



---

# Programa de Arquitectura para entornos de Big Data & Cloud

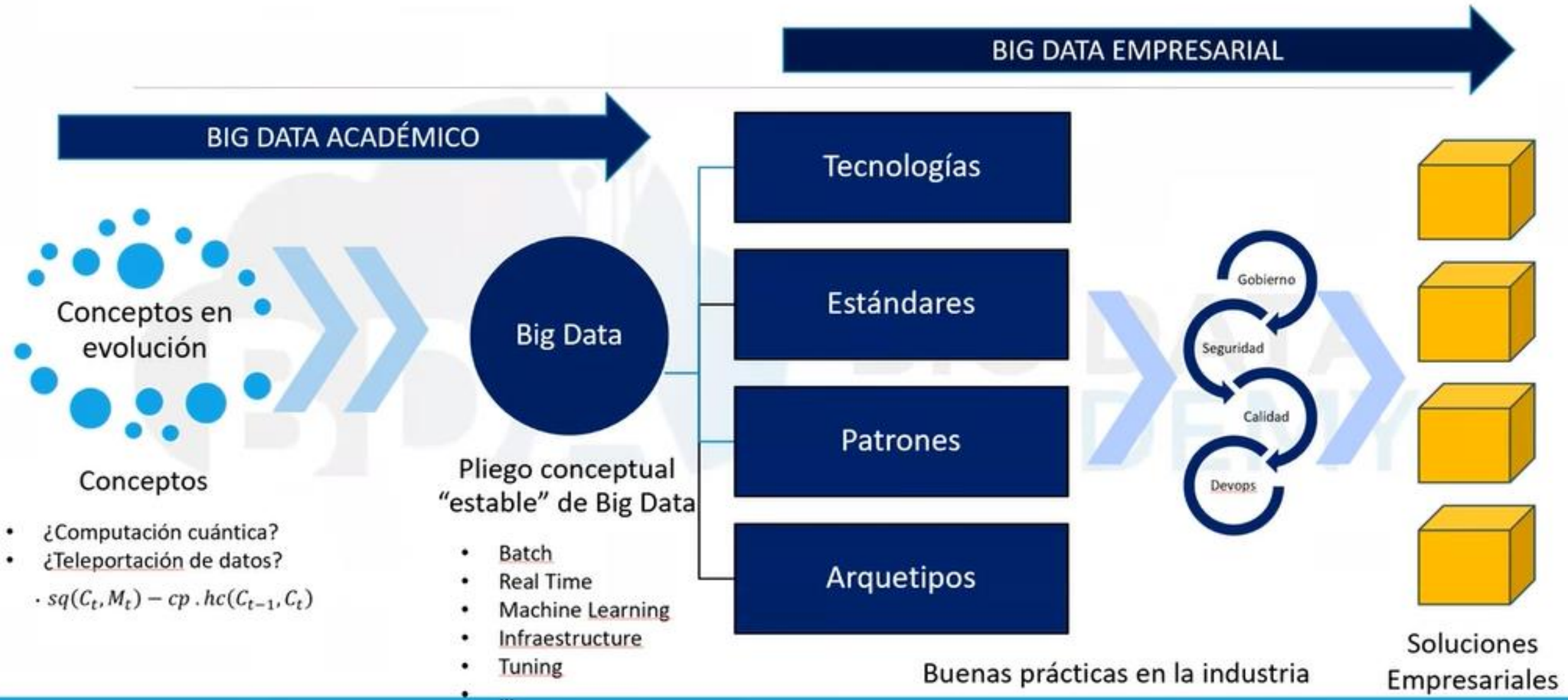
---

BIG DATA ACADEMY

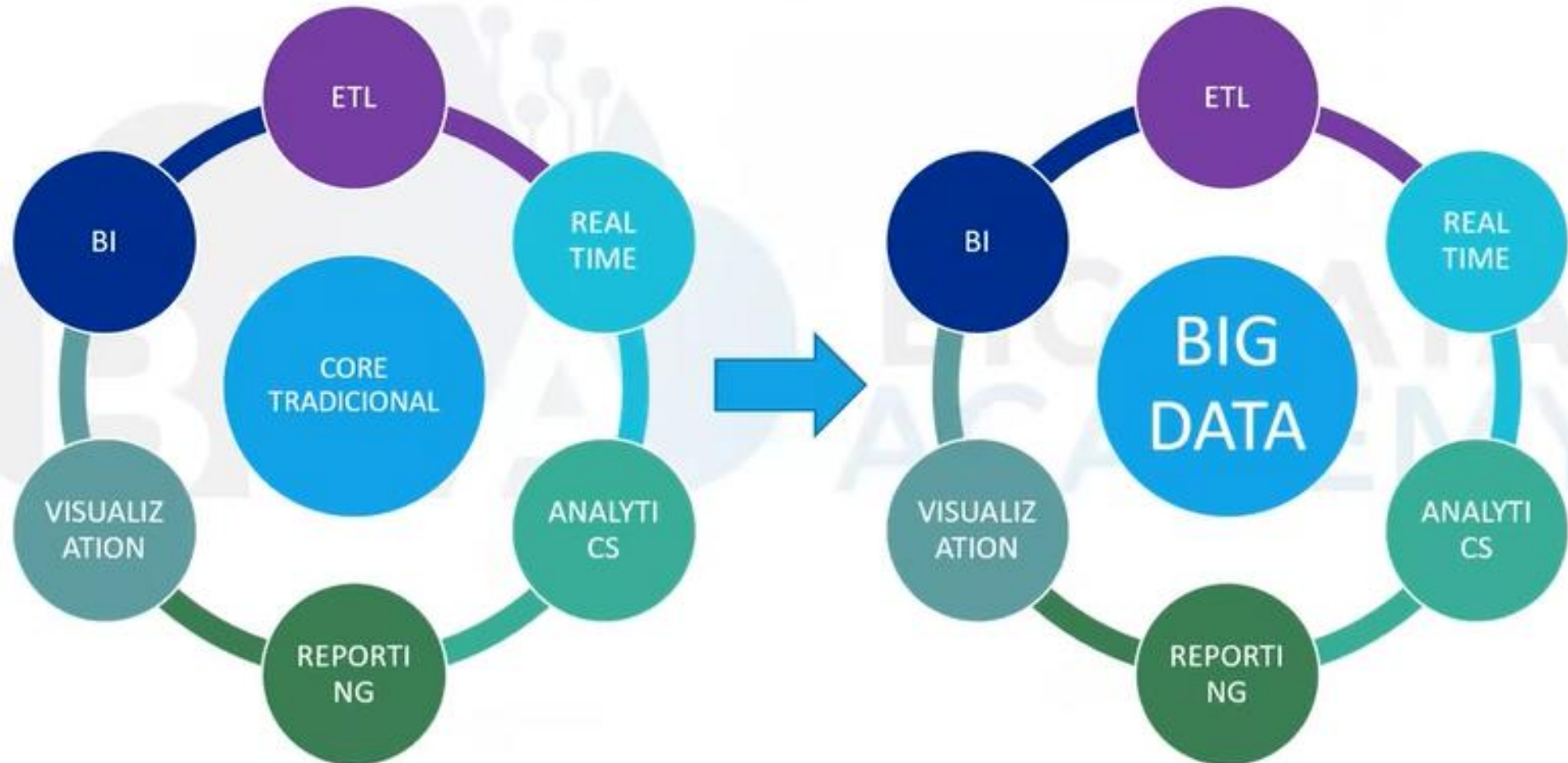
# Big Data aplicado a la empresa

HAGA CLIC PARA AGREGAR TEXTO

# ¿Big Data aplicado a las empresas?



# Concepción de un Proyecto Big Data





# Tecnologías en el mundo del Big Data

En 2016



## En 2020





# Propuesta de mix tecnológico



# Definición arquitectónica

HAGA CLIC PARA AGREGAR TEXTO



# Vistas arquitectónicas

Vista conceptual

Vista tecnológica

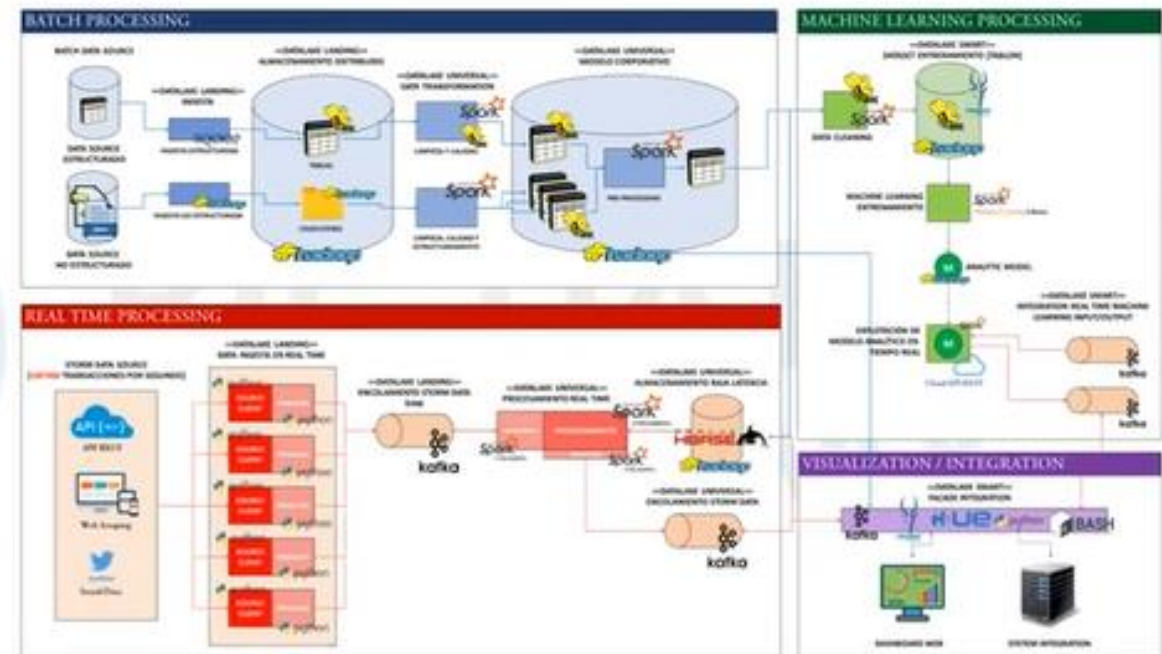
Vista de patrones de diseño

Vista de infraestructura

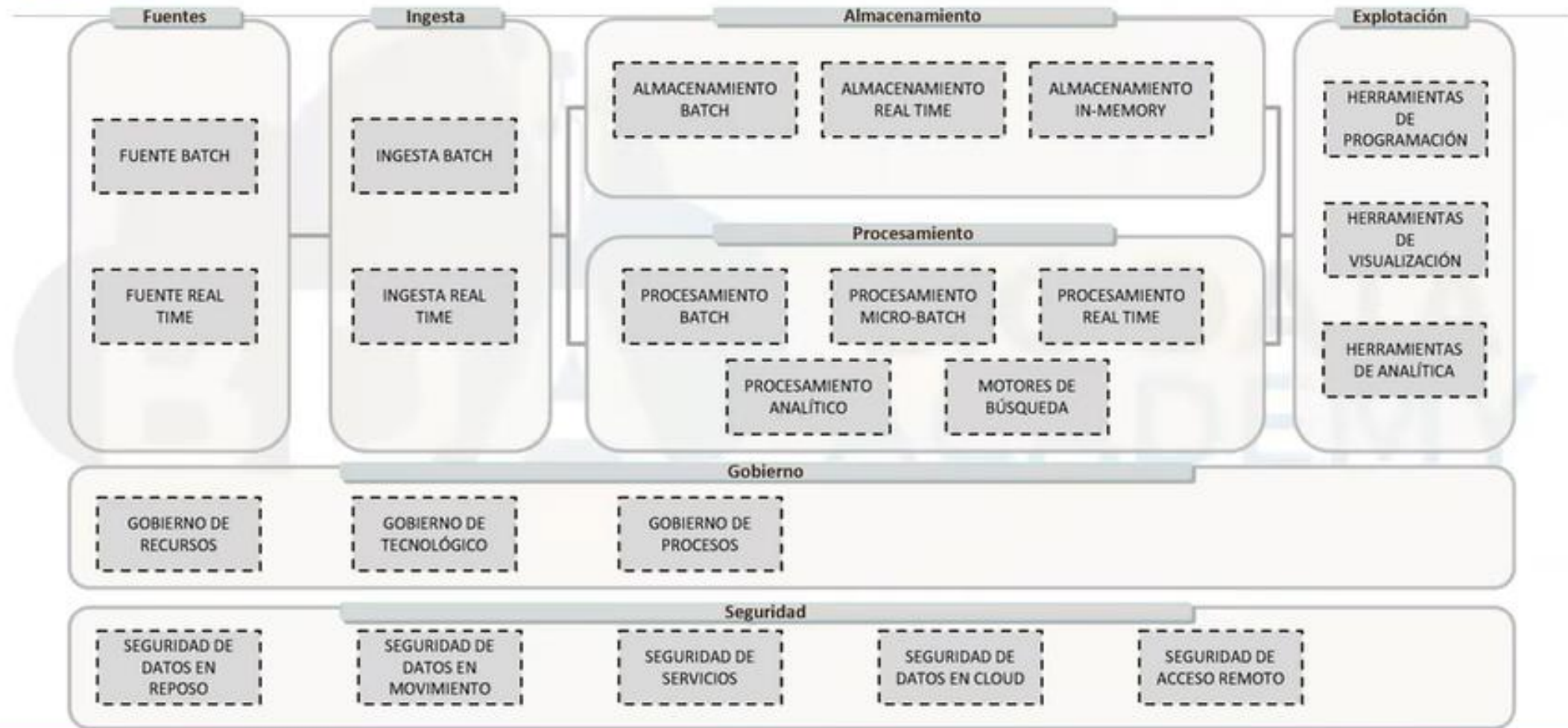
Vista de gobierno

Vista de seguridad

Vista de modelamiento, calidad, devops, ...

