

Big Data, fundamentos y tópicos avanzados

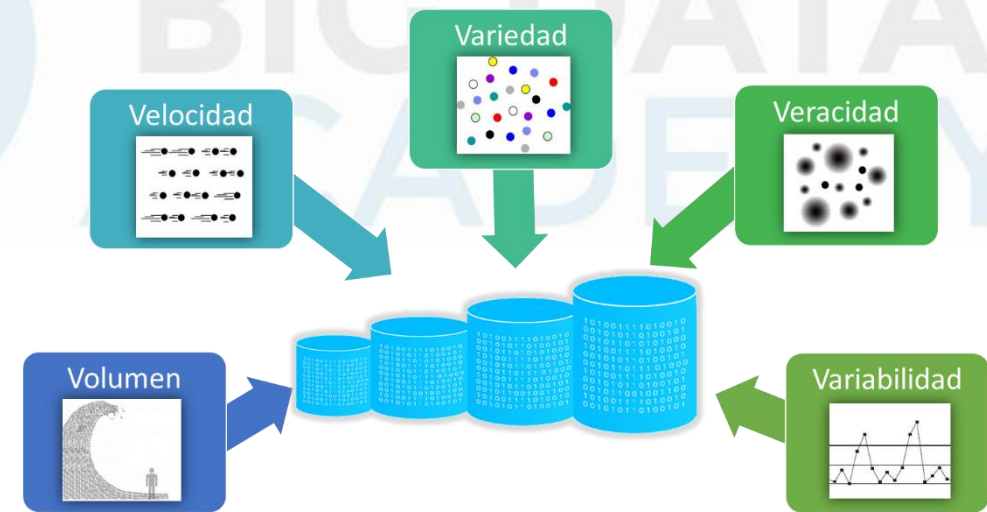
PARTE DE LA FORMACIÓN BIG DATA ENGINEER DE BIG DATA ACADEMY

Concepto

BIG DATA
ACADEMY

Big Data

Es un marco de trabajo (conceptos + tecnologías) que permite procesar grandes volúmenes de datos, de diferentes estructuras o con carencia de estas, que pueden variar en el tiempo, a grandes velocidades que generen valor al negocio.



Objetivo fundamental

Objetivo fundamental del Big Data

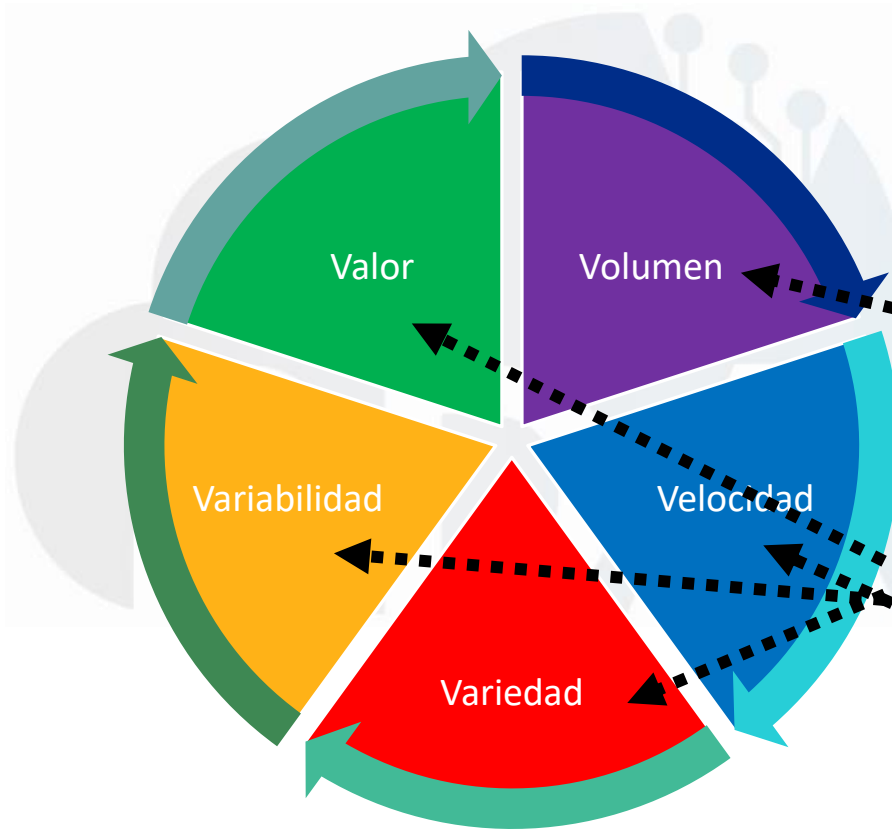
Aumentar el valor del proceso por medio de:

1. **Reducir los tiempos** de procesamiento
2. **Integrar todas las fuentes** de datos disponibles
3. **Reducir los costos** de hardware
4. **Reducir el uso de recursos** computacionales
5. **Crece fácilmente** en potencia computacional
6. **Aumentar la exactitud** en los cálculos
7. **Potenciar otras tecnologías** y marcos de trabajo



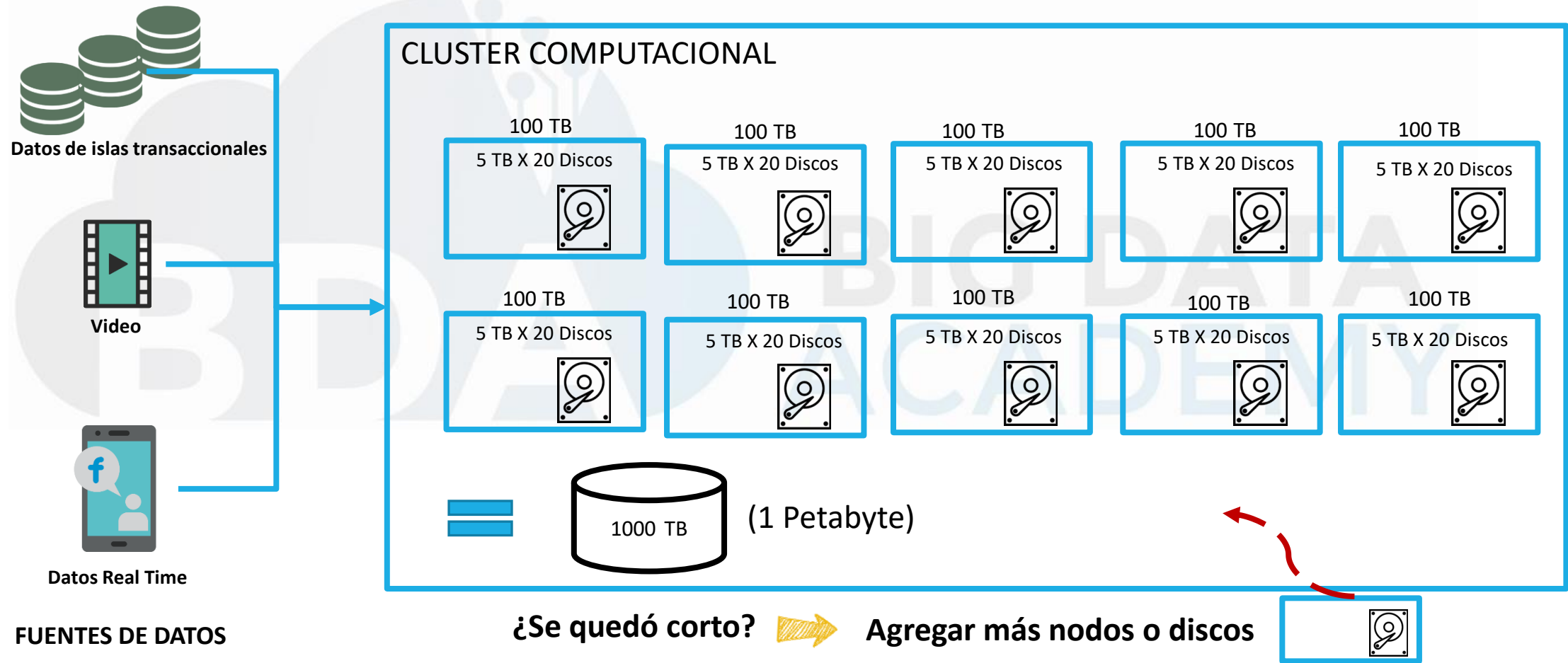
Big Data: Naturaleza

La filosofía de Big Data: Las 5V

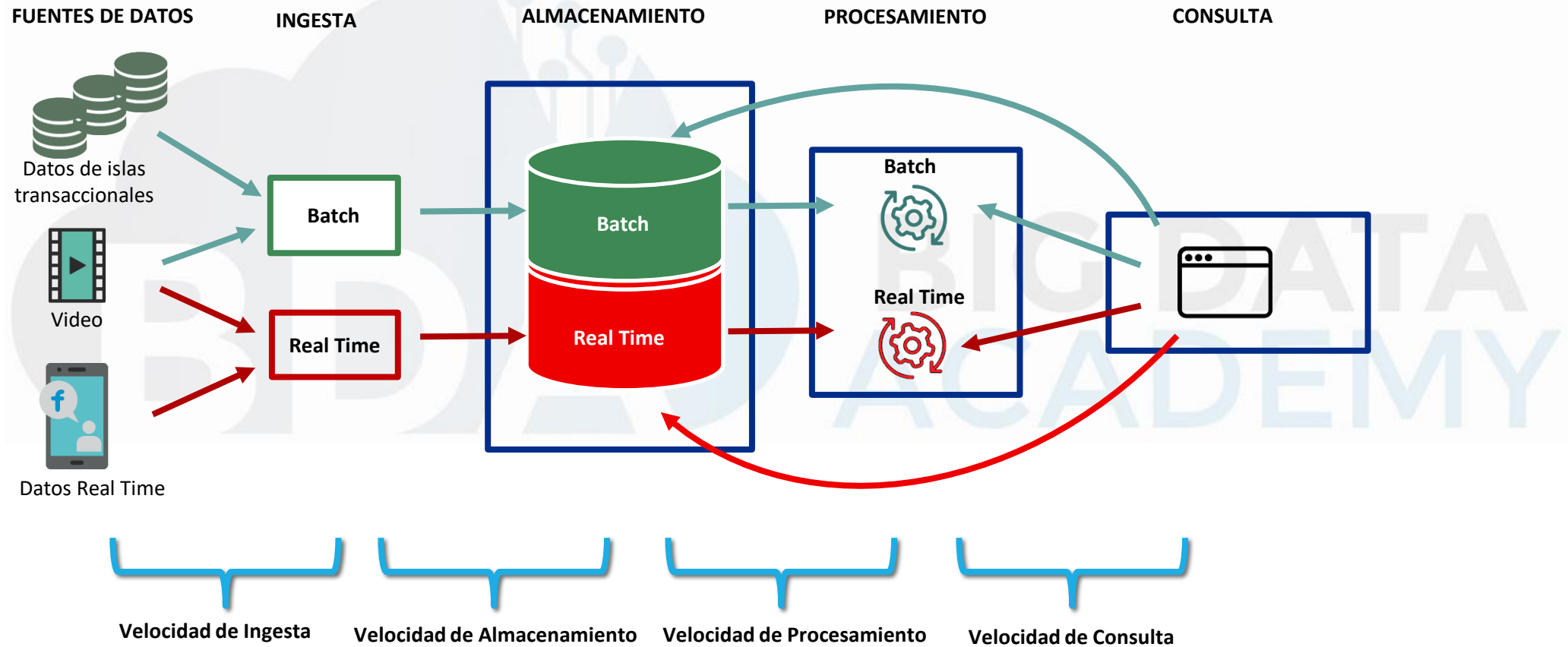


Es un marco de trabajo (conceptos + tecnologías) que permite procesar **grandes volúmenes** de datos, de **diferentes estructuras o con carencia** de estas, que **pueden variar en el tiempo**, a **grandes velocidades** y que generen **valor al negocio**.

