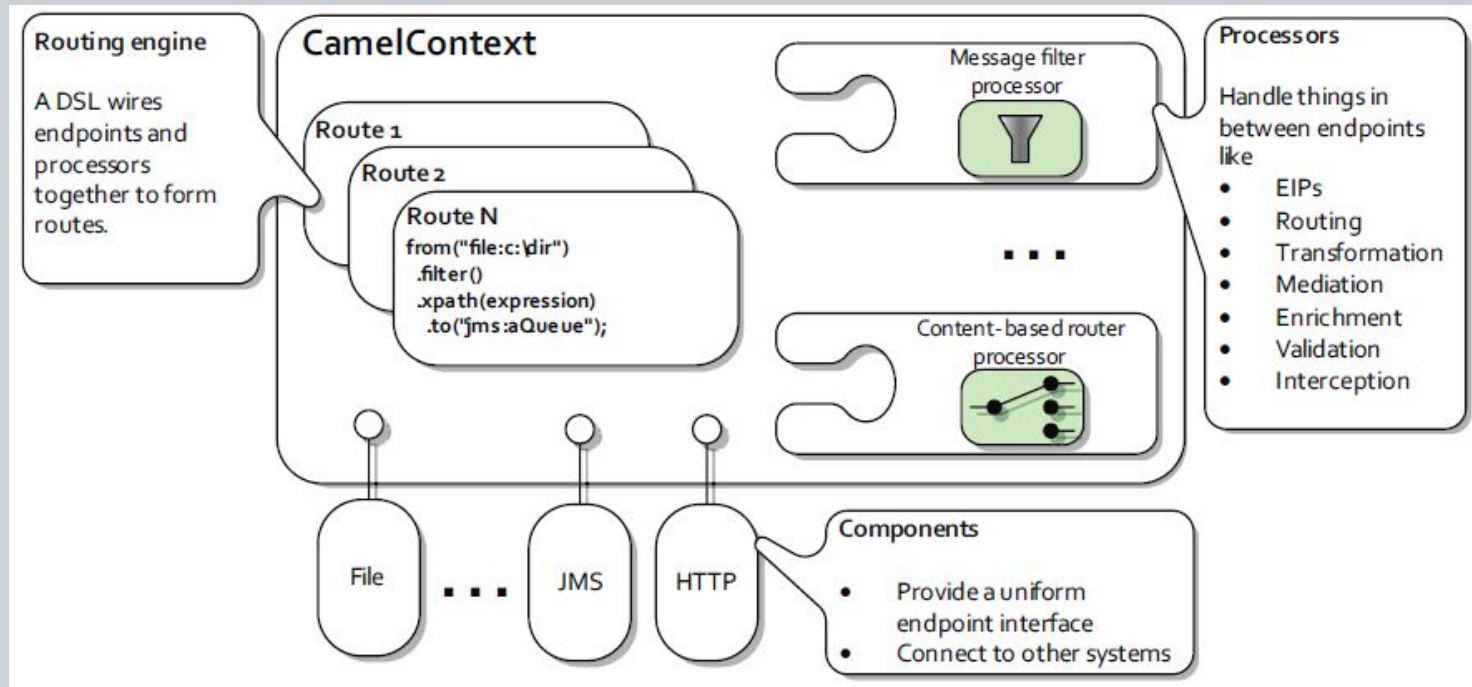https://www.youtube.com/playlist?list=PL55symSEWBbMBSnNW_Aboh2TpYkNIFMgb

