

¿Es Big Data el uso de viejas técnicas matemáticas y estadísticas bajo una mayor capacidad de computación, o es realmente un tema nuevo en el mundo científico?

Big Data

Estrategia para el manejo de información caracterizada por:

- Un alto volumen de datos.
- Alta velocidad en la recepción de los datos
- Variedad: estructurada, semiestructurada, no estructurada

Analytics

Proceso científico de transformación de datos en conocimiento para mejorar el proceso de toma de decisiones [\[Informs\]](#).

I. Problema organizacional

II. Transformación en un problema de analytics

III. Datos

IV. Selección de la metodología

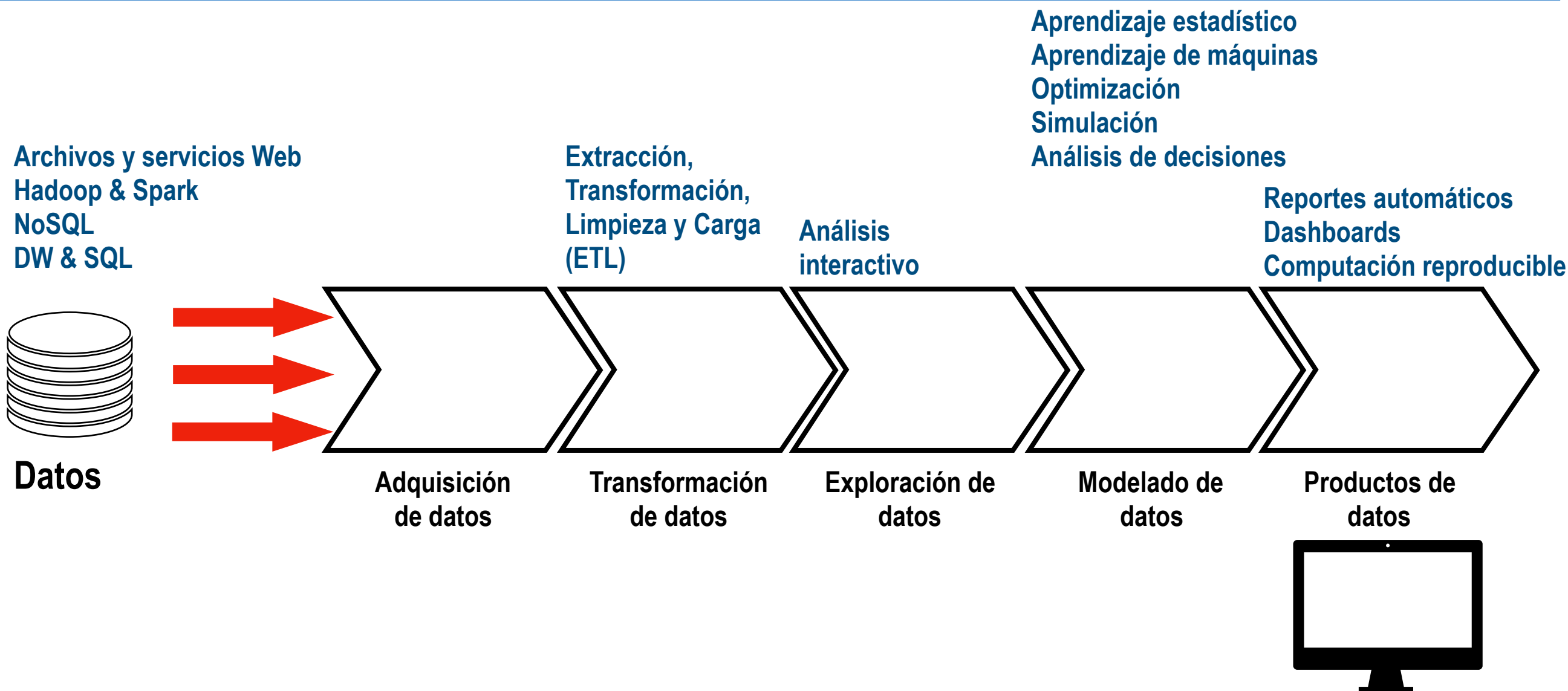
V. Desarrollo del modelo

VI. Puesta en marcha (deploy)

VII. Gestión del ciclo de vida del modelo

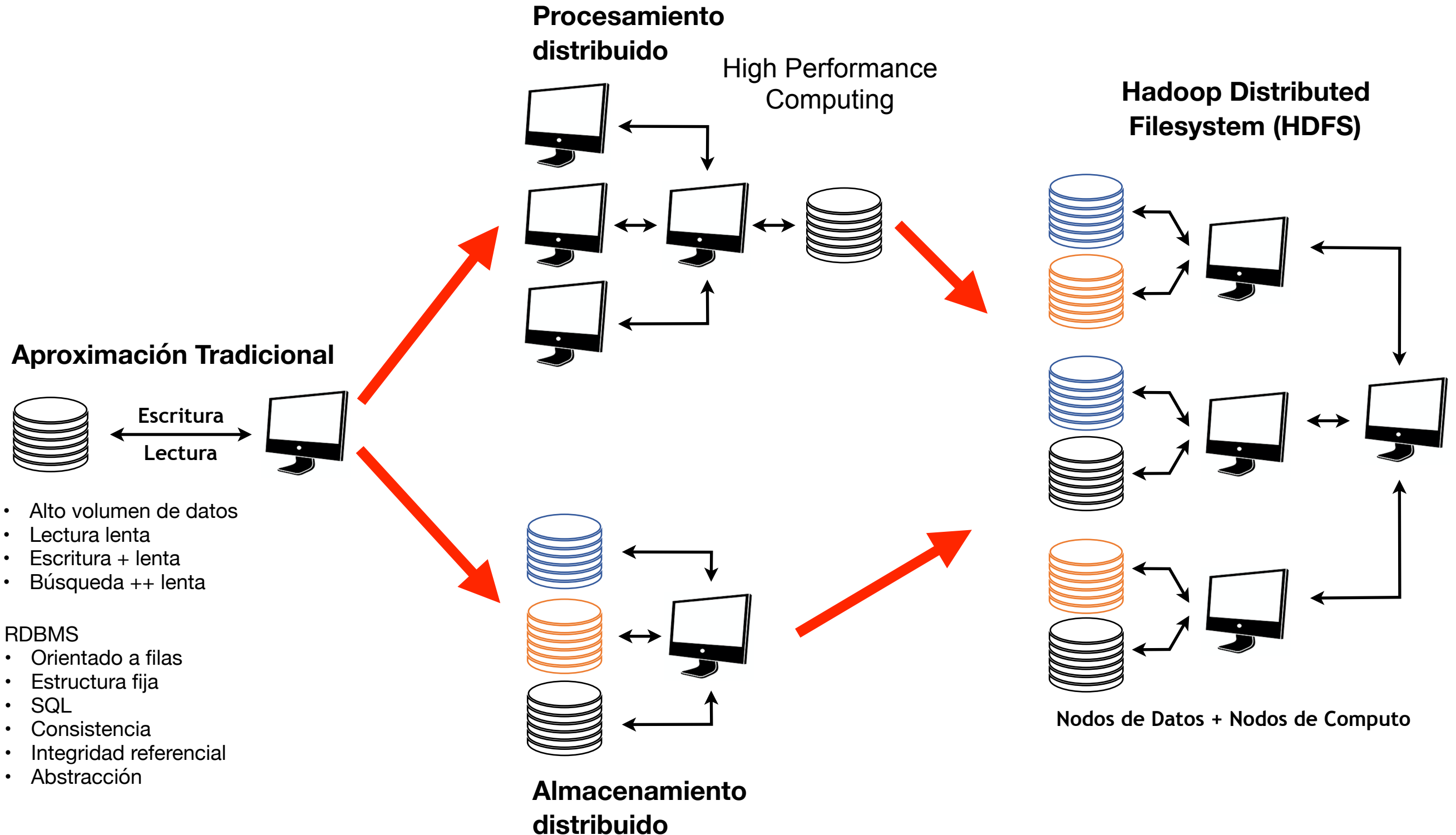
Analytics

Algoritmos, técnicas y metodologías



Infraestructura computacional

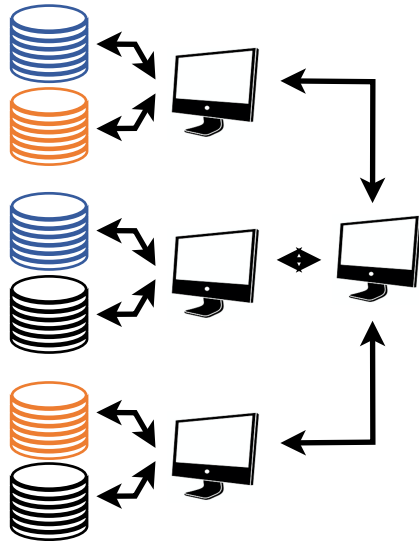
Hadoop / MapReduce



Infraestructura computacional

Computación local

Servidores + red + clientes



Cloud computing / utility computing

Servidores y almacenamiento en la nube + internet + clientes locales