La filosofía de Big Data: Volumen



La filosofía de Big Data: Velocidad



La filosofía de Big Data: Variedad

La variedad hace referencia al tipo de estructura del dato

FUENTES DE DATOS



Datos de islas transaccionales



Datos Real Time



Video

DATA ESTRUCTURADA

Misma estructura para todos los registros



Tablas



CSV

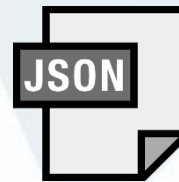


XML



DATA SEMI-ESTRUCTURADA

Cada registro tiene su propia estructura



JSON



XML



DATA NO ESTRUCTURADA

No tiene estructura ni registro



Videos



E-mail



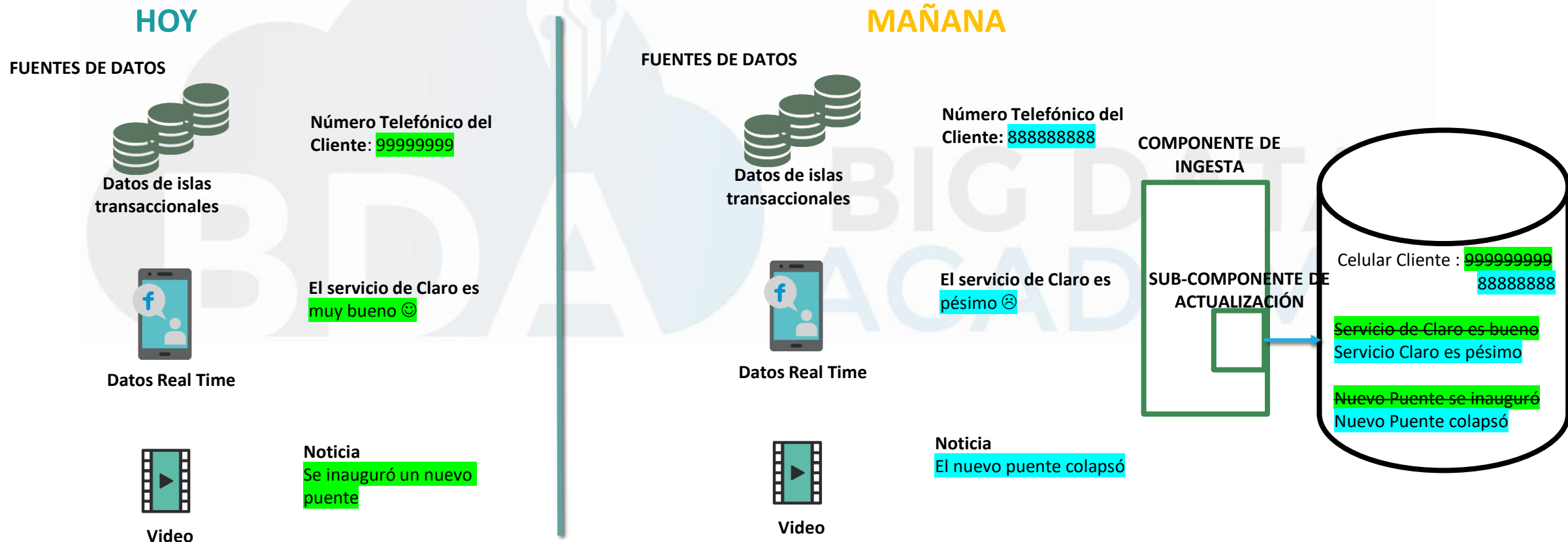
Documentos



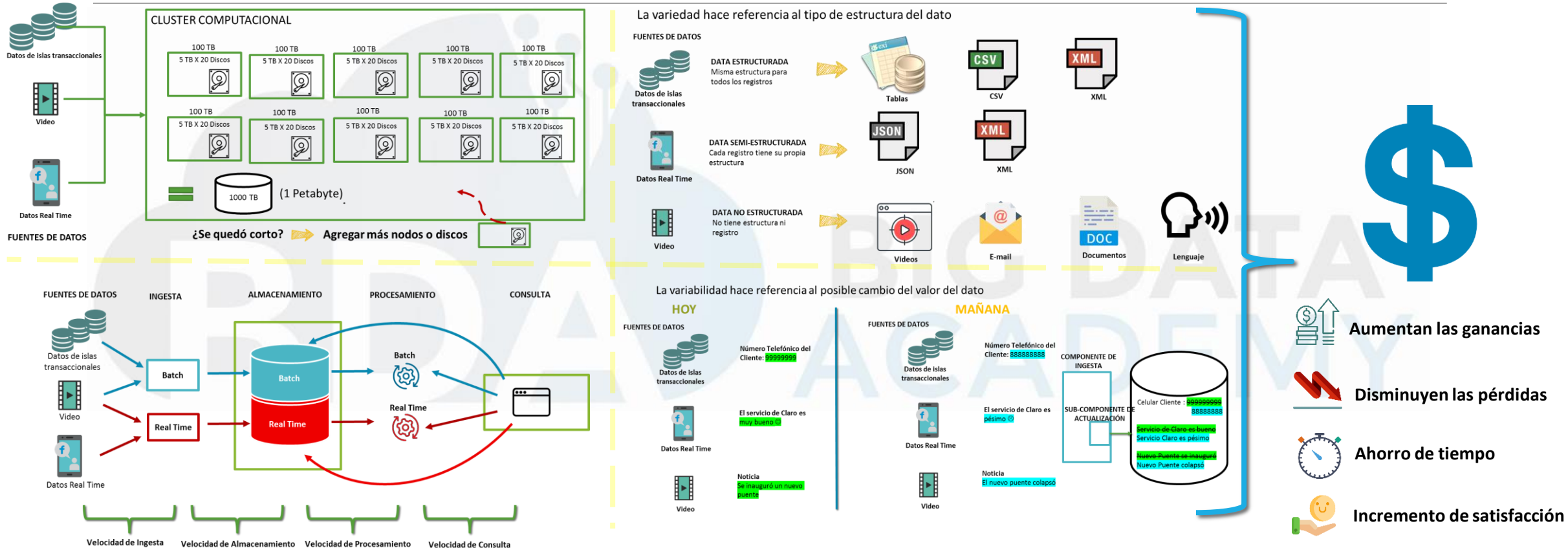
Lenguaje

La filosofía de Big Data: Variabilidad

La variabilidad hace referencia al posible cambio del valor del dato

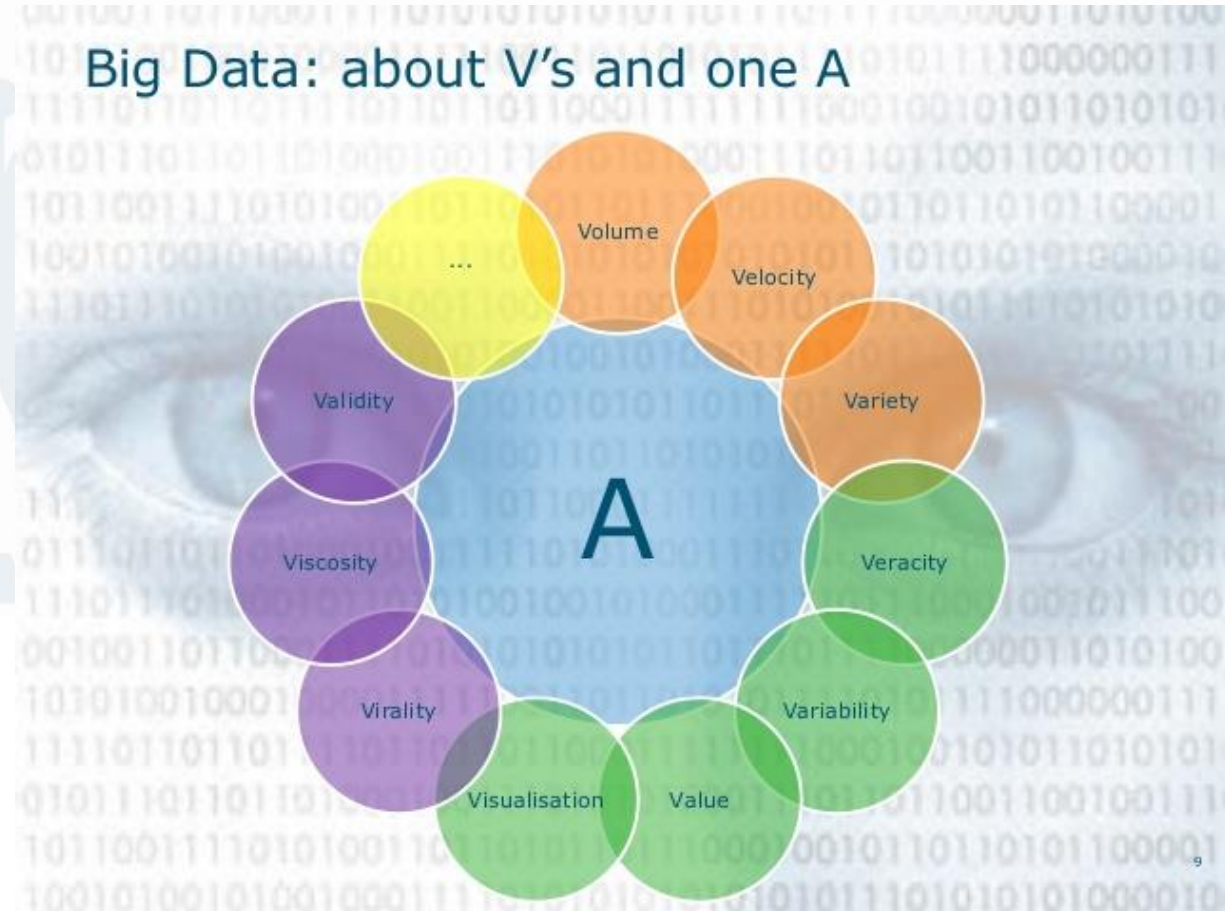
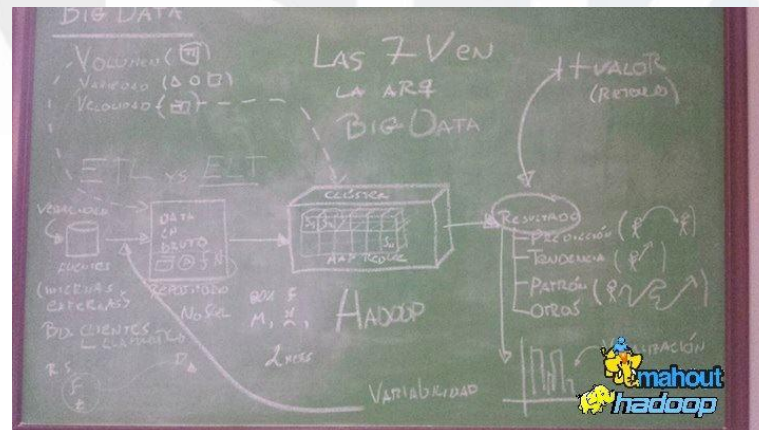
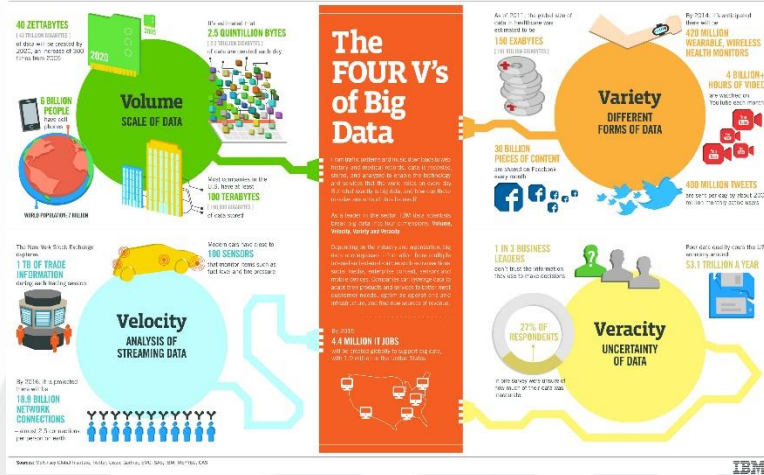


La filosofía de Big Data: Valor



Todos nuestros esfuerzos deben de traducirse en una ganancia para la organización

Otras V



¿Por qué Big Data es un marco de trabajo?