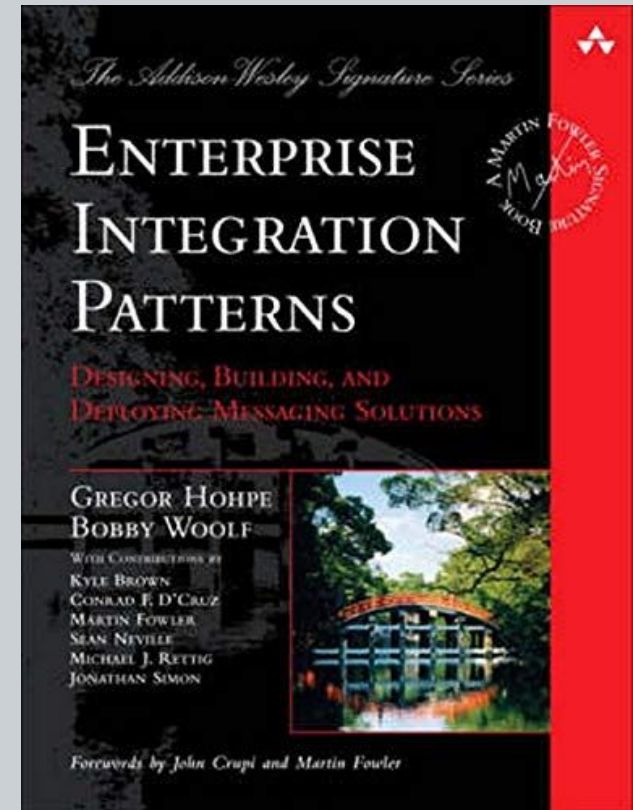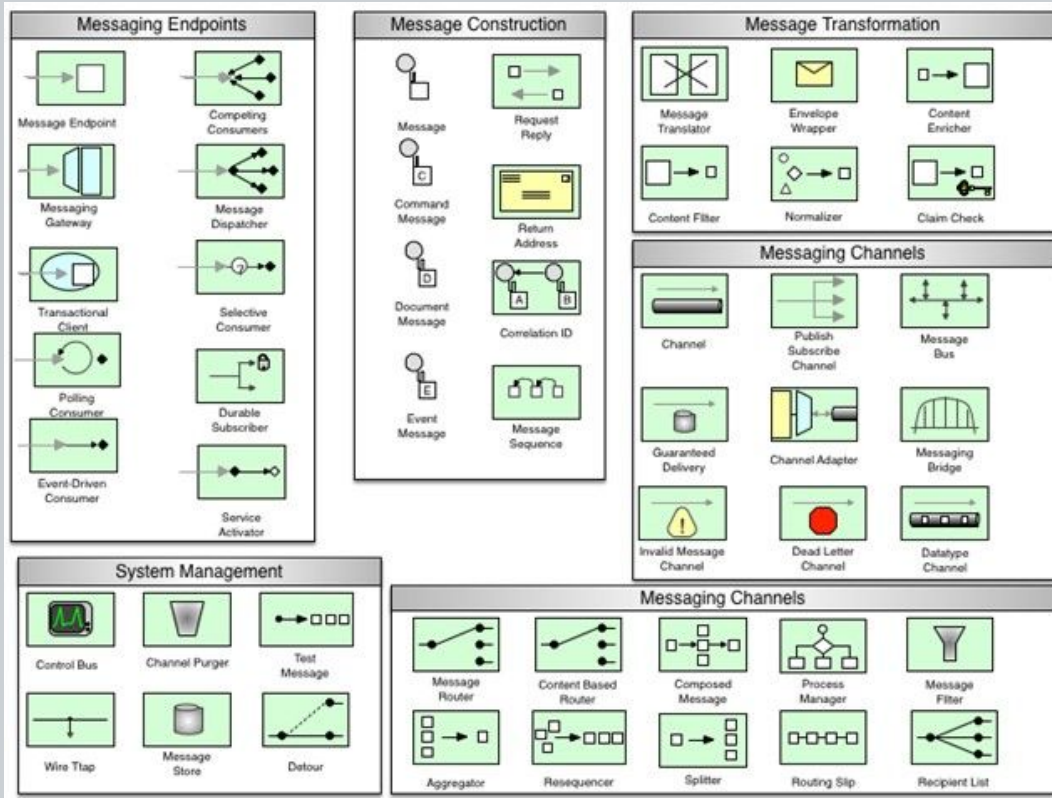UFSC UNIVERSIDADE FEDERAL DE SANTA CATARINA

# Apache Camel

- Ferramenta para integração
- Muito interessante para sistemas legados
- Domain Specific Language (DSL):
  - from("file:data/inbox").to("jms:queue:order");
- Oferece um motor de roteamento altamente escalável

# Apache Camel



**Routing engine**

A DSL wires endpoints and processors together to form routes.