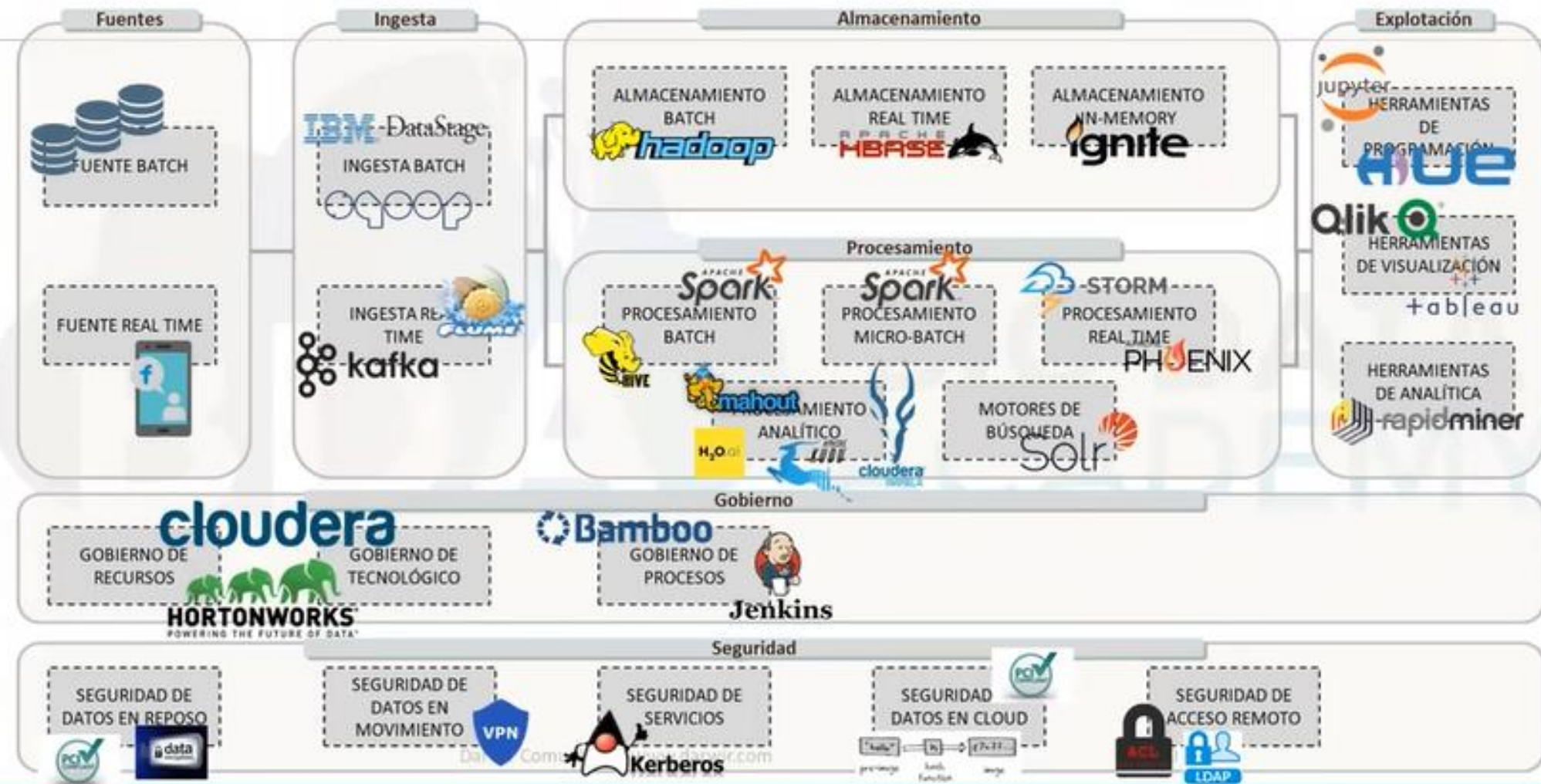


# Arquitectura conceptual general de Big Data



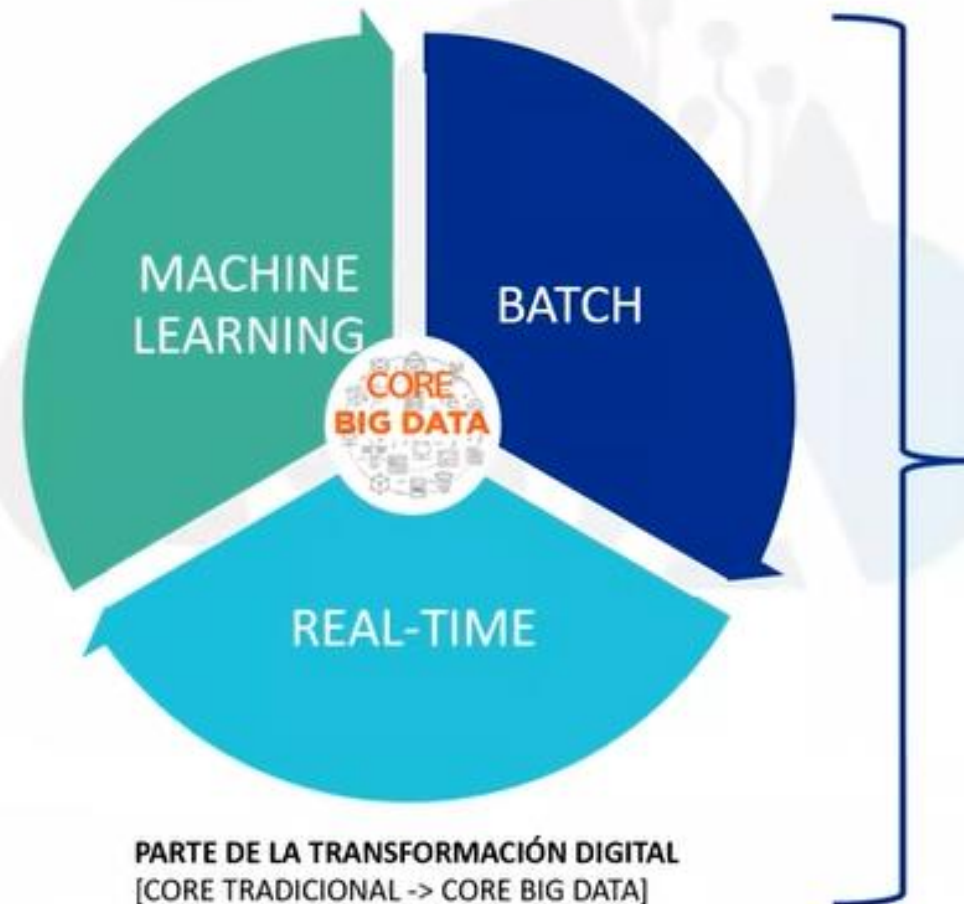


# Arquitectura tecnológica





# Vistas de una arquitectura de Big Data



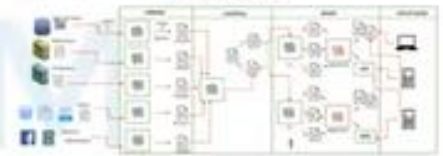
## VISTAS DEL CORE DE BIG DATA

- Vista de arquitectura
- Vista de gobierno
- Vista de infraestructura

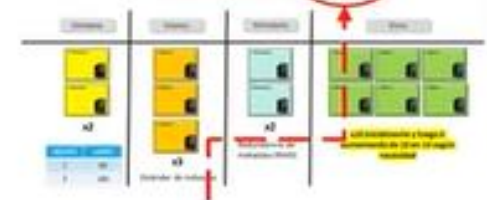
Conceptual y tecnológica



Datalake

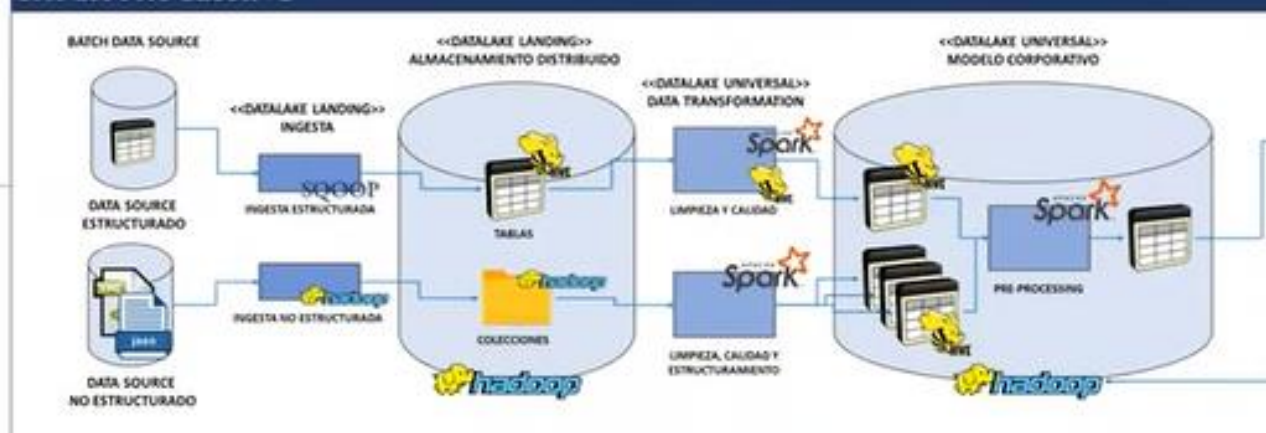


On-Premise y Cloud

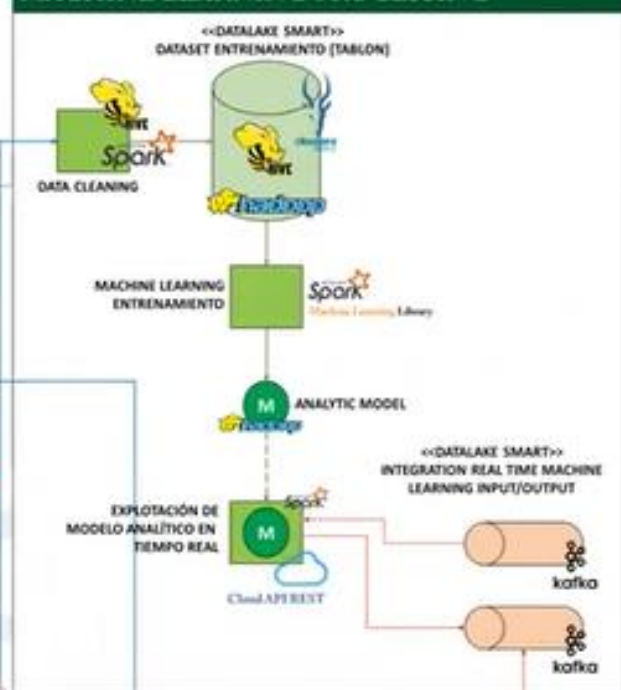


# Vista Conceptual con Patrones de Diseño de Big Data

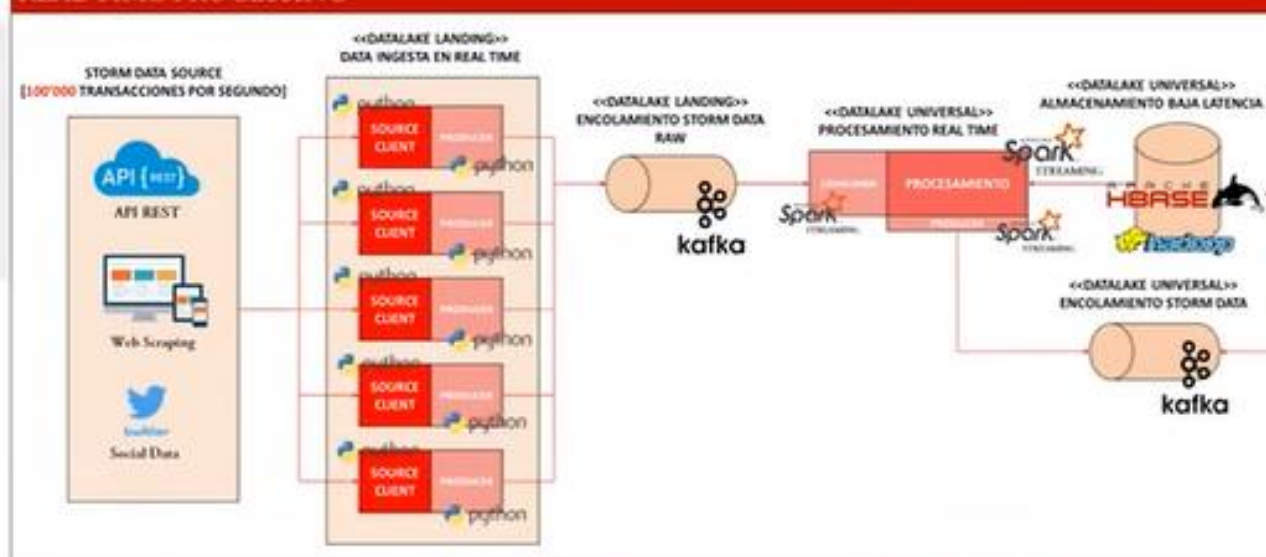