CONCEPTOS

- Las 5V
- Clúster computacional
- Paralelización
- Distribución de carga de trabajo
- Escalabilidad
- Alta disponibilidad
- Seguridad
- Gobierno
- Patrones de diseño

TECNOLOGÍAS

- Hadoop
- Hive
- HBase
- Spark
- Kafka
- Cassandra
- Lenguajes de programación

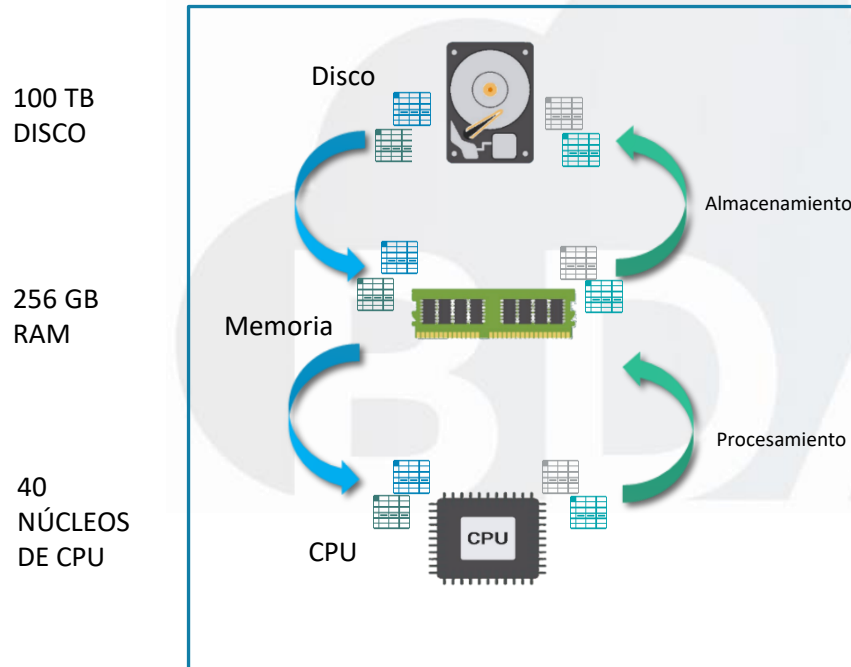
MARCO DE
TRABAJO (*del
inglés
framework*)

BIG DATA

Big Data: Conceptos

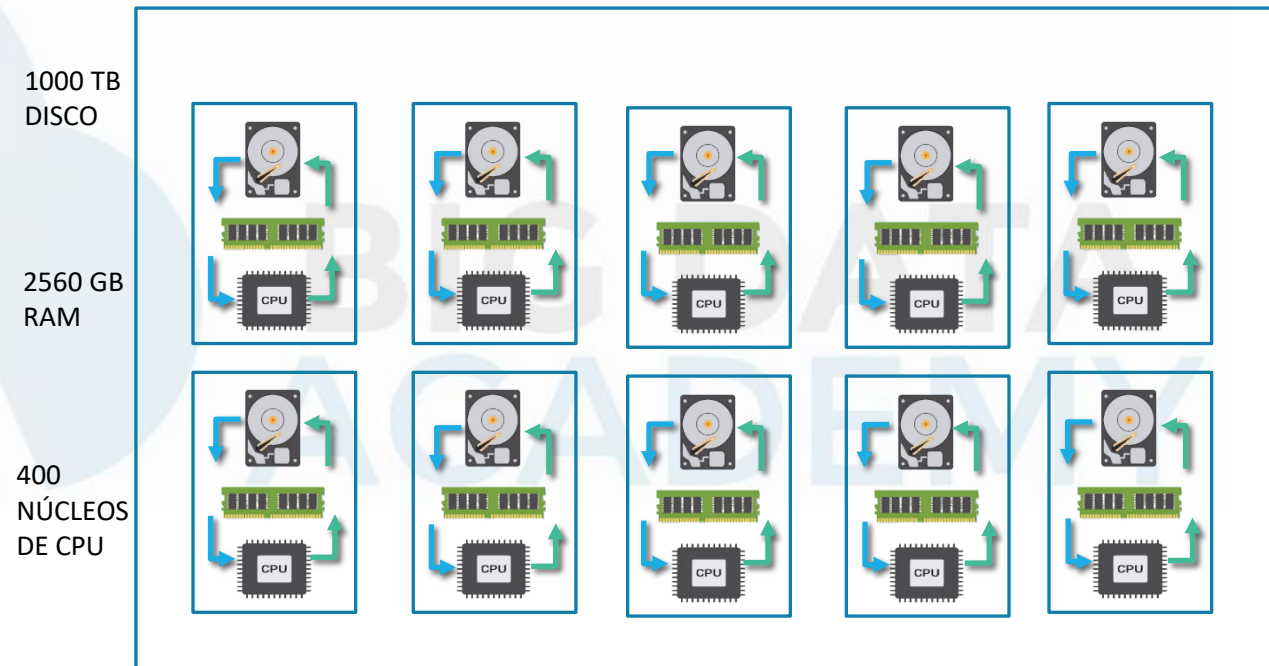
Clúster computacional

¿Cómo trabaja una computadora?



Un servidor común

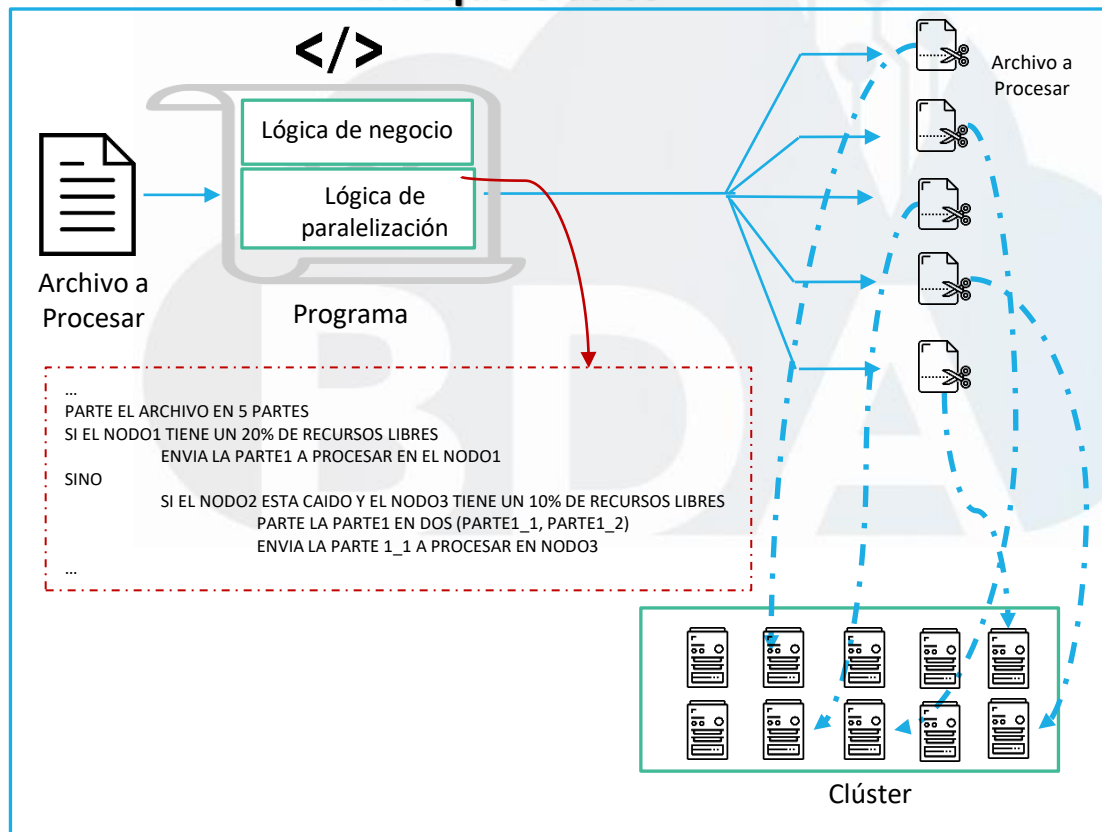
¿Cómo trabaja un clúster?



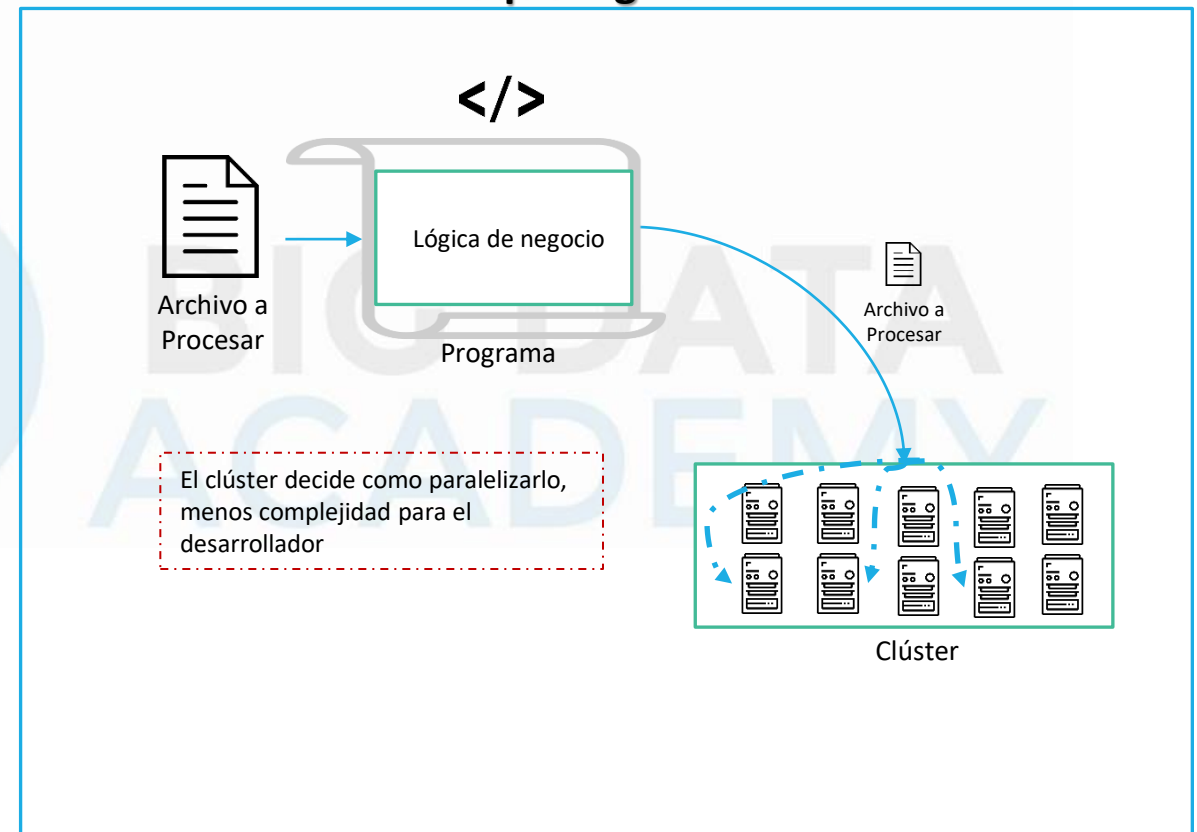
Un clúster es la suma de los recursos computacionales de los servidores que lo conforman, es como si tuviésemos una “super-computadora”