**CamelContext**

Route 1

Route 2

Route N
```
from("file:c:\dir")
 .filter()
 .xpath(expression)
 .to("jms:aQueue");
```

Message filter processor

Content-based router processor

**Processors**

Handle things in between endpoints like
- EIPs
- Routing
- Transformation
- Mediation
- Enrichment
- Validation
- Interception

File  ...  JMS  HTTP

**Components**
- Provide a uniform endpoint interface
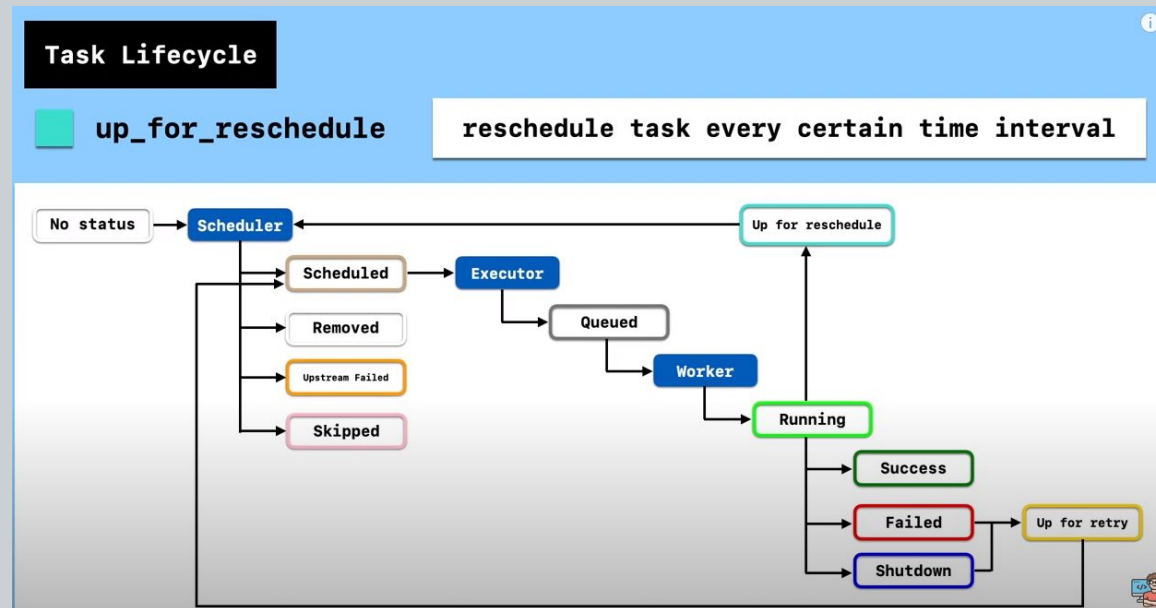- Connect to other systems
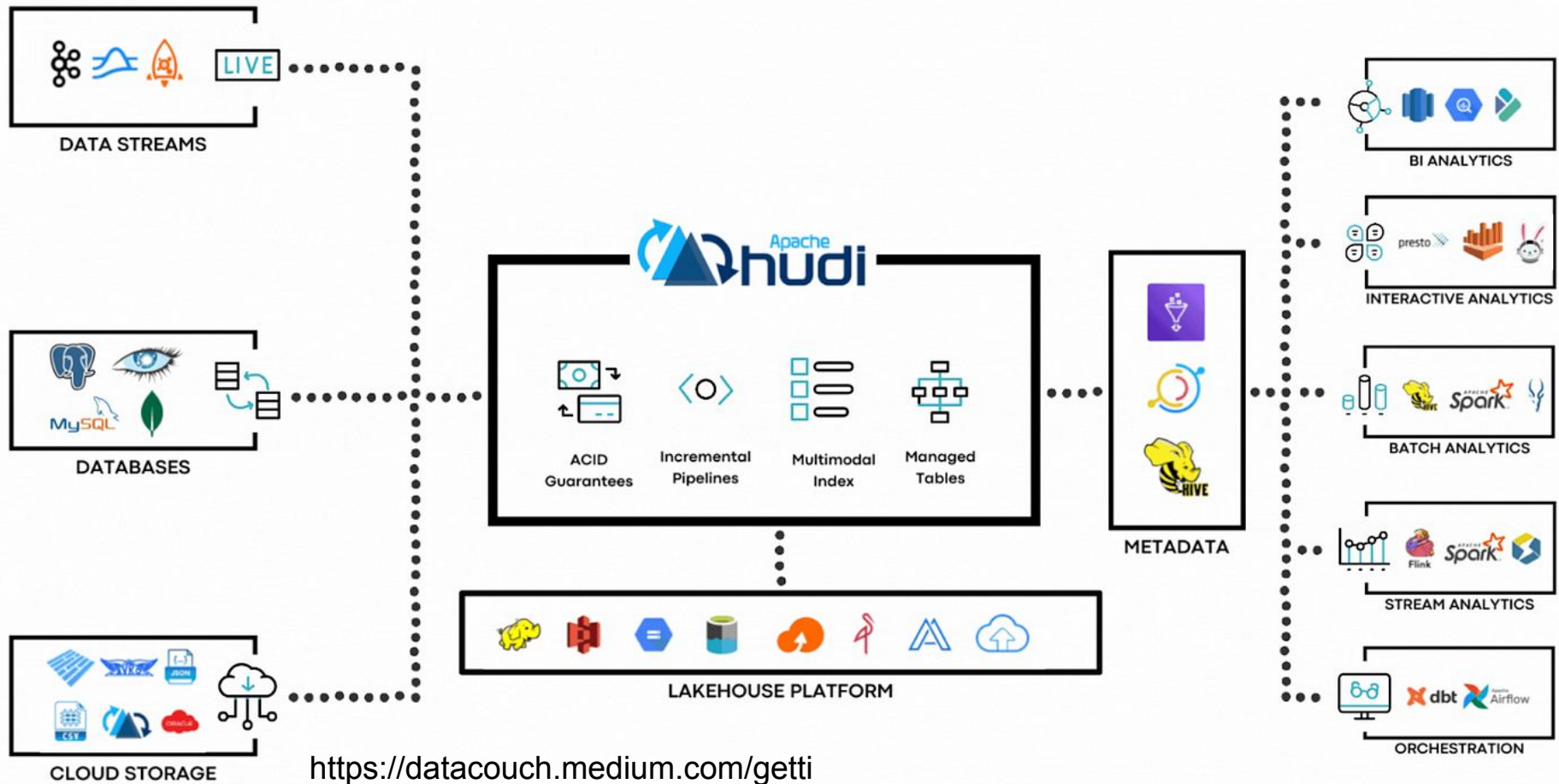
# Apache Camel

- Baseado no EIPs:

# Apache AirFlow

- Criado para o gerenciamento de workflows complexos
- Criado no AirBnb
- Também usa do conceito de DAG
  - Tarefas escritas em Python
  - Operadores



https://www.youtube.com/watch?v=K9AnJ9_ZAXE&list=PLwFJcsJ61oujAqYpMp1kdUBcPG0sE0QMT

# Apache Hudi



https://datacouch.medium.com/getting-started-with-apache-hudi-711b89c107aa

# Delta Lake

- Camada open-source em cima de data lakes
- Trasações ACID
- Time travel

https://www.youtube.com/watch?v=
LJtShrQqYZY

# Delta Lake

UFSC UNIVERSIDADE FEDERAL DE SANTA CATARINA

# Delta Lake

- Camada open-source em cima de data lakes
- Trasações ACID
- Time travel

https://www.youtube.com/watch?v=
LJtShrQqYZY

**Luiz Henrique Zambom Santana**

lhzsantana@gmail.com

https://www.linkedin.com/in/luizsantana/

UNIVERSIDADE FEDERAL
DE SANTA CATARINA