## BATCH PROCESSING



## MACHINE LEARNING PROCESSING



## REAL TIME PROCESSING



## VISUALIZATION / INTEGRATION

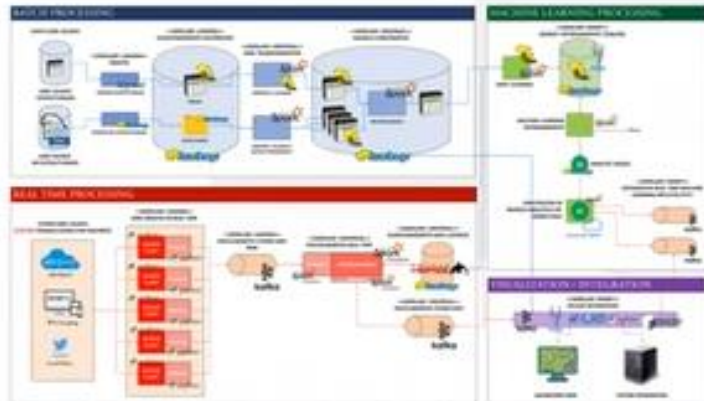


**¡Este diagrama sí lo podemos codificar!**



# Vista de Infraestructura

## Arquitectura tecnológica



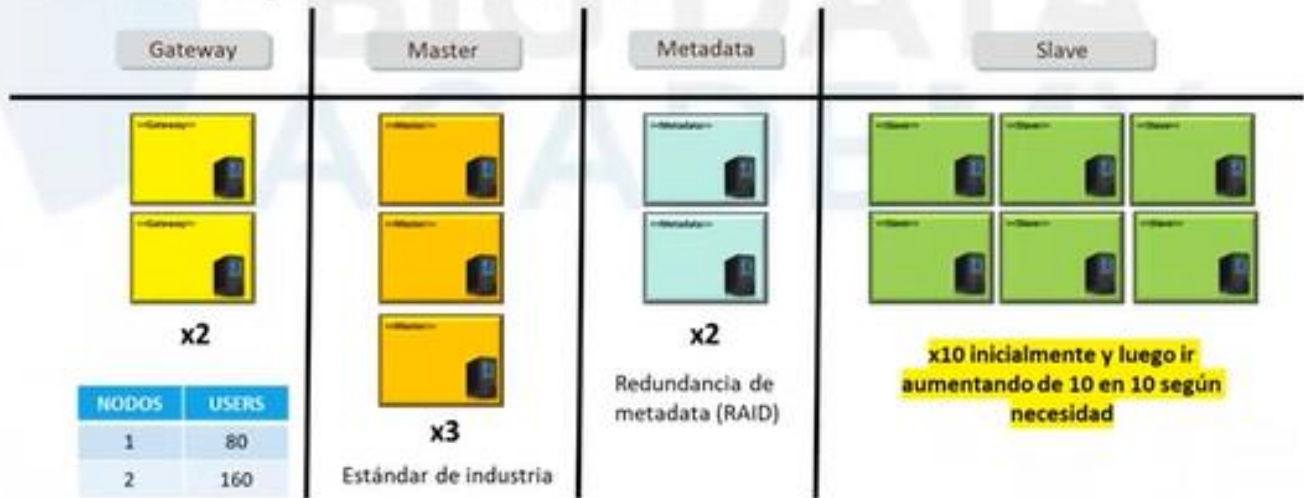
Debemos de traducir la necesidad tecnológica a Hardware:

- On-Premise
- Cloud IaaS
- Cloud PaaS
- Cloud SaaS

¿Qué características Hardware necesitamos?: RAM, CPU, DISCO, Ancho de Banda, GPU



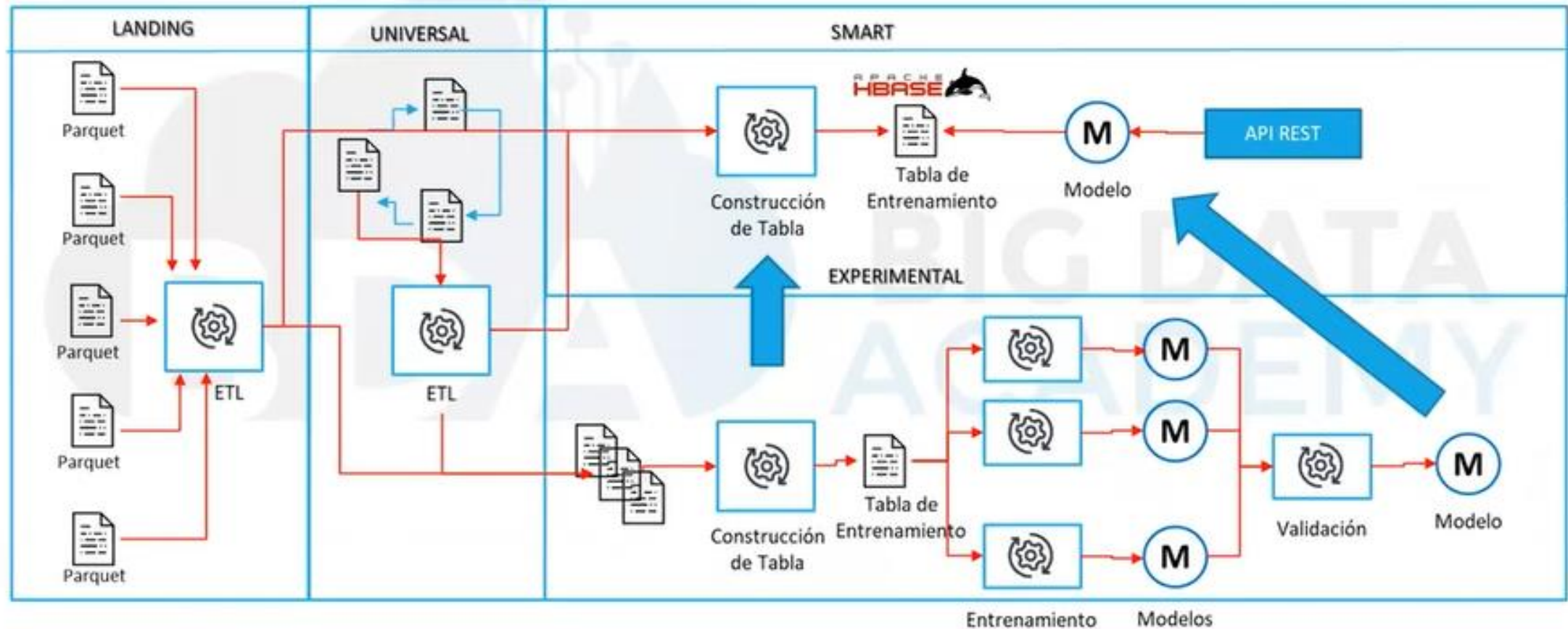
¿Cuántos roles y servidores necesitamos?