La paralelización

Enfoque Clásico

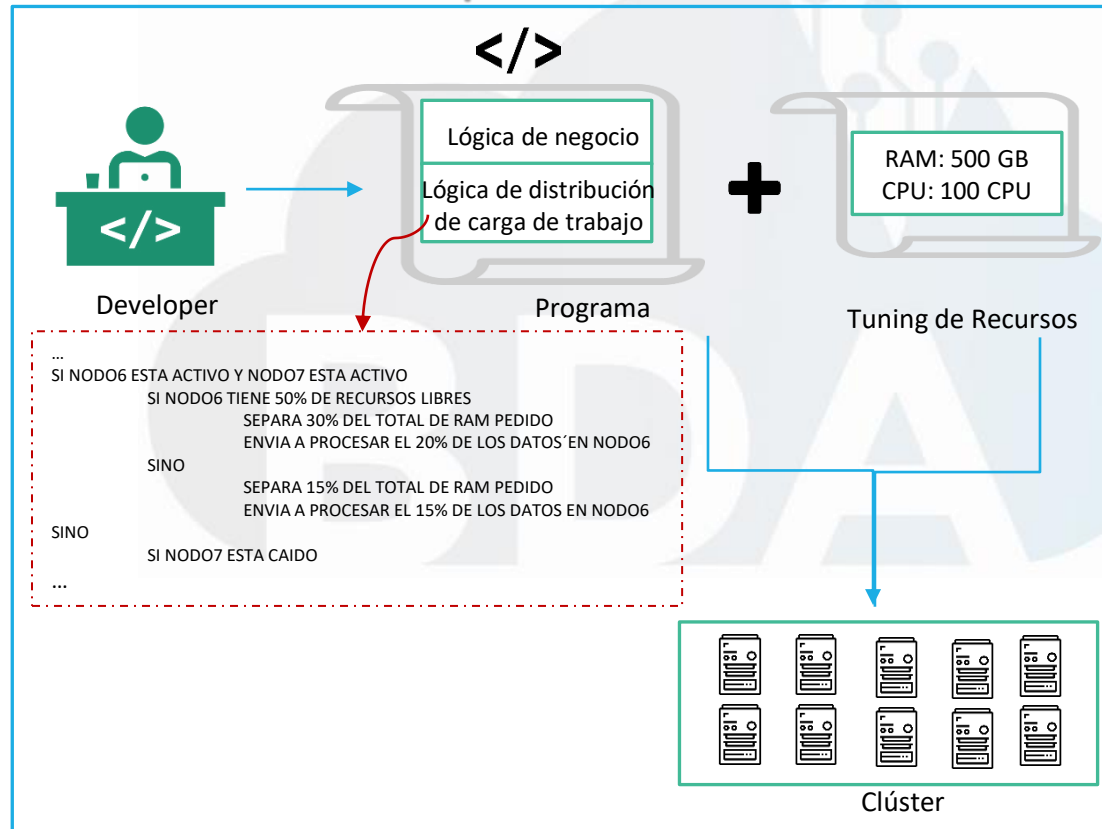


Enfoque Big Data

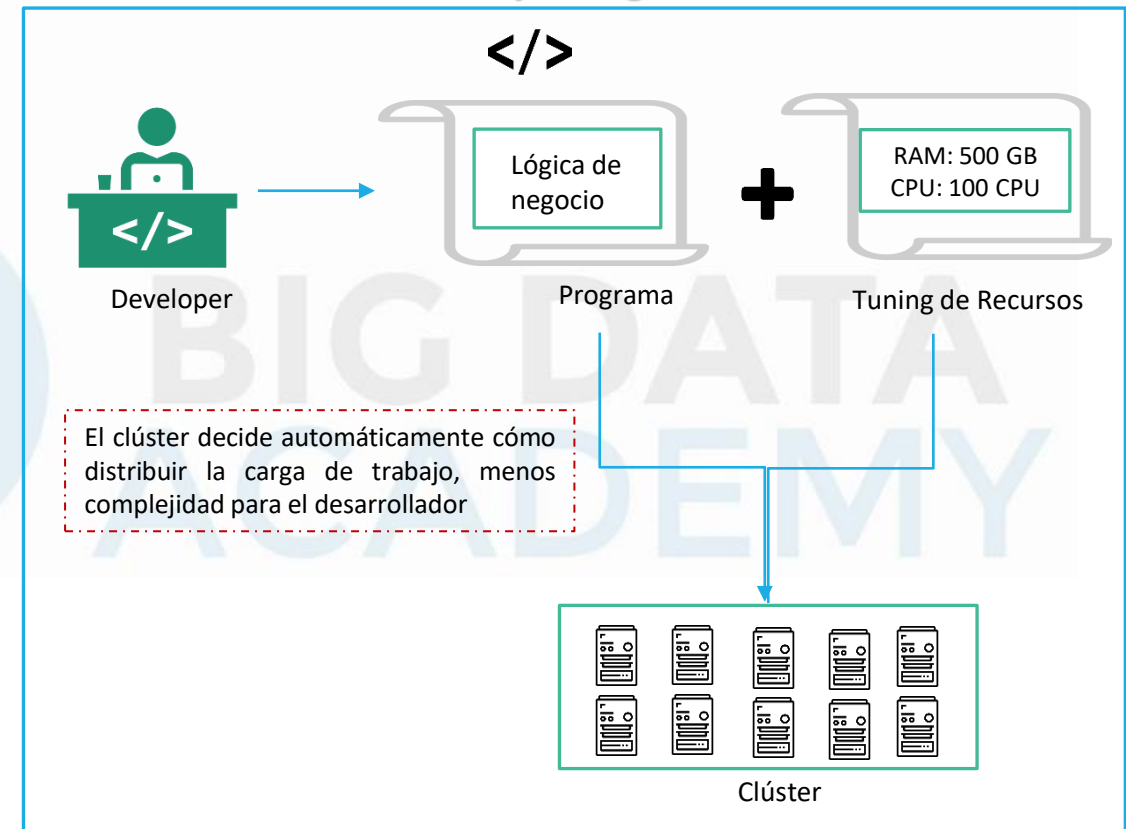


La distribución de carga de trabajo

Enfoque Clásico



Enfoque Big Data



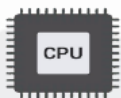
Escalabilidad (de proceso)

Enfoque Clásico



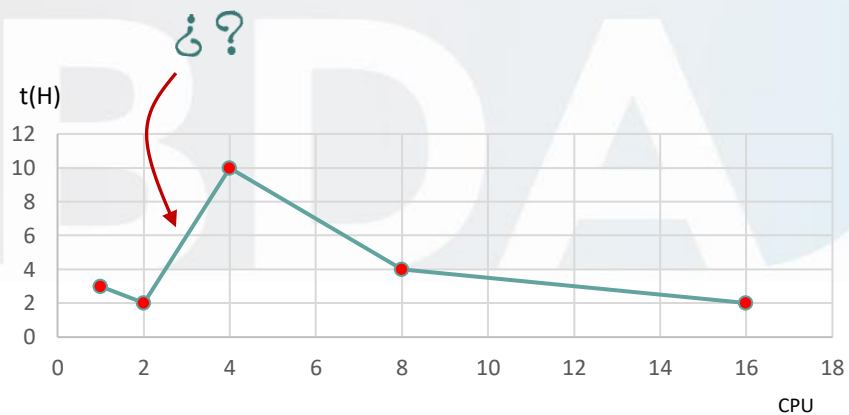
Programa

+



X CPU

CPU	t
1	3H
2	2H
4	10H
8	4H
16	2H

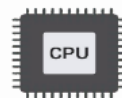


Enfoque Big Data



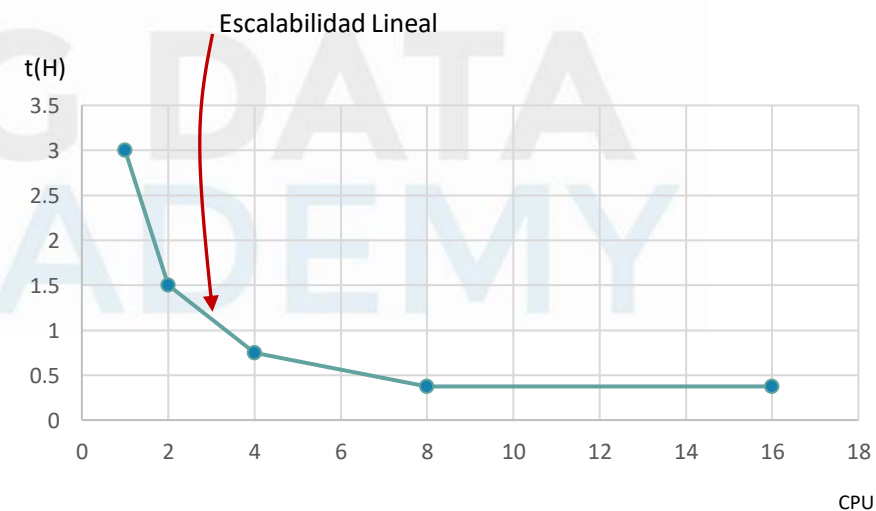
Programa

+



X CPU

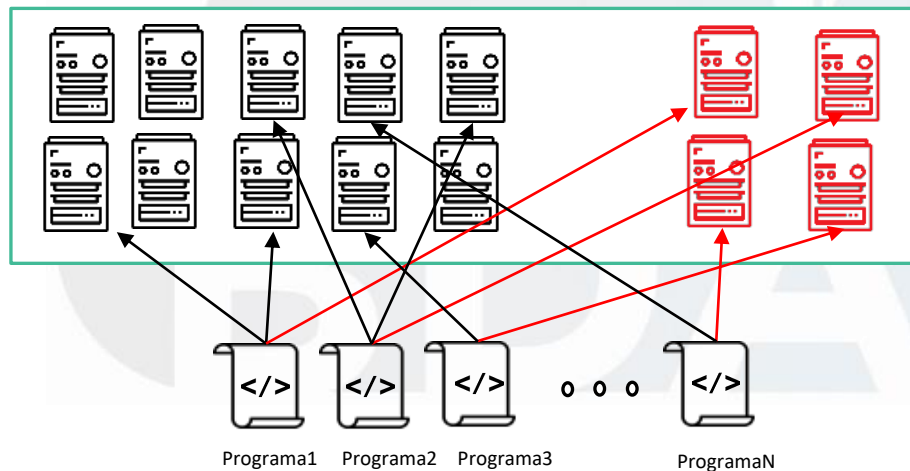
CPU	t
1	3H
2	1.5H
4	0.75H
8	0.375 H
16	0.375H



Escalabilidad (de hardware)

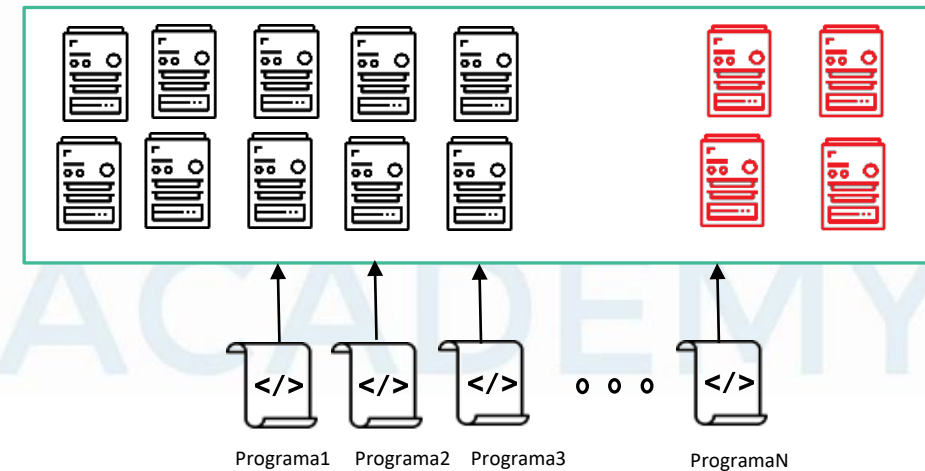
Enfoque Clásico

Agreguemos más nodos al clúster



Modificar los programas para que hagan uso de los nuevos nodos

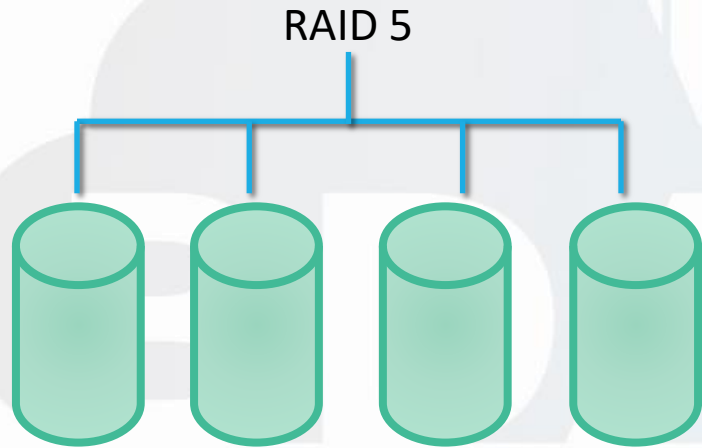
Enfoque Big Data



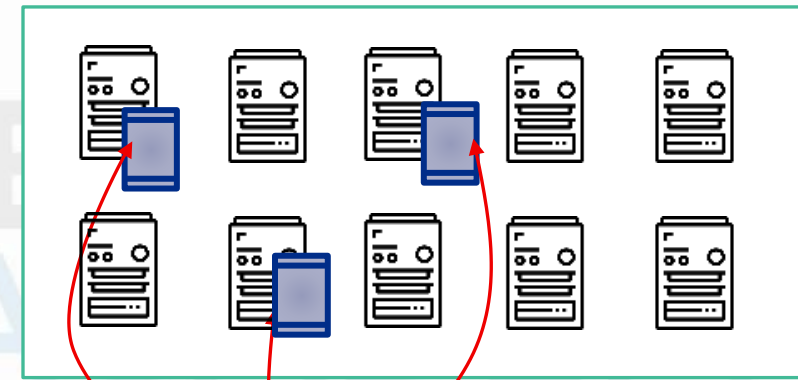
No es necesario modificar los programas, el clúster es visto como un todo