Software as a Service (SaaS)

Software almacenado en máquinas suministradas por un tercero.

Aplicaciones accesadas vía un cliente o la Web.

Orientado a aplicaciones de usuario final.

Platform as a Service (PaaS)

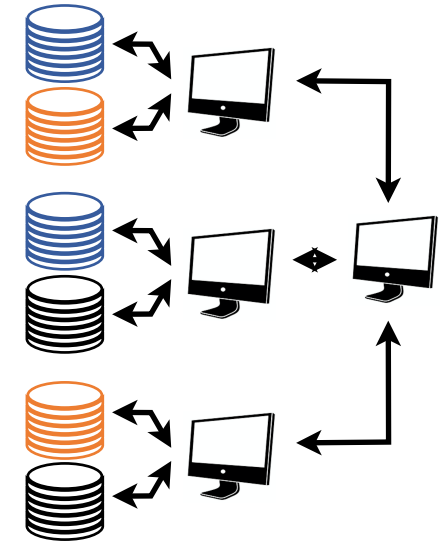
Orientado a desarrolladores.

Ambiente de desarrollo gestionado por un tercero.

Infrastructure as a Service (IaaS)

Bloques básicos para construcción de ambientes manejados por un tercero

Capacidad de procesamiento, almacenamiento, conectividad, seguridad, etc.