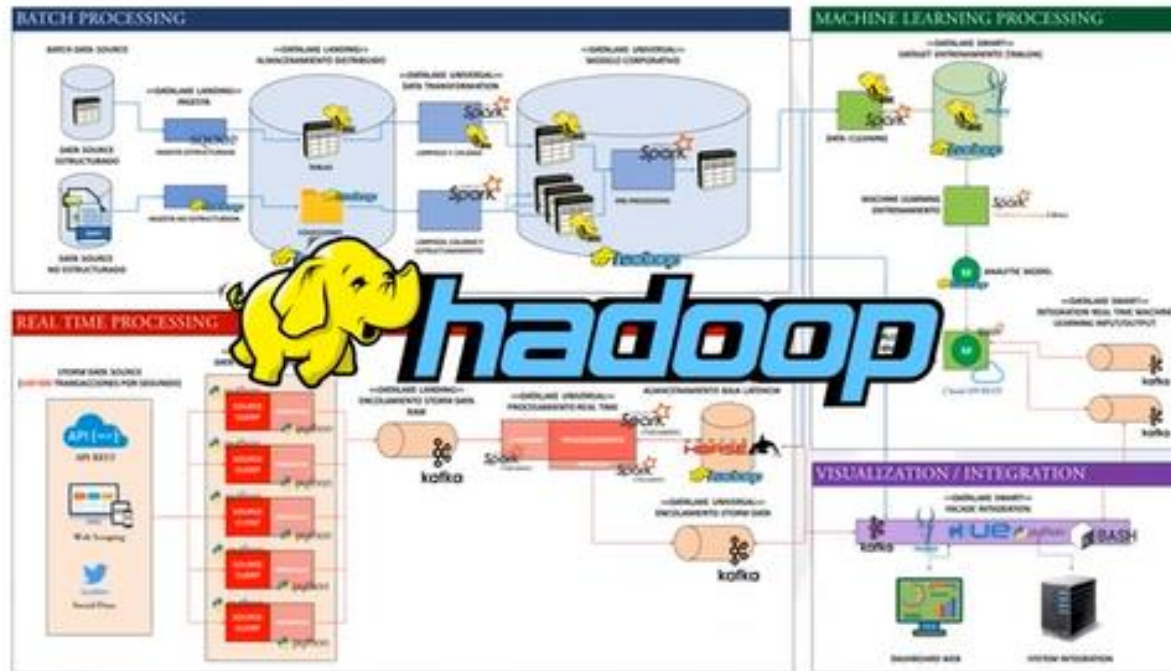


# Vista de Gobierno

## DATALAKE

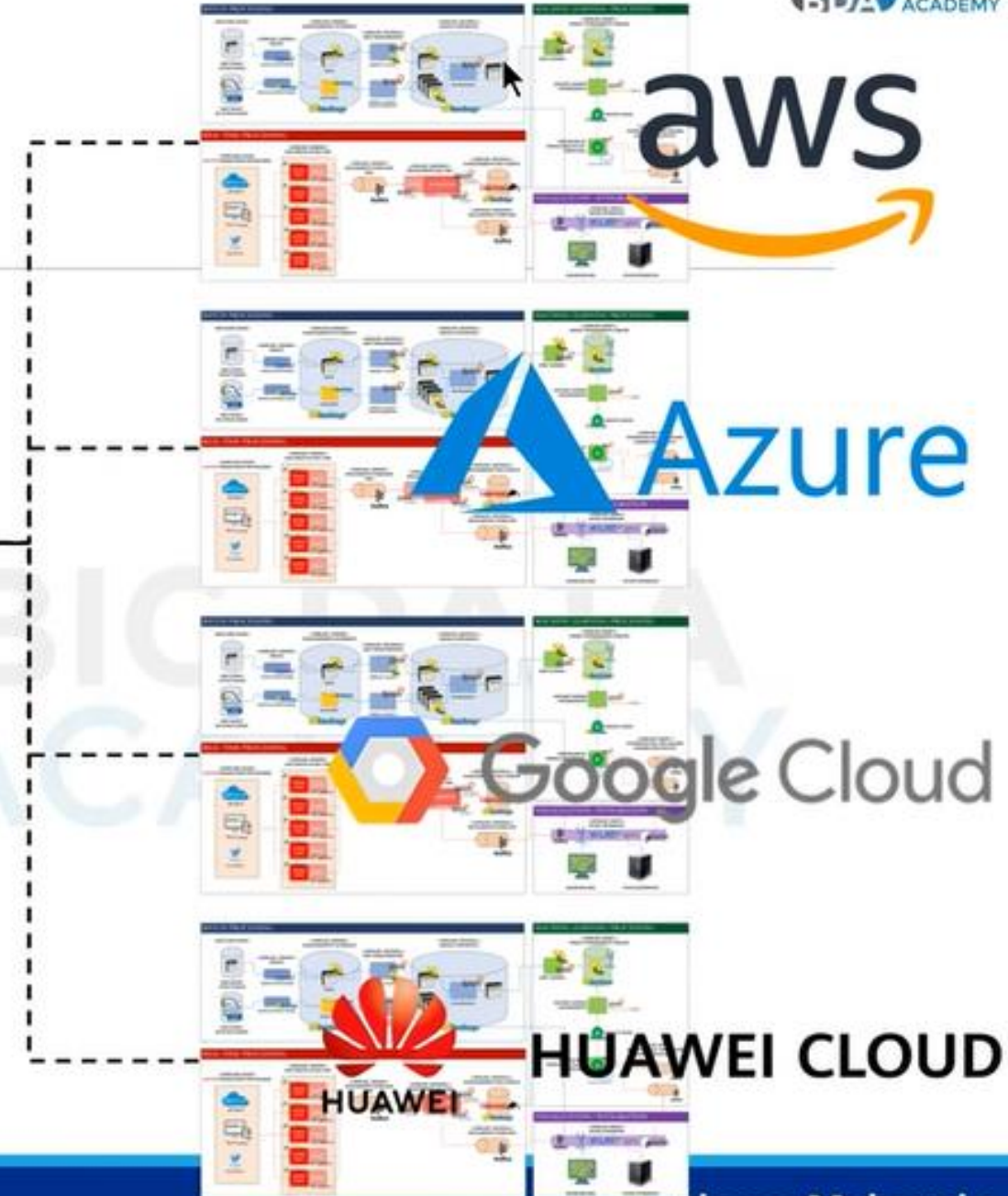


# Vista Tecnológica



Basado en el Ecosistema **Estándar Hadoop**

Los patrones de diseño son agnósticos al  
ecosistema tecnológico





# ¿Arquitecturas tecnológicas Cloud?

HAGA CLIC PARA AGREGAR TEXTO



# Servicios Cloud

## Google Cloud



## amazon



## Azure

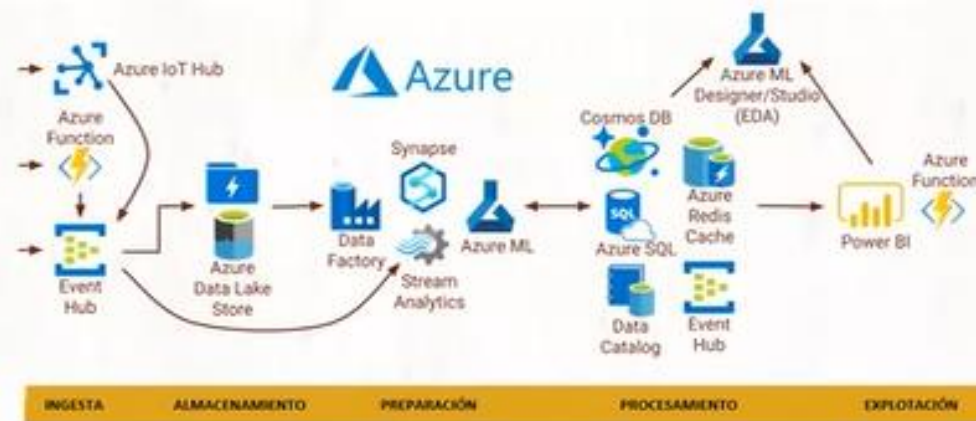


¡Existen muchos servicios de Cloud!

**Lo que nos importan** son los Servicios de Big Data on Cloud

# El flujo del Big Data on Cloud

Cada nube ofrece su propio mix tecnológico, pero **el flujo de procesamiento es estándar y agnóstico a la nube:**





# Herramientas estándar

Almacenamiento



Interfaz SQL



ETL



Soluciones



Workflow



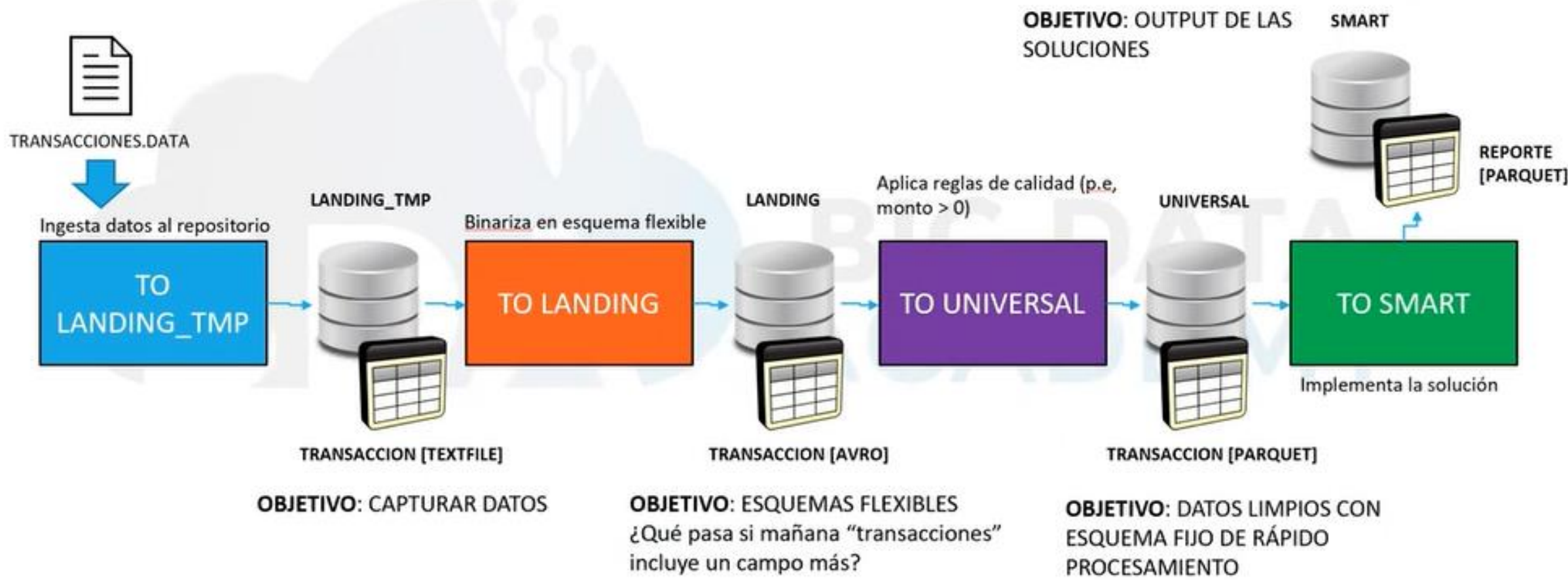


# Herramientas cloud equivalentes

	ECOSISTEMA ESTÁNDAR	aws	Google Cloud Platform	Azure
Almacenamiento	<b>hadoop</b>	<b>amazon S3</b>	Cloud Storage	Azure Blob Storage
Interfaz SQL	<b>HIVE</b>	Amazon Athena	Google Big Query	Azure Synapse Analytics
ETL	<b>HIVE</b>	AWS Glue	Data Fusion	DATAFLOW
Soluciones	<b>APACHE Spark</b>	<b>amazon EMR Spark</b>	<b>Cloud Dataproc</b>	<b>spark</b>
<u>Workflow</u>	<b>oozie</b>	<b>Airflow</b> aws	Google Cloud Composer	Azure Data Factory



# Flujo estándar sobre un Datalake



# Traducción de estándares a cada cloud de cualquier flujo que llegue hasta UNIVERSAL





# ¿Y qué pasa con las soluciones?



Transaccion en "Universal"

```
#Importamos el objeto para instanciar una sesion
from pyspark.sql import SparkSession

#Instanciamos una sesion