Alta disponibilidad (de datos)

Enfoque Clásico



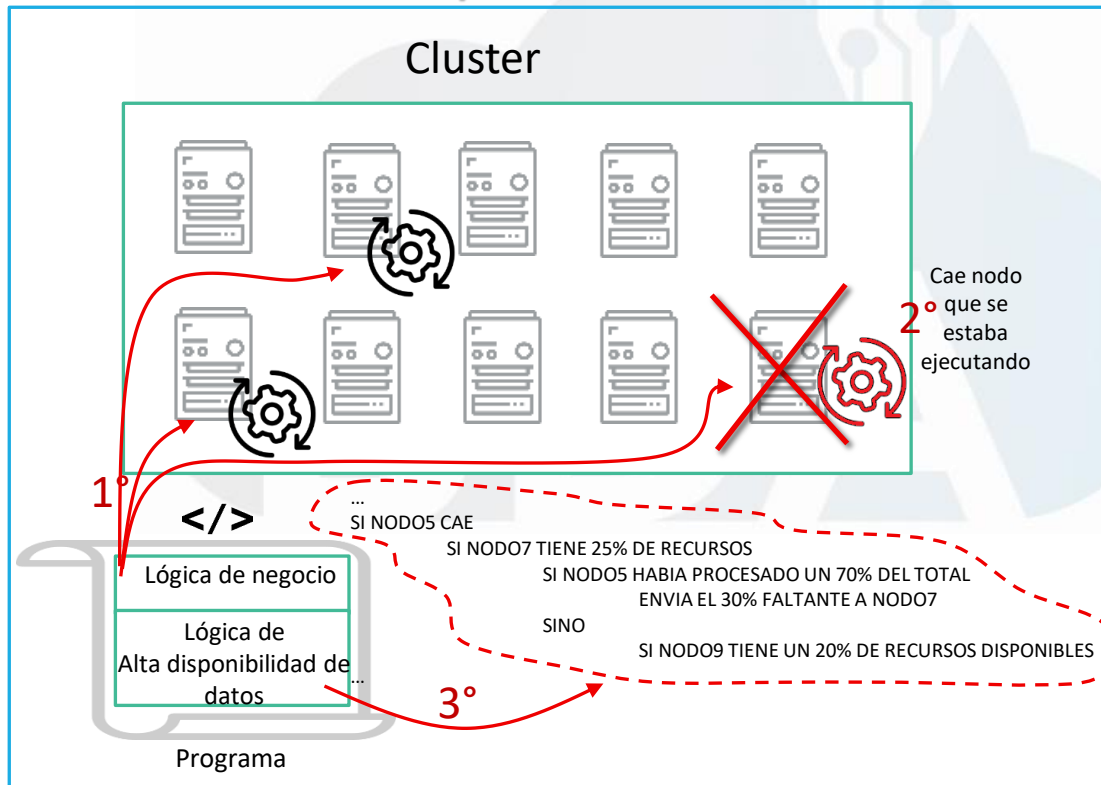
Enfoque Big Data



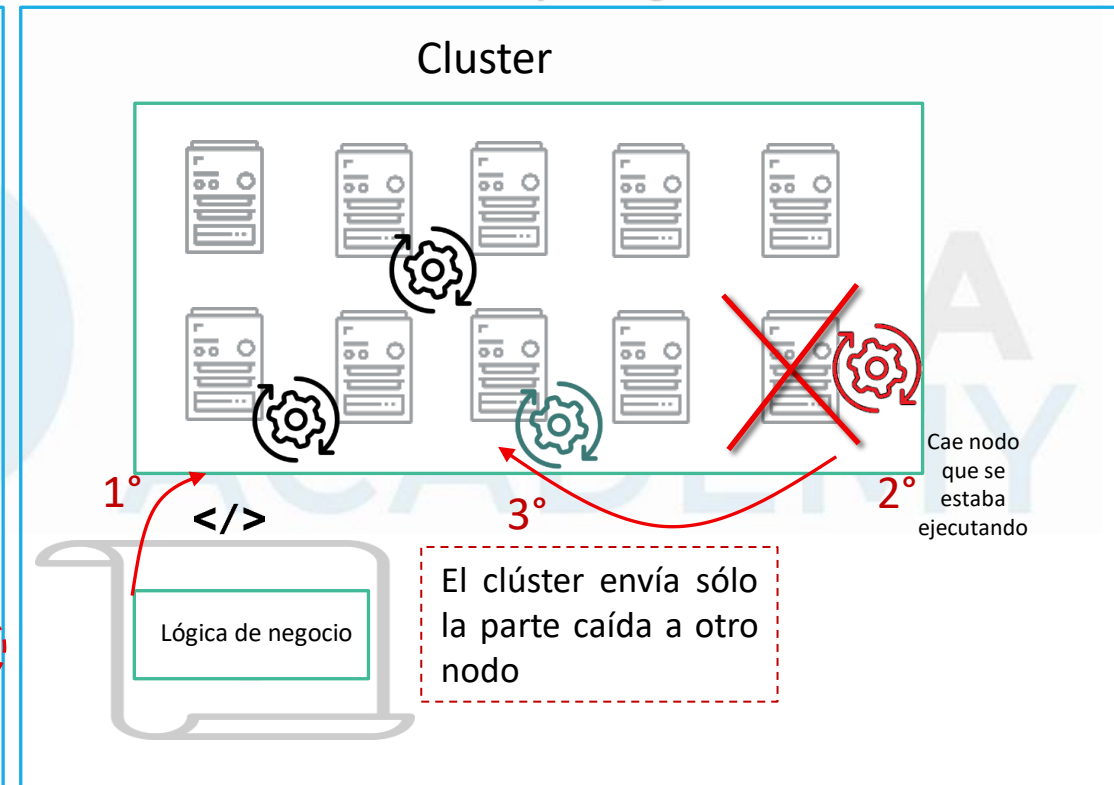
Replicación de datos

Alta disponibilidad (de proceso)