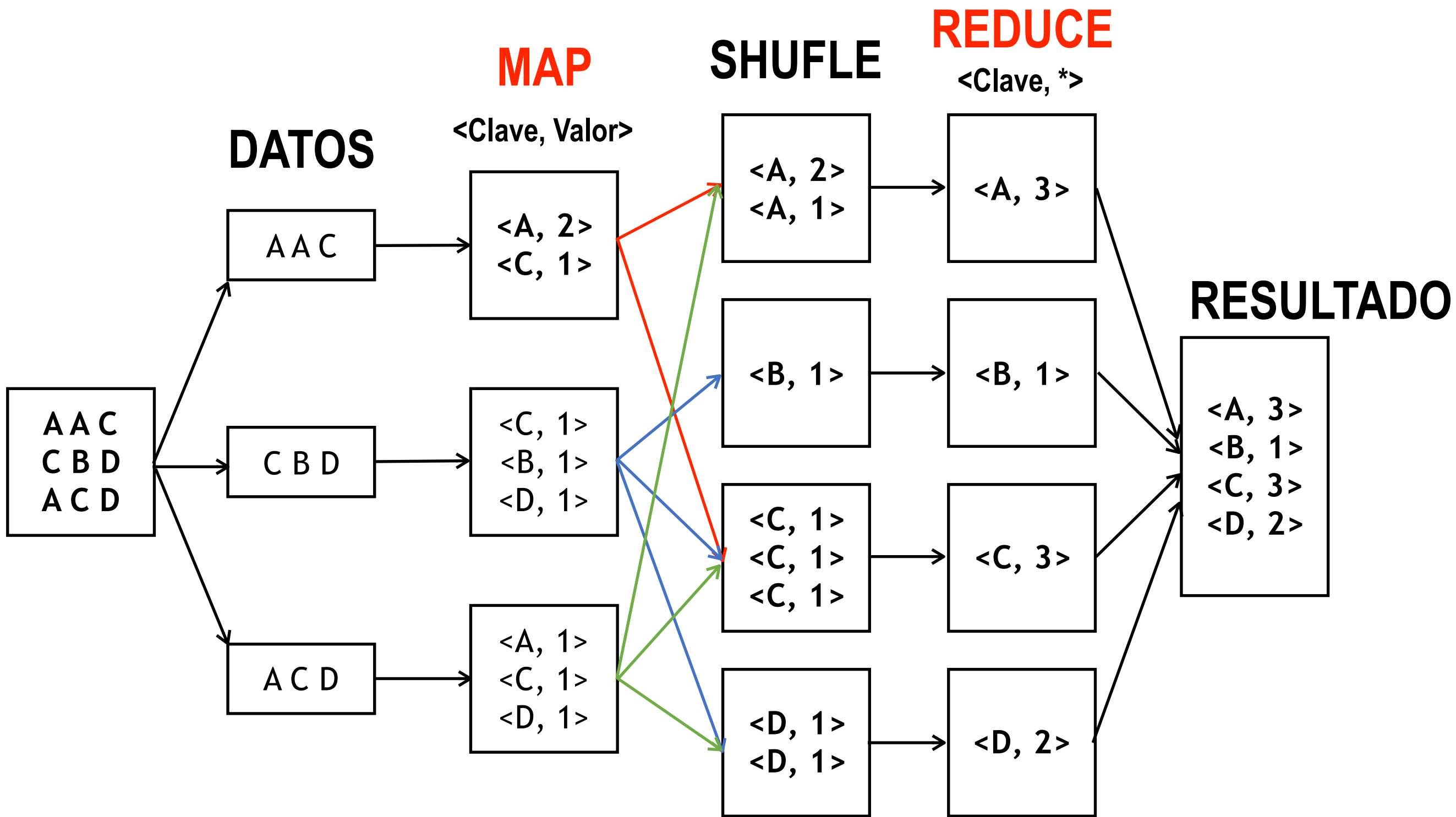
Nube

Internet



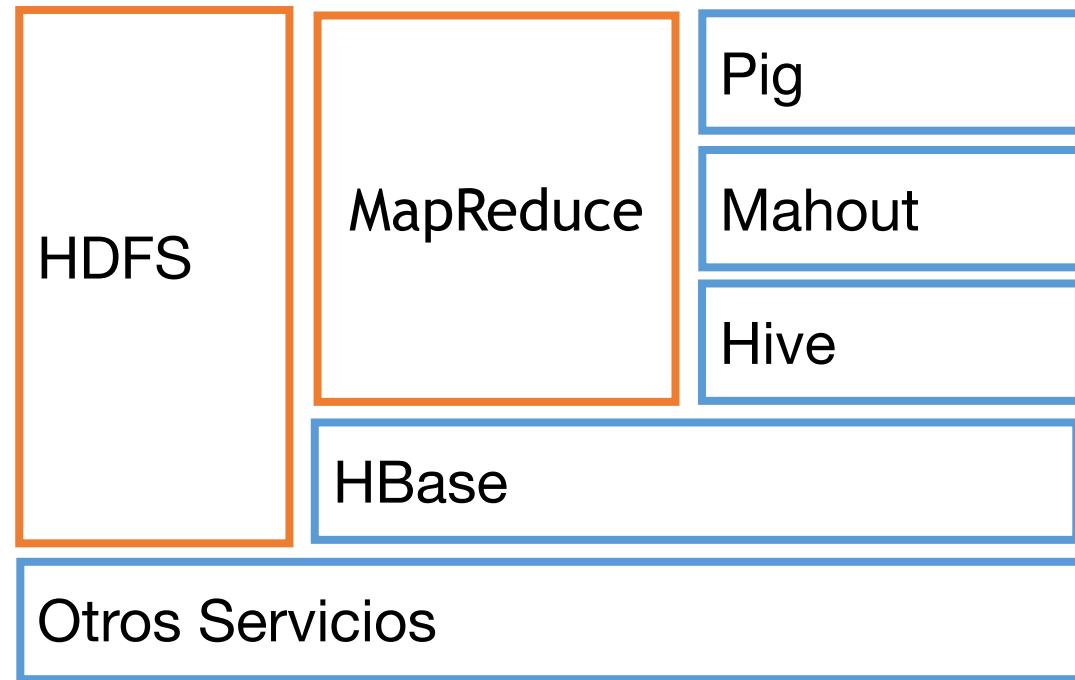
Máquina Local
(Cliente)

MapReduce



Ecosistema Apache Hadoop

Hadoop 1.x



HDFS