#Dependiendo de la potencia del cluster, tomara hasta 3 minutos en asignarle recursos
spark = SparkSession.builder.getOrCreate()

#Leemos los datos que procesaremos
dfTransaccion = spark.read.format("parquet").load("/objects/inputs/transaccion")

#Almacenamos el dataframe como vista temporal
dfTransaccion.createOrReplaceTempView("dfTransaccion")

#Procesamos el dataframe
dfReporte = spark.sql("""
    SELECT
        T.ID_PERSONA ID_PERSONA,
        COUNT(T.ID_PERSONA) CANTIDAD_TRANSACCIONES,
        SUM(T.MONTO) MONTO_TRANSACCIONES
    FROM
        dfTransaccion T
    GROUP BY
        T.ID_PERSONA
""")

#Guardamos el resultado en un directorio temporal del cluster
dfReporte.write.mode("overwrite").format("parquet").save("/objects/outputs/reporte")
```



Reporte resultante en "Smart"

Codificación agnóstica en SPARK sobre las tres nubes



amazon  
EMR



Cloud DataProc



# Pueden ser soluciones más complejas como Redes Neuronales



Transaccion en "Universal"



```
'''
#section Aplicamos modelos de Deep Learning [RED NEURONAL]
'''

#Importamos el algoritmo
from pyspark.ml.classification import MultilayerPerceptronClassifier

#Configuramos el algoritmo
algoritmo3 = MultilayerPerceptronClassifier(
    maxIter=1000,
    layers=[9, 50, 50, 2],
    blockSize=128,
    seed=1234,
    labelCol='Survived',
    featuresCol='features'
)

#Generamos el modelo con el algoritmo seleccionado
modelo3 = algoritmo3.fit(dfTrain)

#Validamos el modelo
dfPredicciones3 = modelo3.transform(dfTest)
dfPredicciones3.show()
```

Codificación agnóstica en SPARK sobre las tres nubes



Red Neuronal

¡La implementación de nuestras soluciones sí puede ser agnóstica!