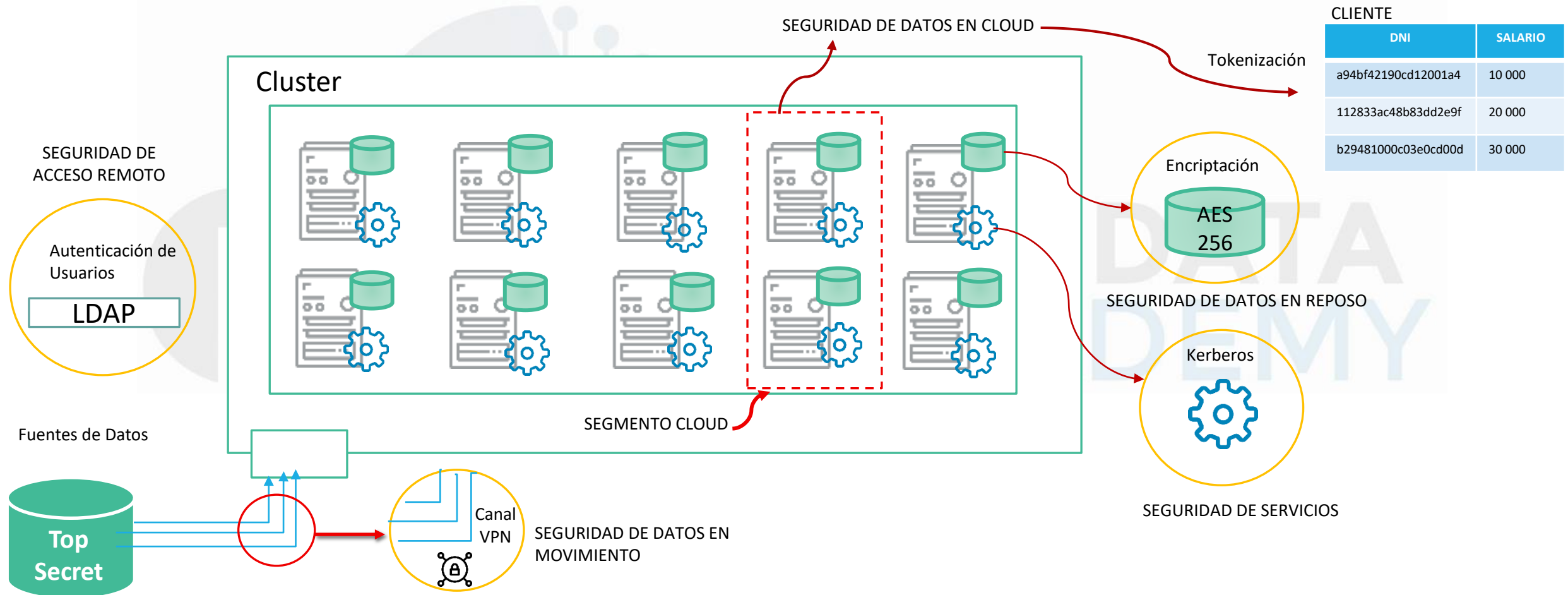
Enfoque Clásico



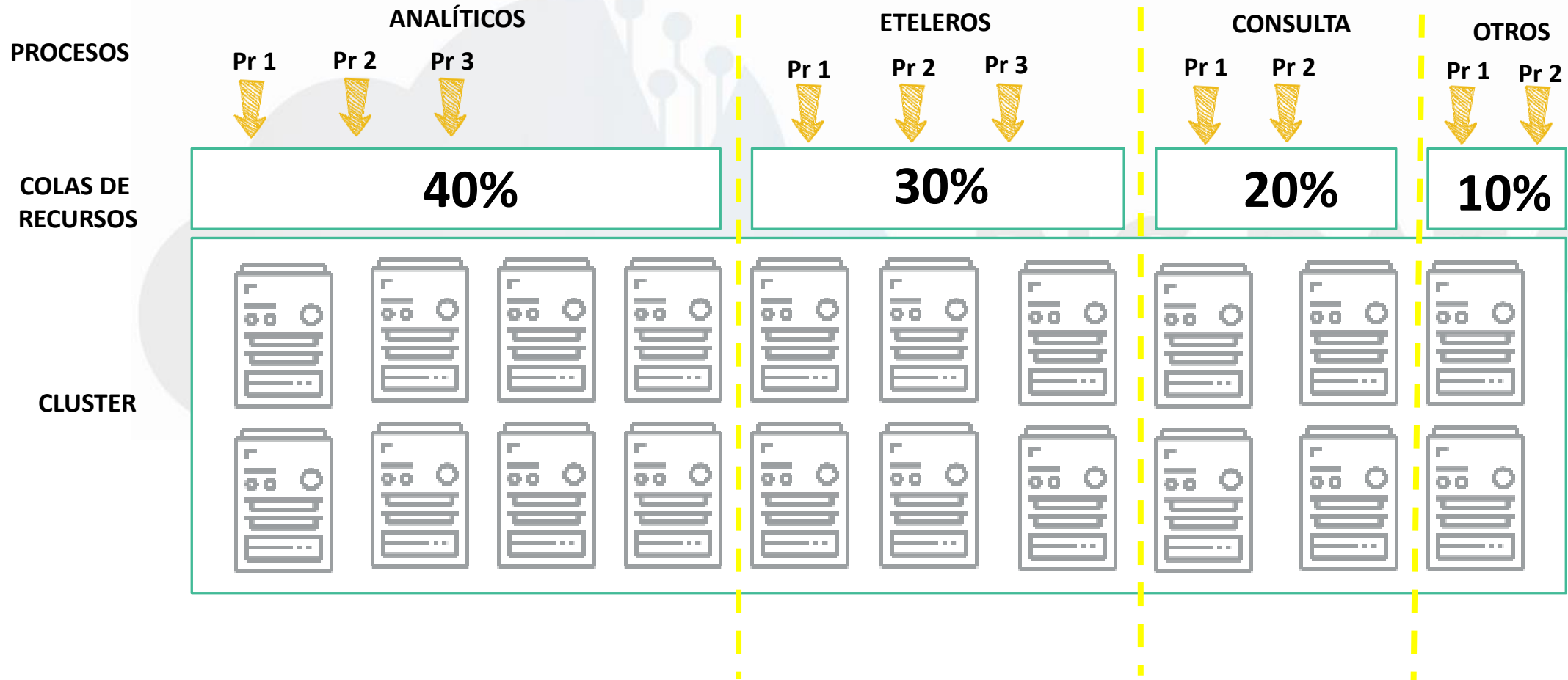
Enfoque Big Data



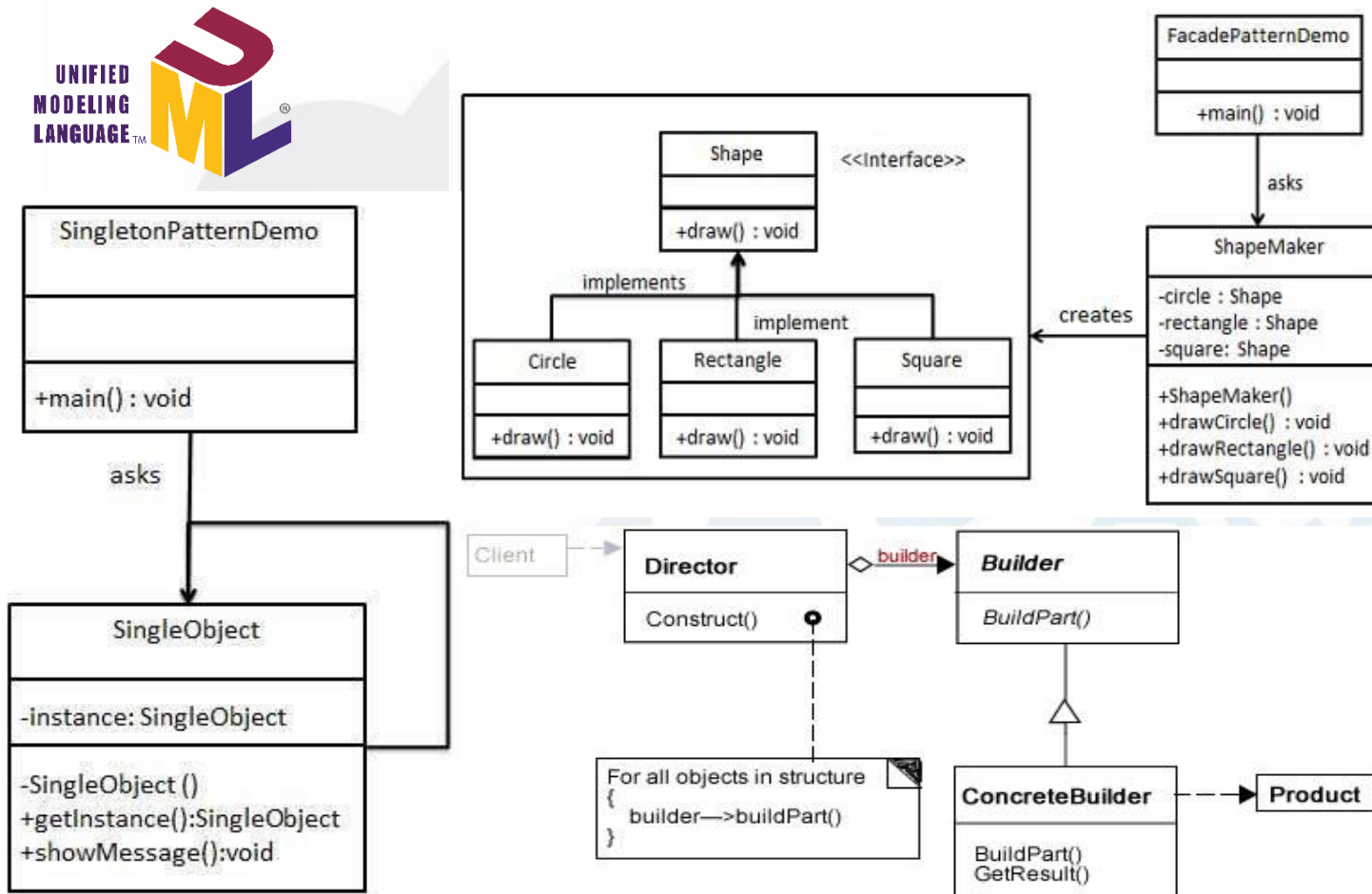
Seguridad



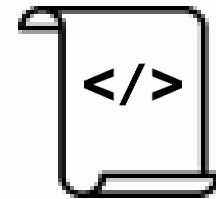
Gobierno



Patrones de diseño



PROGRAMACIÓN

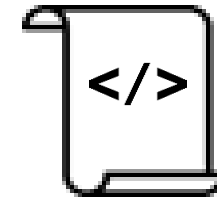


P1



P2

...



PN

Big Data: Tecnologías

A large, faint, light blue watermark of the Big Data Academy logo is visible in the background of the slide, behind the main title.