Operaciones básicas del sistema de archivos

MapReduce
Algoritmo de procesamiento (Java)

HBase
Base de datos orientada a columnas.

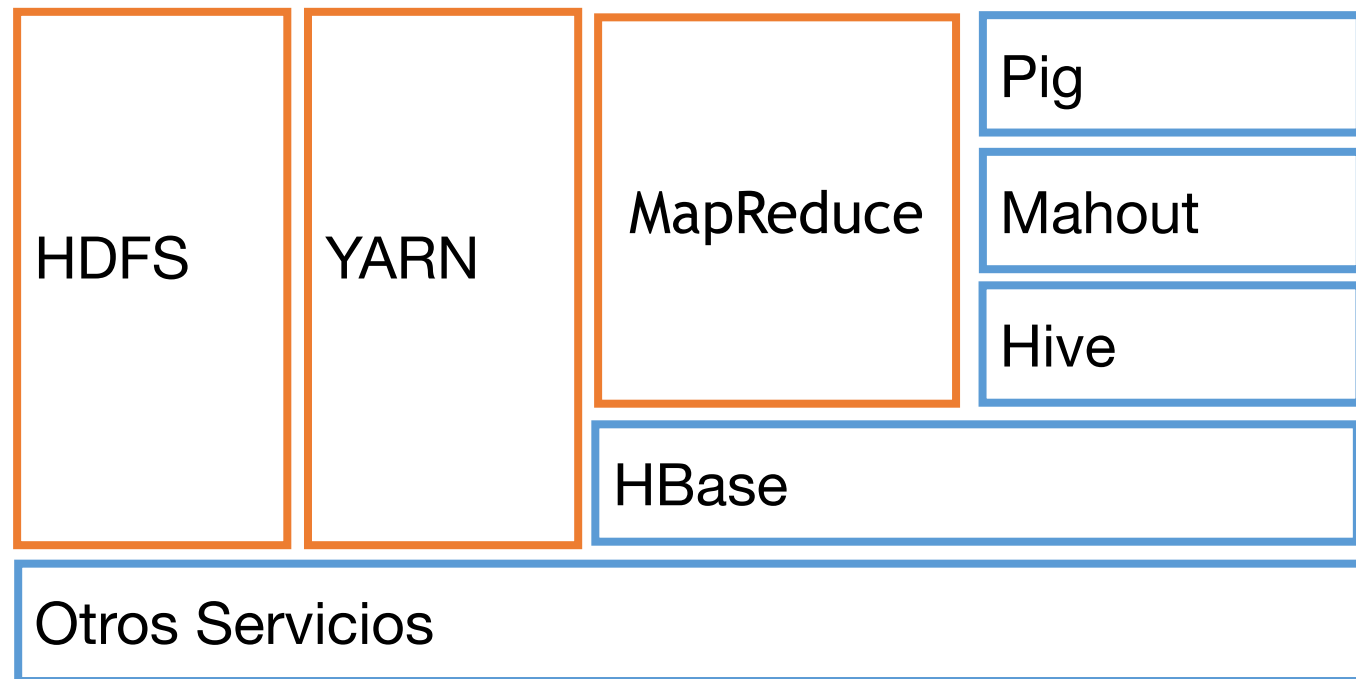
Pig
Lenguaje para el procesamiento de datos.