# Vista de patrones de diseño [Proceso BATCH]

HAGA CLIC PARA AGREGAR TEXTO

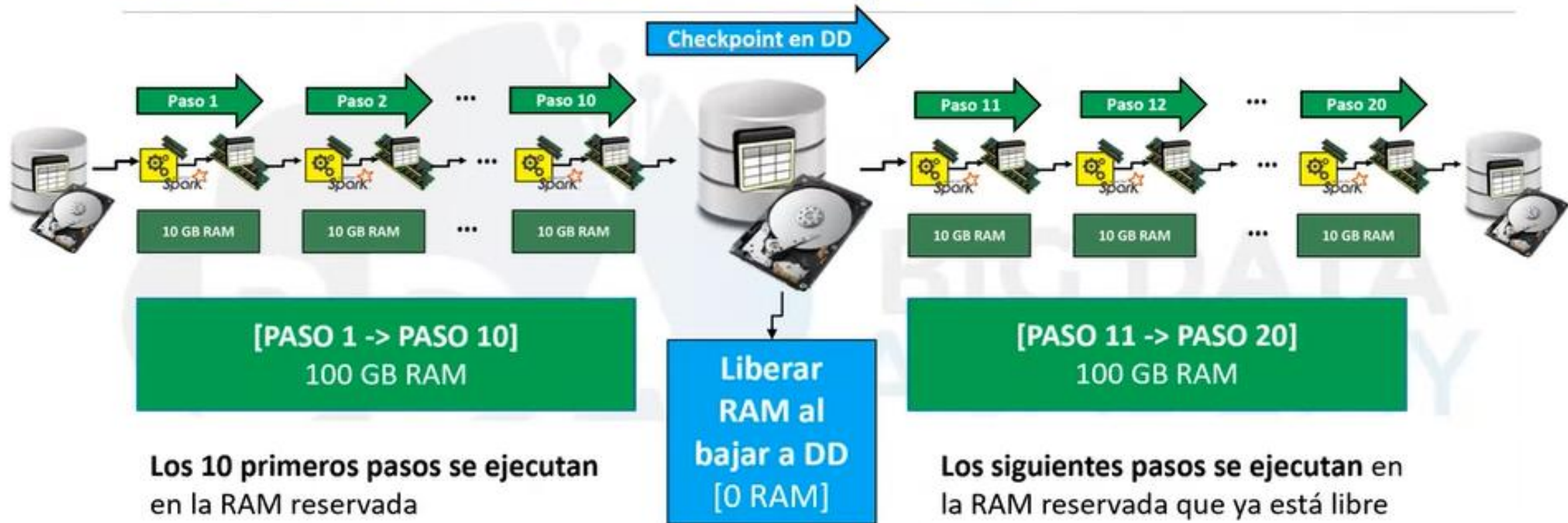


# Problema: Mi script tiene muchos pasos y no tengo suficiente RAM en el clúster



Si nuestro proceso requiere de muchos pasos y ya no tenemos RAM, **¿cómo podemos programarlo sin que colapse?**

# Solución: Patrón checkpoint

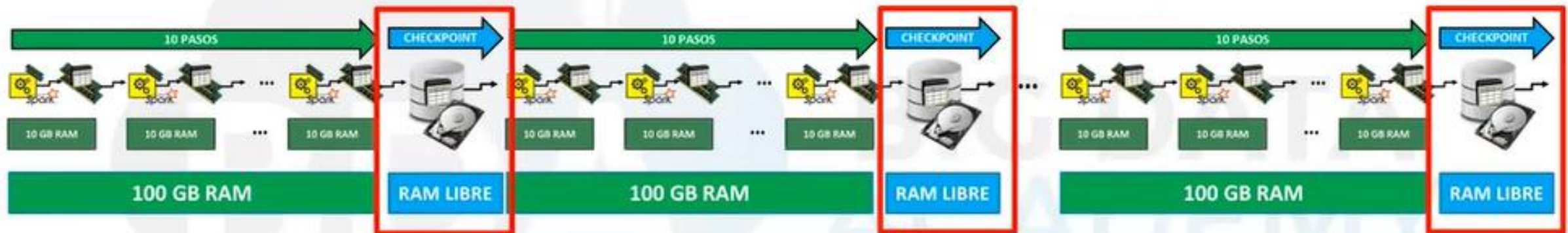


**Almacenamos en disco duro el dataframe del paso 10 para liberar toda la ram** de los dataframes del paso 1 al paso 9 que lo generaron, leemos nuevamente el dataframe y continuamos procesando



# ¿Y si aumentamos más pasos en el proceso?

Cuando la memoria RAM esté a punto de llenarse, crearemos un checkpoint, de esa manera nuestro proceso podrá tener todos los pasos necesarios



Gracias al “checkpoint”, evitamos que la memoria colapse.

¿Y tiene alguna **desventaja**? Sí, ya que agrega un paso de escritura al disco duro de la resultante de la cadena de procesos

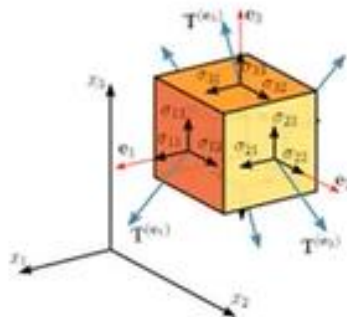


# Vista de patrones de diseño [Proceso MACHINE LEARNING]

HAGA CLIC PARA AGREGAR TEXTO

# ¿Qué es un tensor?

## Objetos hiper-dimensionales



## Cálculo tensorial

Principles of Tensor Calculus

$$g_{ij} = \mathbf{E}_i \cdot \mathbf{E}_j$$

$$\mathbf{E}_i = \frac{\partial \mathbf{r}}{\partial x^i}$$

$$g_{ij}^{1-2} = \begin{bmatrix} g_{11} & g_{12} \\ g_{21} & g_{22} \end{bmatrix}$$

$$g_{ij}^{1-2} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$g_{ij}^{1-2} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

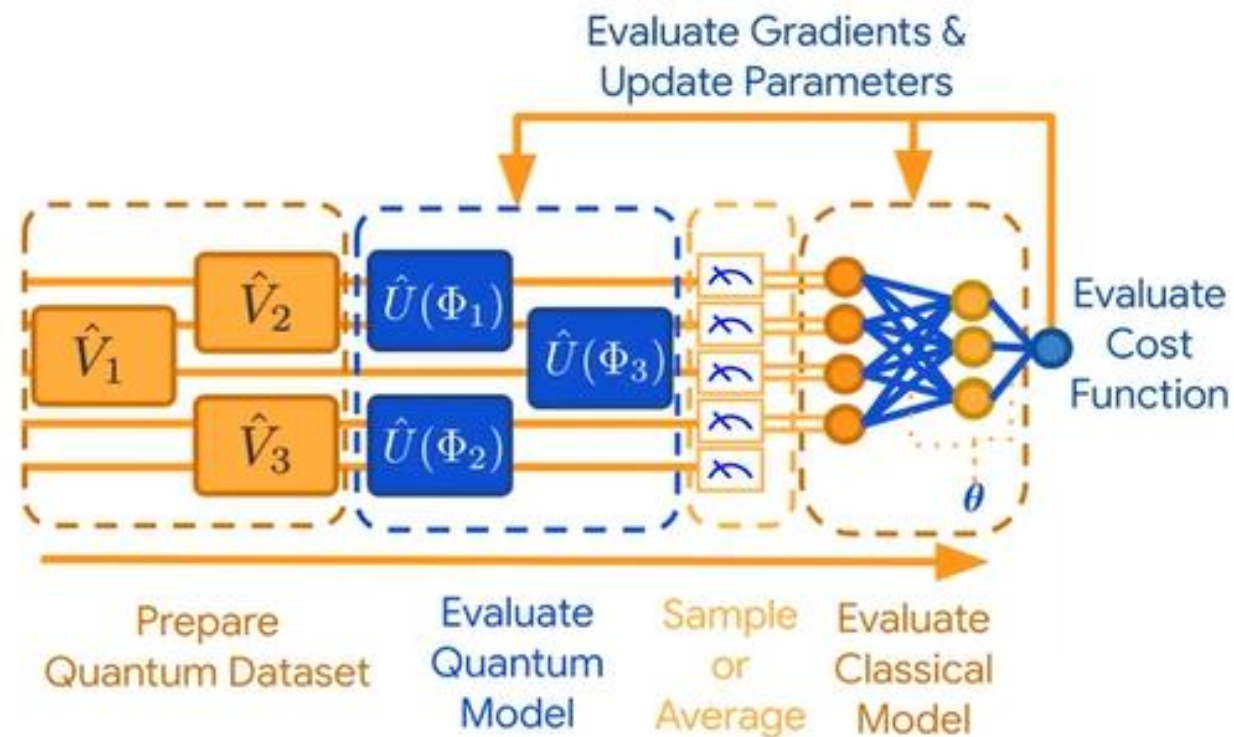
$$\mathbf{E}^i = \nabla u^i$$

$$g^{ij} = \mathbf{E}^i \cdot \mathbf{E}^j$$

$$\frac{\delta A^i}{\delta t} = \frac{dA^i}{dt} + \Gamma_{ij}^k A^j \frac{dx^i}{dt}$$

Taha Sochi

## Aplicaciones para computación cuántica

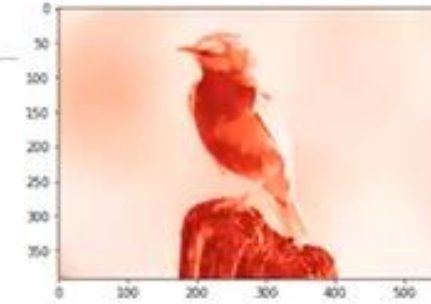


Pongamos un ejemplo simple para entender el concepto

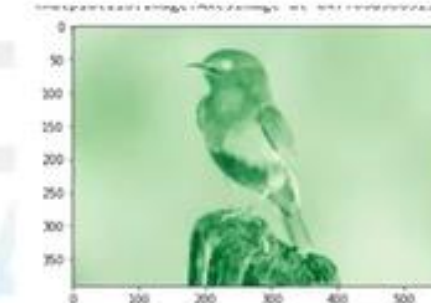
# ¿Qué es un tensor?