Tipos de tecnologías

INGESTA



ALMACENAMIENTO



PROCESAMIENTO



EXPLORACIÓN



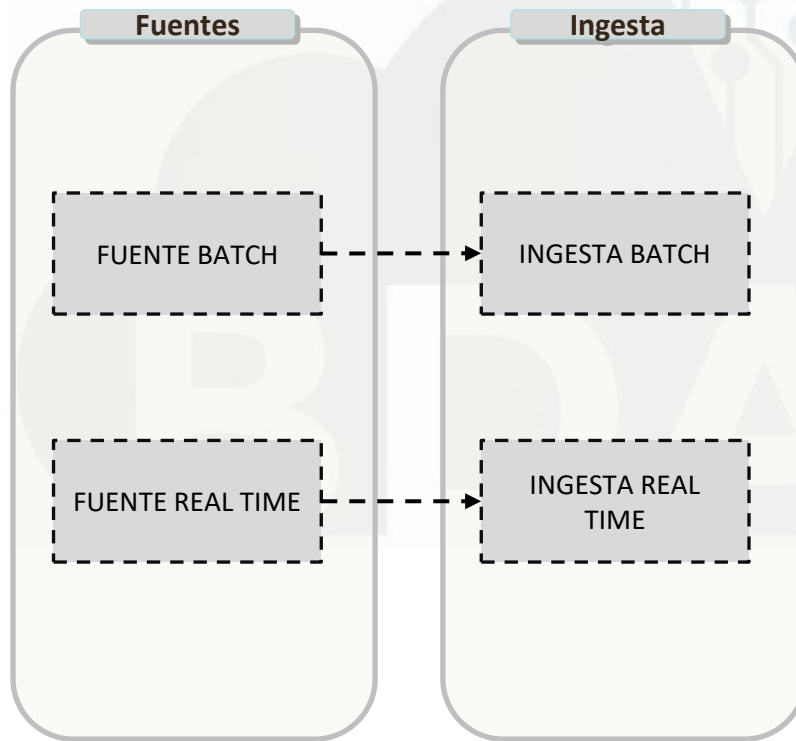
GOBIERNO



SEGURIDAD



Tecnologías de ingesta



Quiero ingestar de forma batchera

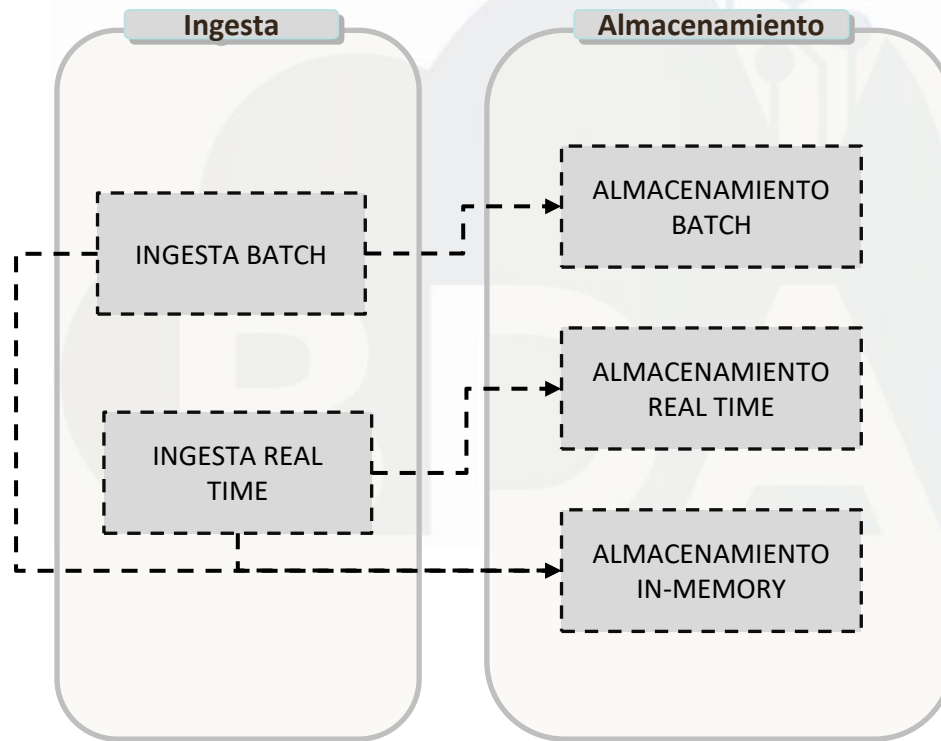


IBM DataStage

Quiero ingestar de forma real time



Tecnologías de almacenamiento



Quiero almacenar de forma batchera



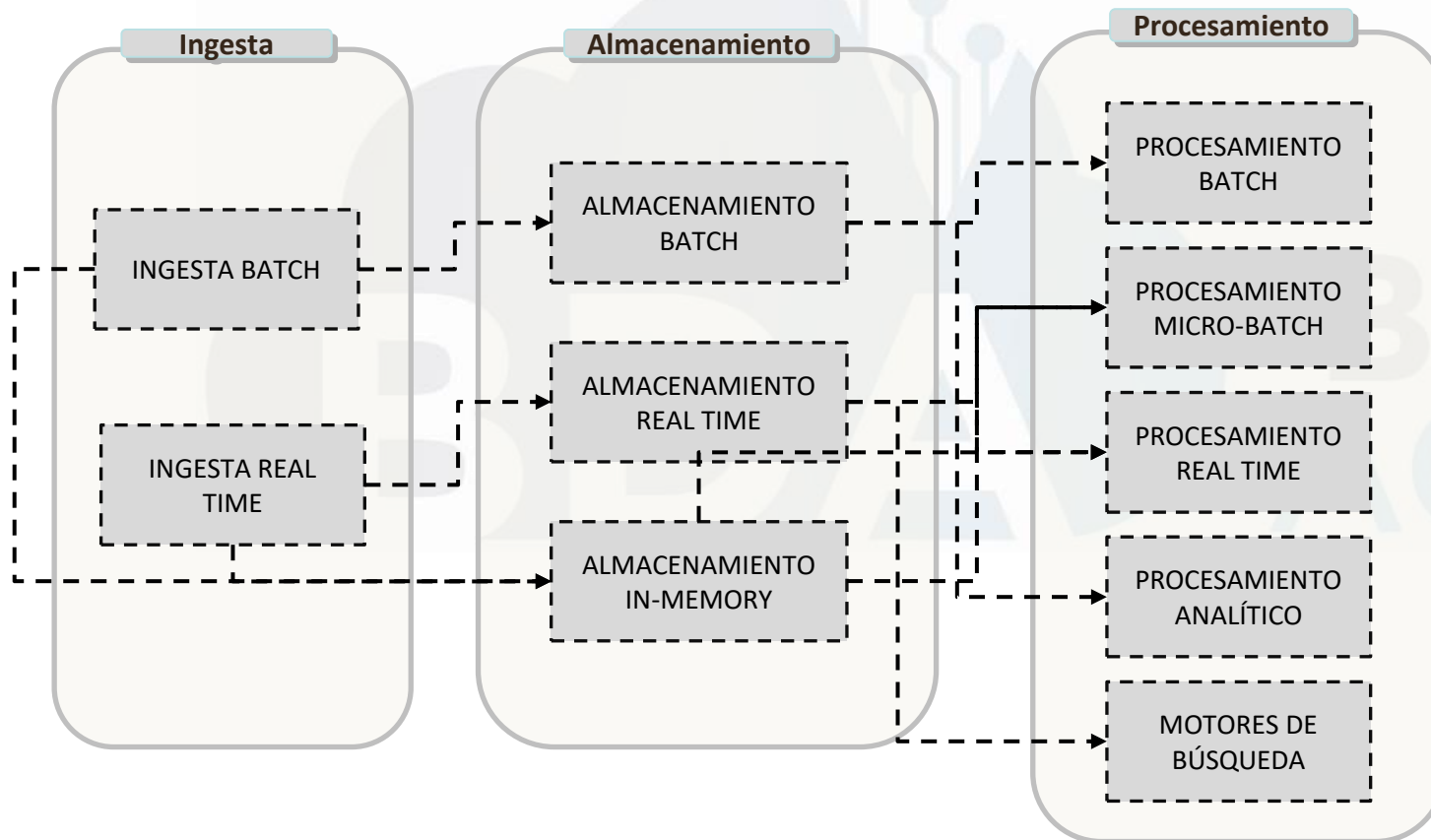
Quiero almacenar de forma real time



Quiero almacenar in-memory



Tecnologías de procesamiento



Quiero procesar de forma batchera



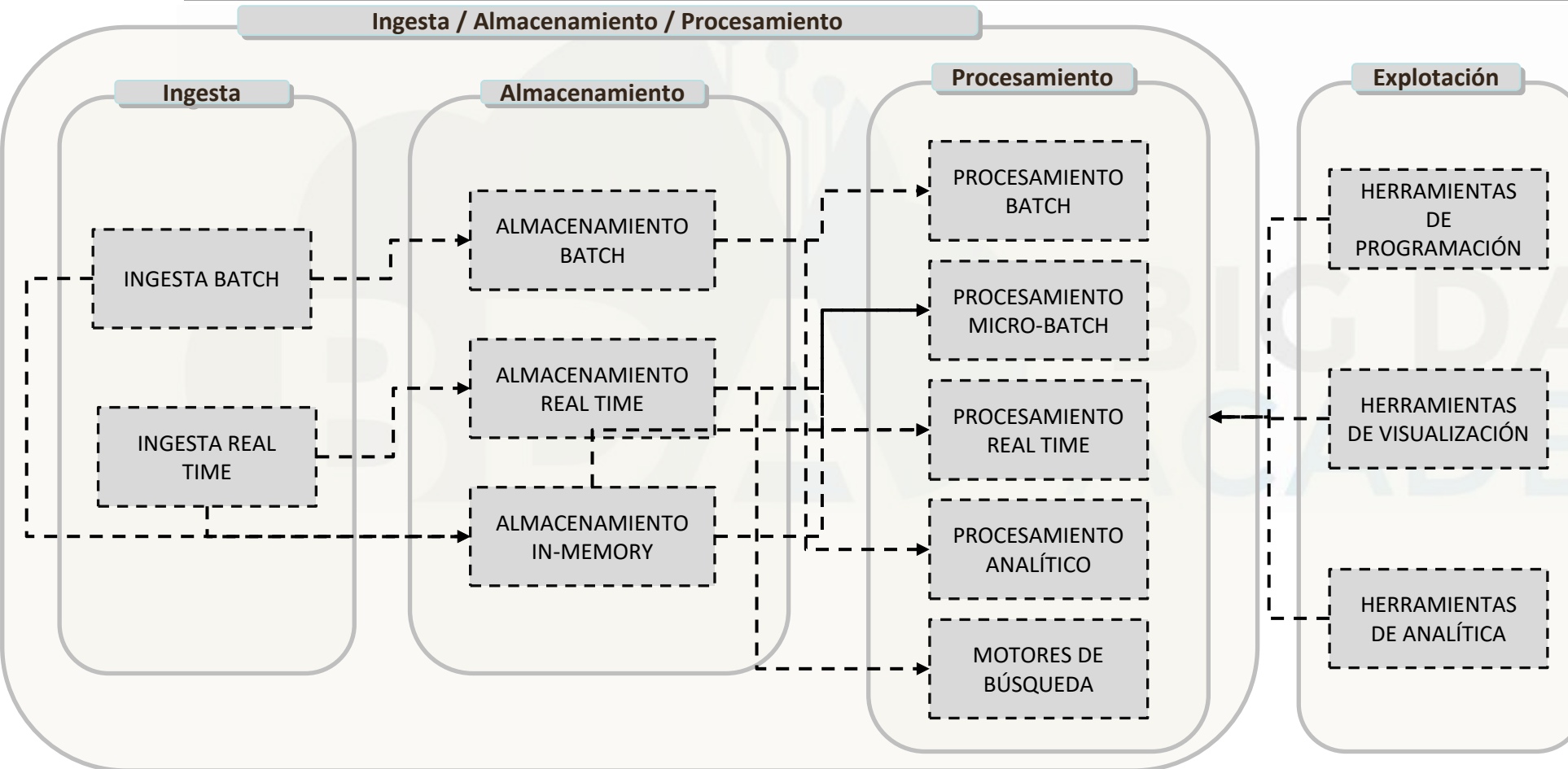
Quiero procesar de forma real time



Quiero procesar de forma analítica



Tecnologías de explotación



Quiero programar



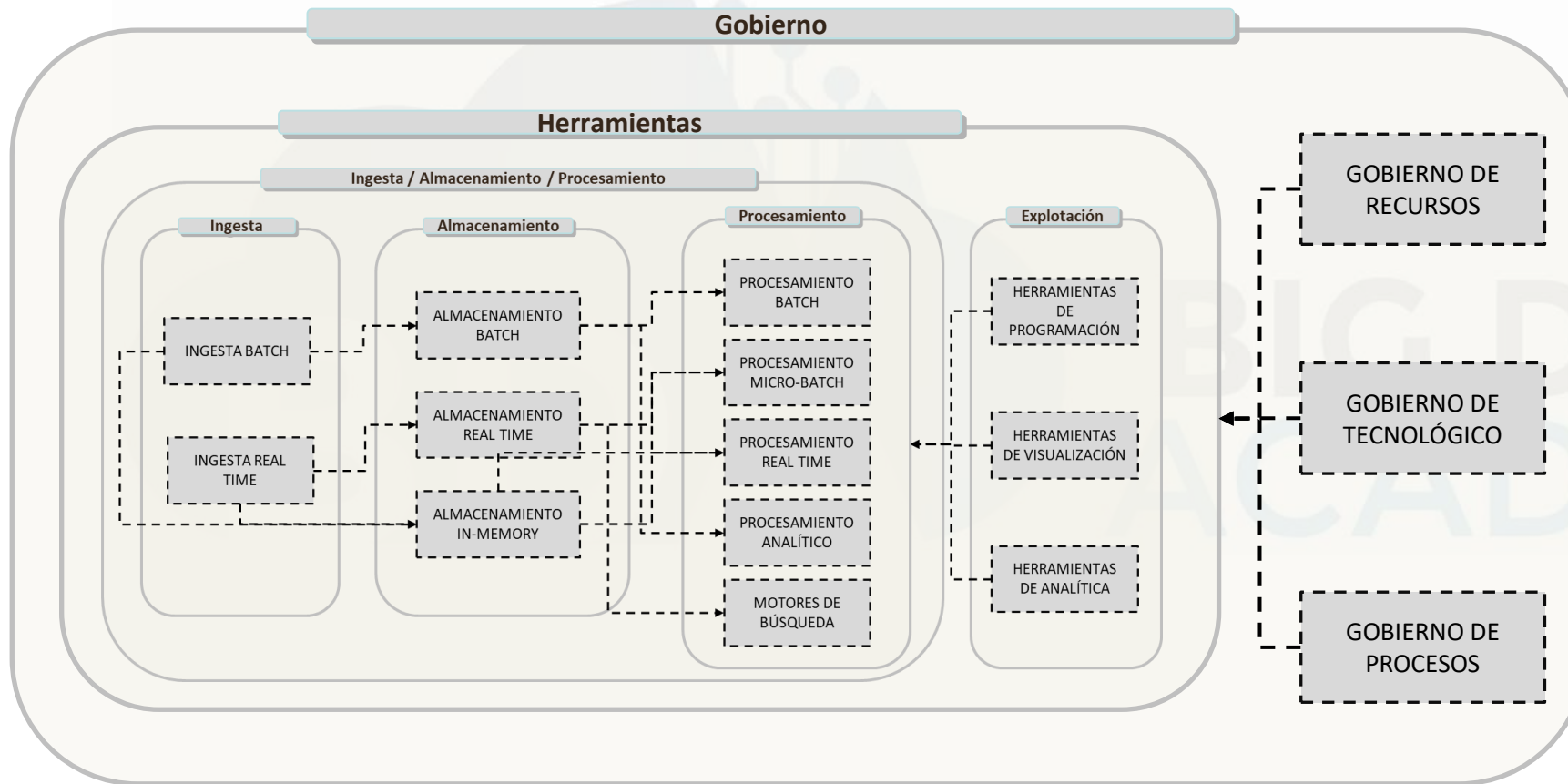
Quiero pintar gráficos y reportes



Quiero analizar datos



Tecnologías de gobierno



Quiero administrar el acceso a RAM, CPU, prioridades



Quiero administrar las herramientas instaladas



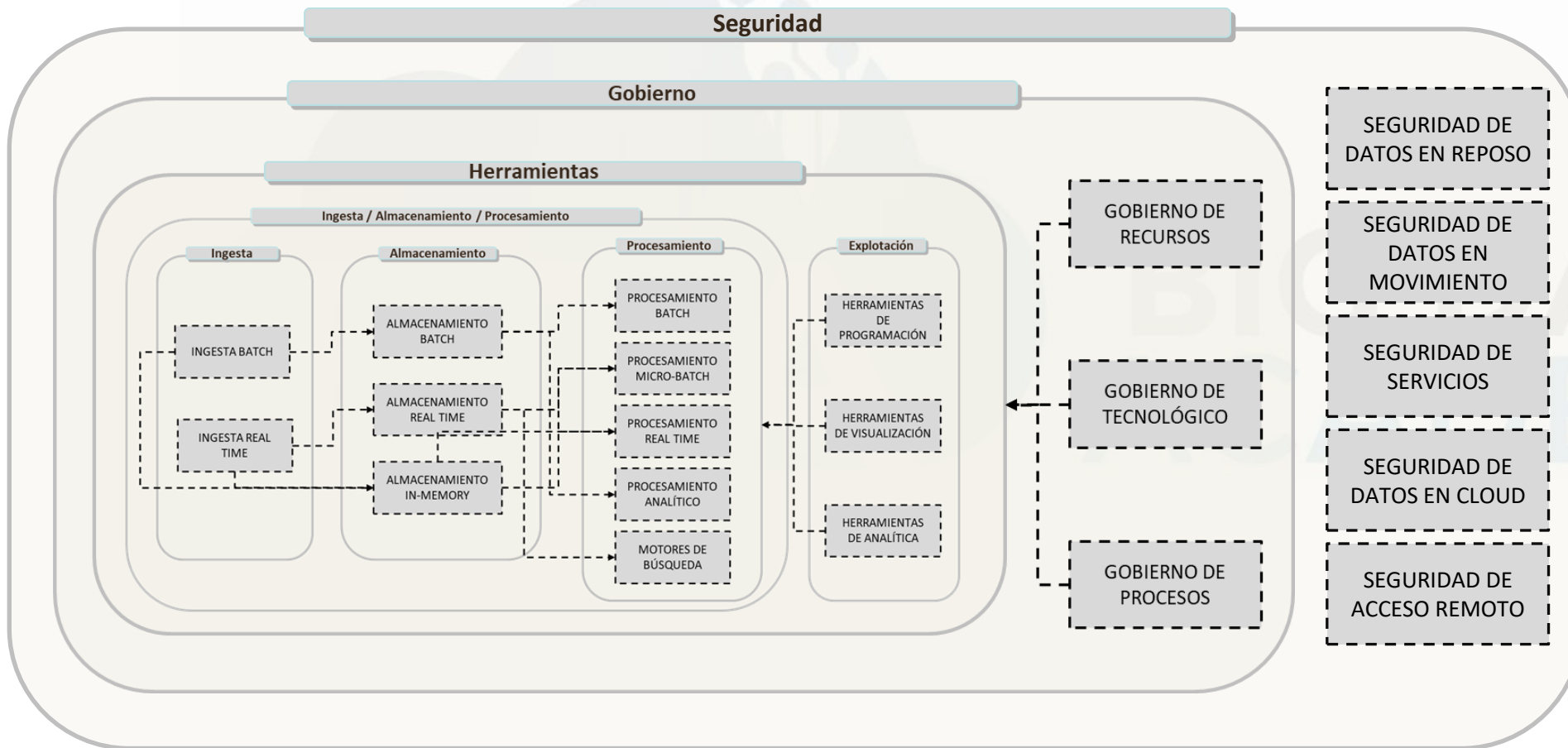
Quiero administrar el ciclo de desarrollo y puesto en producción