Mahout
Machine Learning usando MarReduce

Hive
Lenguaje de consultas (~SQL).

YARN
Manejo y organización de recursos del cluster

Hadoop 2.x



Ecosistema Apache Hadoop

Ejemplo de Pig

```
records = LOAD 'sample.txt' AS (year:chararray, temperature:int, quality:int);
filtered_records = FILTER records BY temperature;
grouped_records = GROUP filtered_records BY year;
max_temp = FOREACH grouped_records GENERATE group, MAX(filtered_records.temperature);
DUMP max_temp;
```

Ejemplo de Hive

```
CREATE TABLE records (year STRING, temperature INT, quality INT)
  ROW FORMAT DELIMITED FIELDS TERMINATED BY '\t';
LOAD DATA LOCAL INPATH 'sample.txt' OVERWRITE INTO TABLE records;
SELECT year, MAX(temperature) FROM records GROUP BY year;
```


Ecosistema Apache Hadoop

HDFS

Operaciones básicas del sistema de archivos

MapReduce

Algoritmo de procesamiento (Java)

HBase

Base de datos orientada a columnas.

Pig

Lenguaje para el procesamiento de datos.

Mahout

Machine Learning usando MapReduce

Hive

Lenguaje de consultas (~SQL).

YARN

Manejo y organización de recursos del cluster

Sqoop

ZooKeeper

Cassandra

Avro

Zeppelin

Oozie

Tez

Flume

Ambari

Kafka

Storm

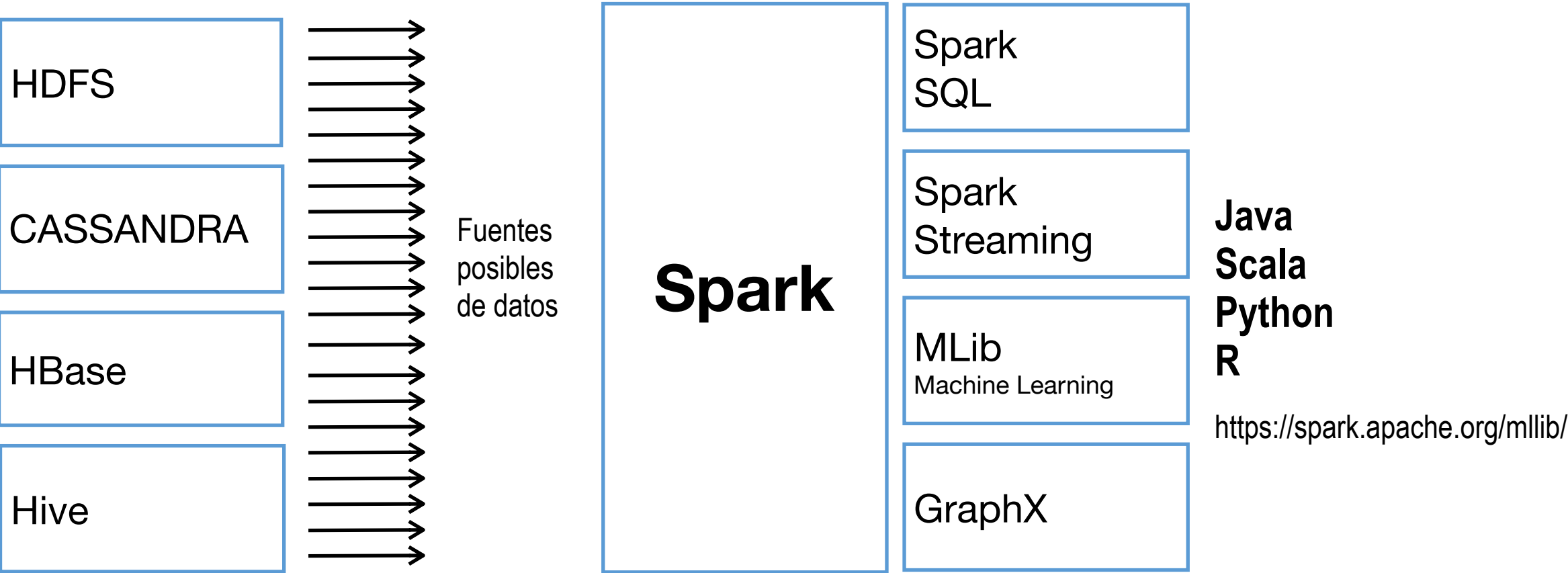
Acummulo

Solr

Knox

HAWQ

Mahout, MLib, R y Python



Mahout vs MLlib

MLib Machine Learning

Algorithms

MLlib contains many algorithms and utilities.

ML algorithms include:

- Classification: logistic regression, naive Bayes,...
- Regression: generalized linear regression, survival regression,...
- Decision trees, random forests, and gradient-boosted trees
- Recommendation: alternating least squares (ALS)
- Clustering: K-means, Gaussian mixtures (GMMs),...
- Topic modeling: latent Dirichlet allocation (LDA)
- Frequent itemsets, association rules, and sequential pattern mining

ML workflow utilities include:

- Feature transformations: standardization, normalization, hashing,...
- ML Pipeline construction
- Model evaluation and hyperparameter tuning
- ML persistence: saving and loading models and Pipelines

Other utilities include:

- Distributed linear algebra: SVD, PCA,...
- Statistics: summary statistics, hypothesis testing,...

Collaborative Filtering *with CLI drivers*

User-Based Collaborative Filtering

Item-Based Collaborative Filtering

Matrix Factorization with ALS

Matrix Factorization with ALS on Implicit Feedback

Weighted Matrix Factorization, SVD++

Classification *with CLI drivers*

Logistic Regression - trained via SGD

Naive Bayes / Complementary Naive Bayes

Hidden Markov Models

Clustering *with CLI drivers*

Canopy Clustering

k-Means Clustering

Fuzzy k-Means

Streaming k-Means

Spectral Clustering

Dimensionality Reduction *note: most scale reduction algorithms are available through the* **Core Library for all engines**

Singular Value Decomposition

Lanczos Algorithm

Stochastic SVD

PCA (via Stochastic SVD)

QR Decomposition

Topic Models

Latent Dirichlet Allocation

Miscellaneous

RowSimilarityJob

Collocations

Sparse TF-IDF Vectors from Text

XML Parsing

Email Archive Parsing

Evolutionary Processes

Ahora la respuesta a la pregunta

Si y No, pero tal vez

Elementos teóricos y conceptuales.

Mucha teoría, técnicas y herramientas desarrolladas para un puñado de datos.

Elementos prácticos

Los resultados numéricos y tiempos de cómputo son afectados por el gran volumen de datos.

Los nuevos algoritmos deben aprovechar la infraestructura computacional existente.

Elementos de la práctica profesional

Los profesionales deben dominar las herramientas computacionales del estado del arte.

¿Es Big Data el uso de viejas técnicas matemáticas y estadísticas bajo una mayor capacidad de computación, o es realmente un tema nuevo en el mundo científico?