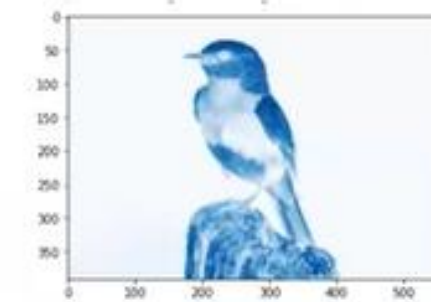
Imagen original



Canal R



Canal G

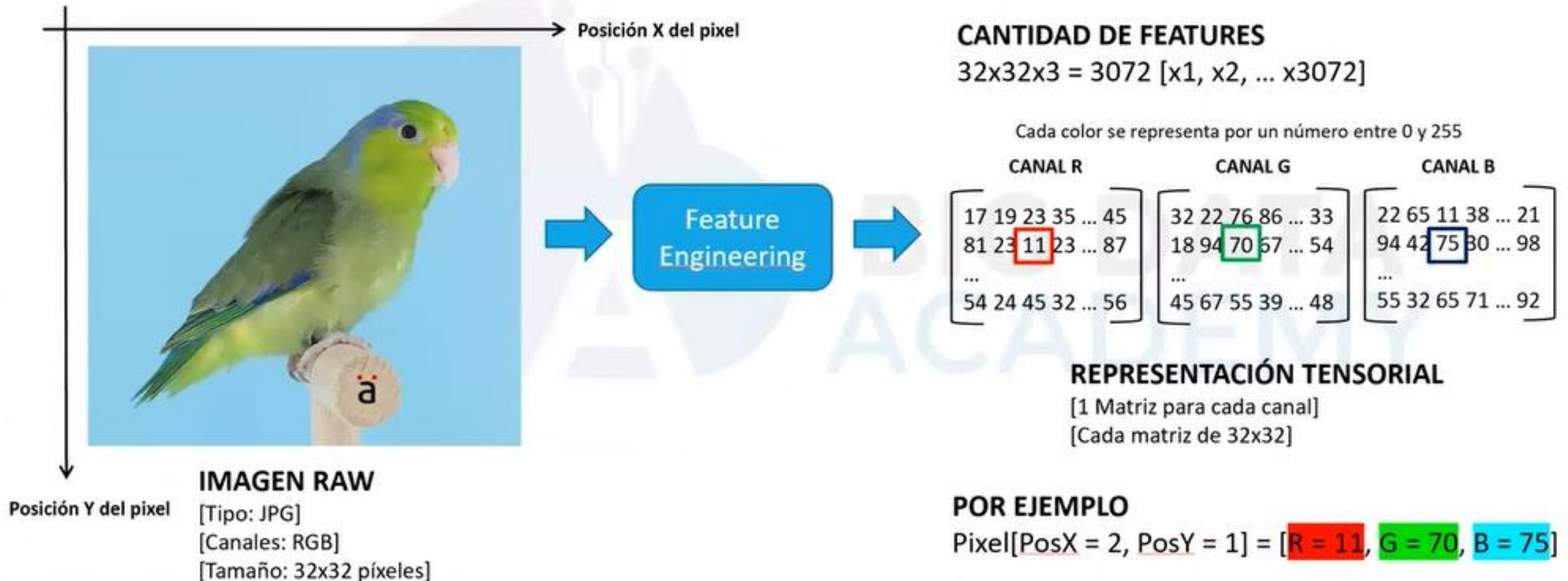


Canal B

Representación tensorial  
[alto, ancho, canal]



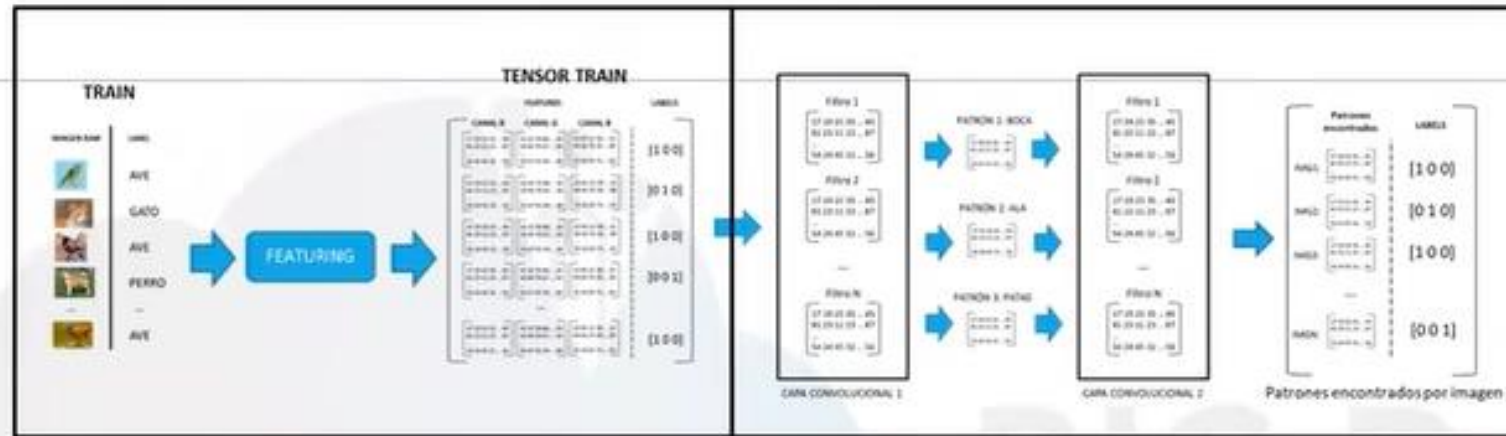
# Representación tensorial de una imagen