Jenkins



Tecnologías de seguridad



Quiero asegurar datos en mis discos duros



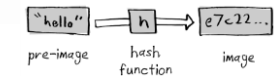
Quiero asegurar los datos que viajan por la red



Quiero asegurar el acceso a servicios y procesos



Quiero asegurar mis datos en la nube



Quiero asegurar el acceso al clúster



Lenguajes de programación

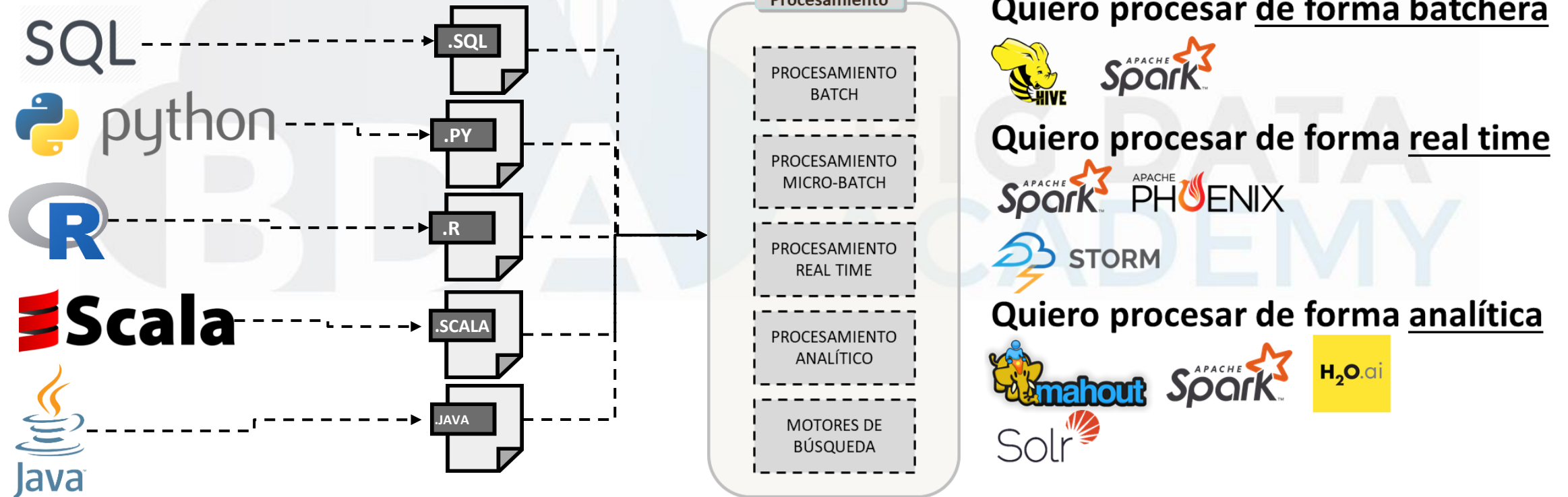


Lenguajes de programación

LENGUAJES

PROGRAMAS

MOTORES DE PROCESAMIENTO



Características

SQL

- Orientado a datos estructurados
- Aprendizaje fácil
- Estándar para consultas a bases de datos

 python

- Orientado a la programación funcional
- Interpretado, no tipado y muy flexible
- Librerías casi para todo

 R

- Orientado a la estadística
- Interpretado
- Gráficos avanzados simples de realizar

 Scala

- Orientado a la programación funcional y de objetos
- Compilado y no tipado
- Fork de Java

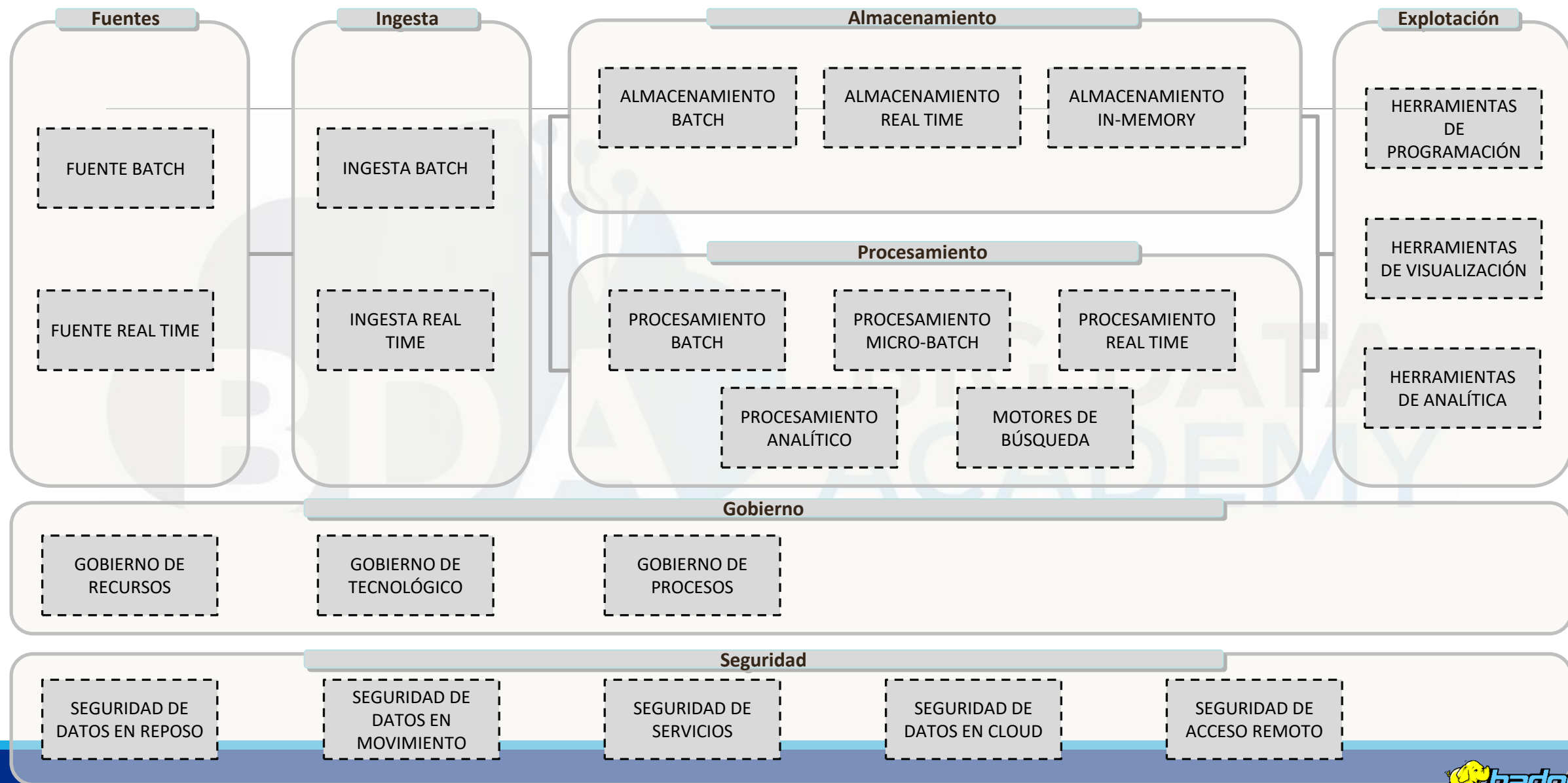
 Java

- Orientado a la programación de objetos
- Compilado y tipado
- Aprovecha muy bien los servidores con grandes recursos computacionales

Arquitectura general de Big Data



Arquitectura conceptual



Tipos de tecnologías

INGESTA



ALMACENAMIENTO



PROCESAMIENTO



EXPLOTACIÓN



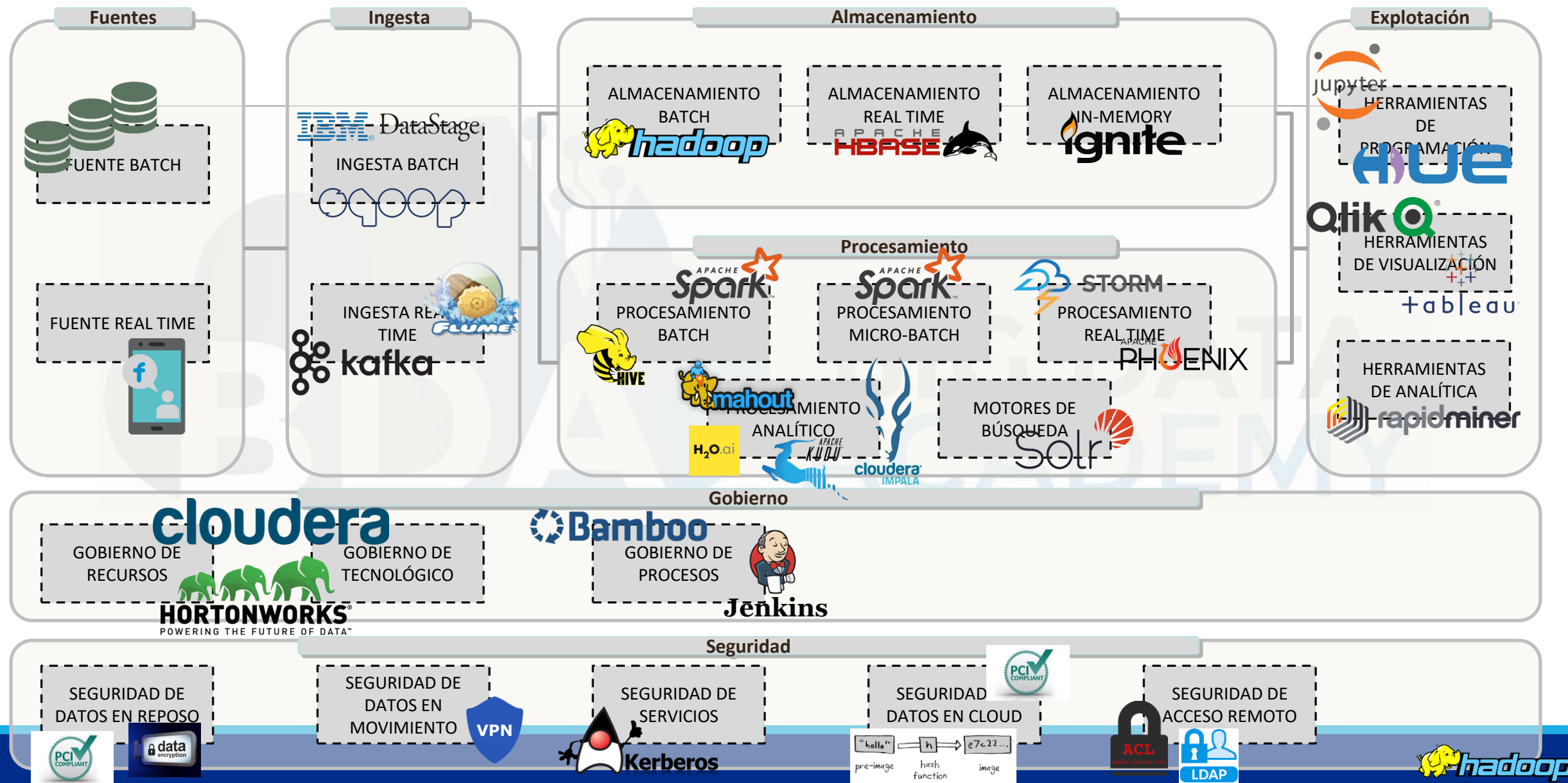
GOBIERNO



SEGURIDAD



Arquitectura tecnológica



Ejercicios teóricos

Ejercicios teóricos

1. ¿Qué es Big Data?
2. ¿Cuál es el objetivo del Big Data?
3. ¿Cuáles son las 5V?
4. ¿Cuál es la diferencia entre variedad y variabilidad?
5. ¿Qué es un clúster computacional?
6. ¿Qué es la paralelización?
7. ¿Qué es la escalabilidad?
8. ¿Cuáles son los puntos claves en la seguridad de un clúster?
9. ¿Qué tipos de tecnologías existen en el Big Data?
10. ¿Qué motores de procesamiento conoce?
11. ¿Qué lenguaje de programación es mejor para programar sobre Big Data?
12. ¿Por qué se prefiere un proceso micro-batch sobre uno real time?
13. Dibuje la arquitectura conceptual de Big Data

Resumen

Hablemos...