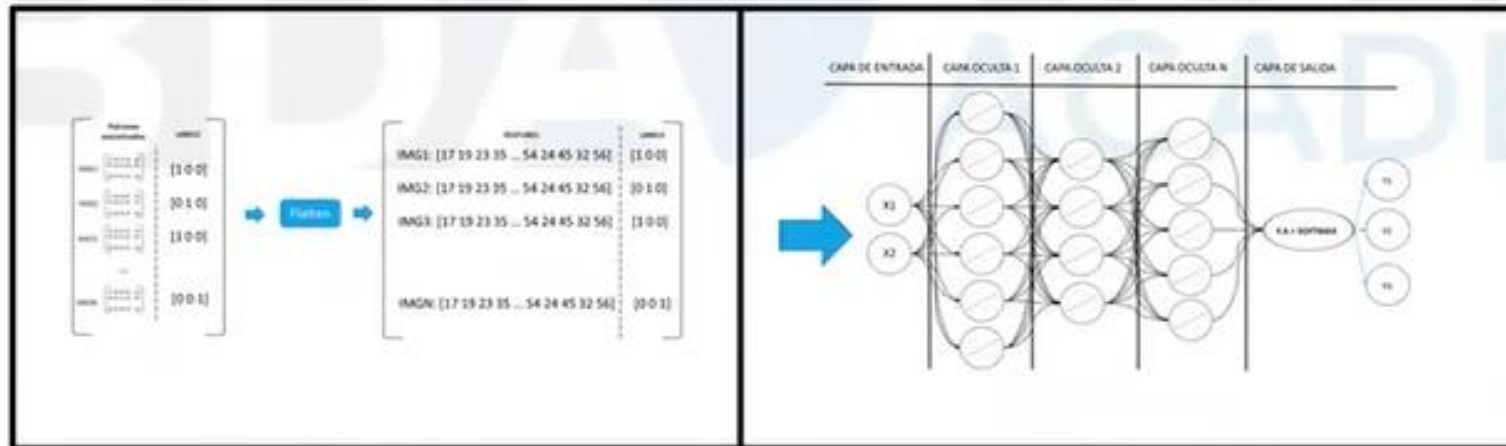
## Representación tensorial del Dataset

# Arquitectura para Visión Artificial



Feature Engineering

Red Convolutional

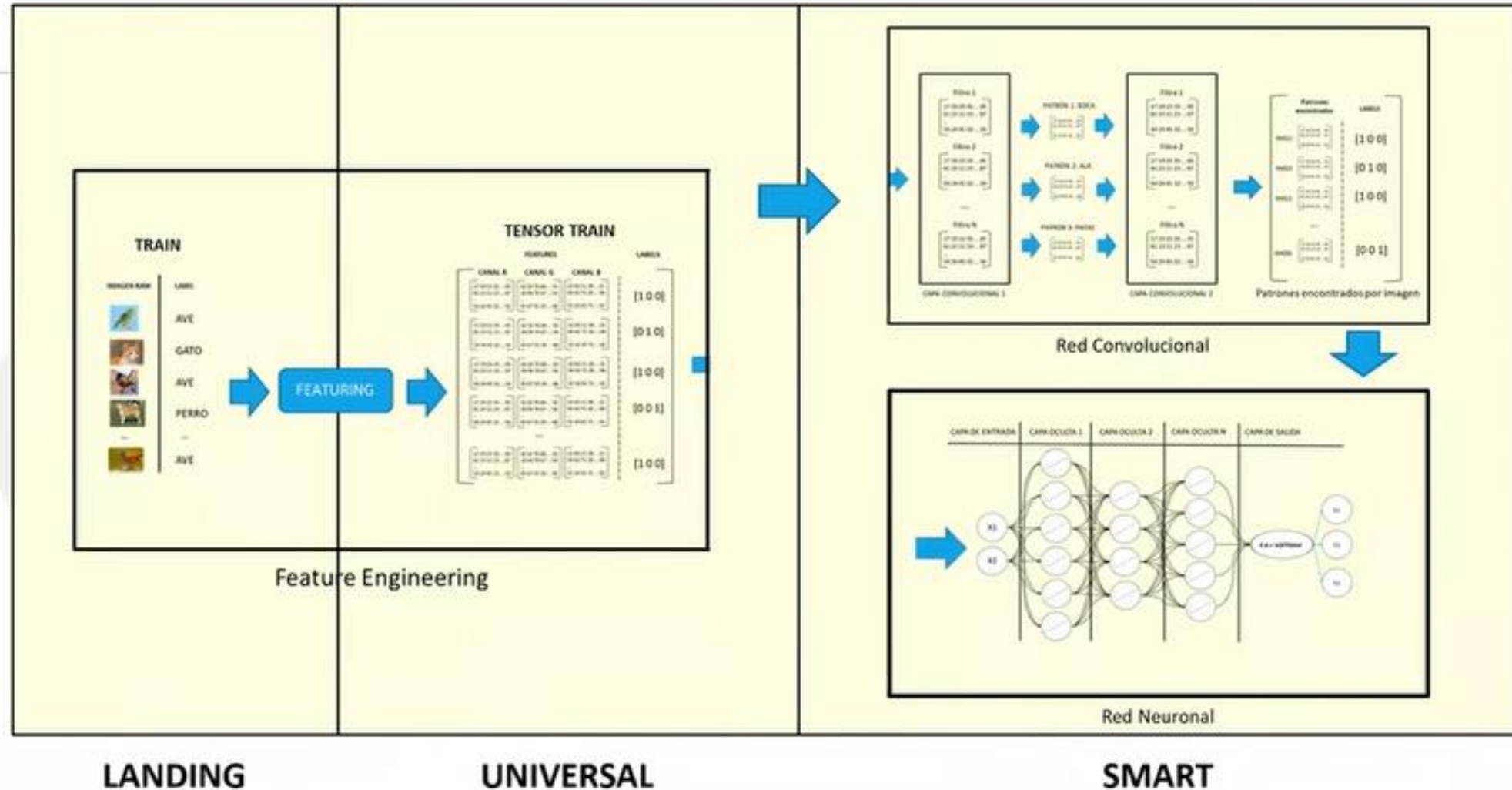


Flatten

Red Neuronal



# Machine Learning on Datalake



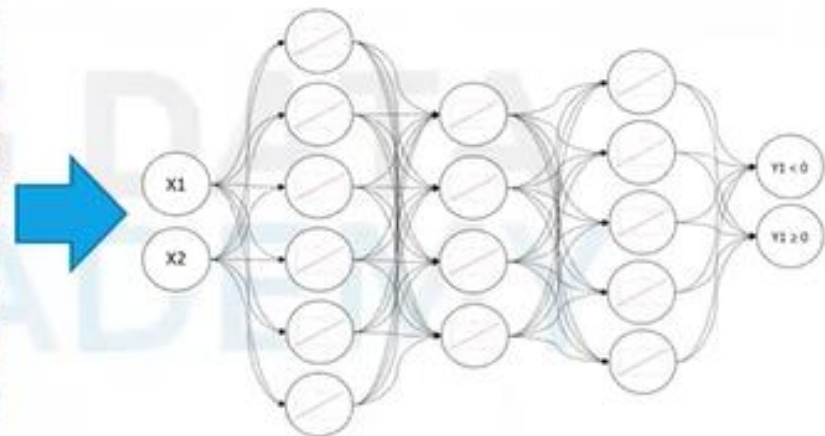
# Uso de GPU para procesos matriciales

Ya que la data la tenemos en forma matricial y usamos el álgebra lineal, podemos paralelizar aún más el procesamiento usando GPU

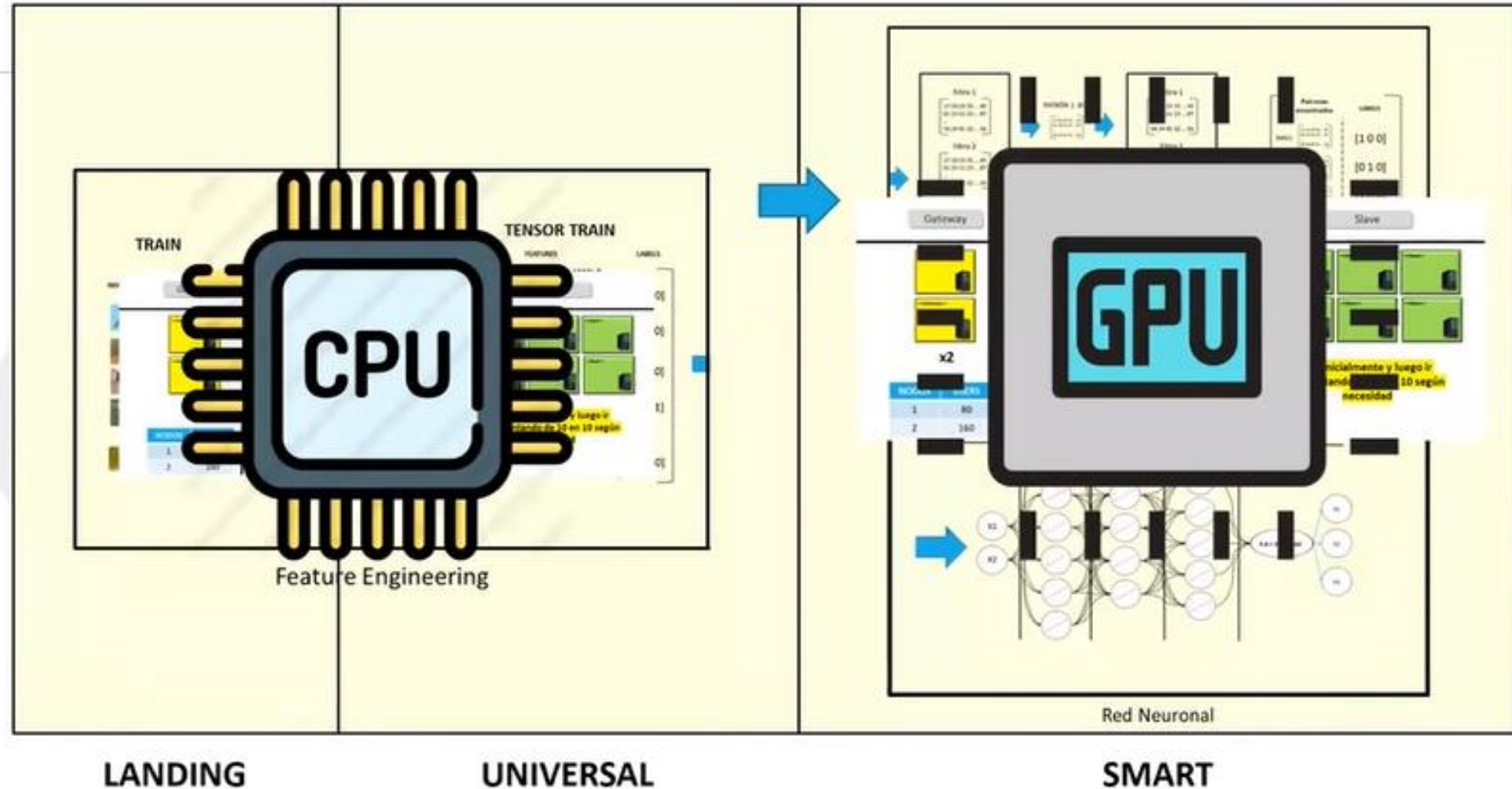
$$X = \begin{matrix} & X_1 & X_2 & X_3 & X_4 & X_5 \\ \begin{matrix} P_1 \\ P_2 \\ P_3 \\ P_4 \\ P_5 \\ P_6 \\ P_7 \end{matrix} & \begin{bmatrix} 0.3 & 0.1 & 0.2 & 0.5 & 0.1 \\ 0.7 & 0.9 & 0.2 & 0.7 & 0.8 \\ 0.2 & 0.6 & 0.2 & 0.8 & 0.6 \\ 0.9 & 0.3 & 0.9 & 0.9 & 0.1 \\ 0.1 & 0.4 & 0.4 & 0.4 & 0.9 \\ 1.0 & 0.8 & 0.8 & 0.2 & 0.7 \\ 0.9 & 0.7 & 0.7 & 0.2 & 0.6 \end{bmatrix} \end{matrix}$$
$$Y = \begin{matrix} & Y_1 \\ \begin{matrix} P_1 \\ P_2 \\ P_3 \\ P_4 \\ P_5 \\ P_6 \\ P_7 \end{matrix} & \begin{bmatrix} 0.0 \\ 1.0 \\ 1.0 \\ 1.0 \\ 0.0 \\ 1.0 \\ 0.0 \end{bmatrix} \end{matrix}$$



¡La GPU es especialista en realizar cálculos de productos matriciales!, **factor de aumento en velocidad de entre x2 a x10**



# Infraestructura para Machine Learning on Datalake

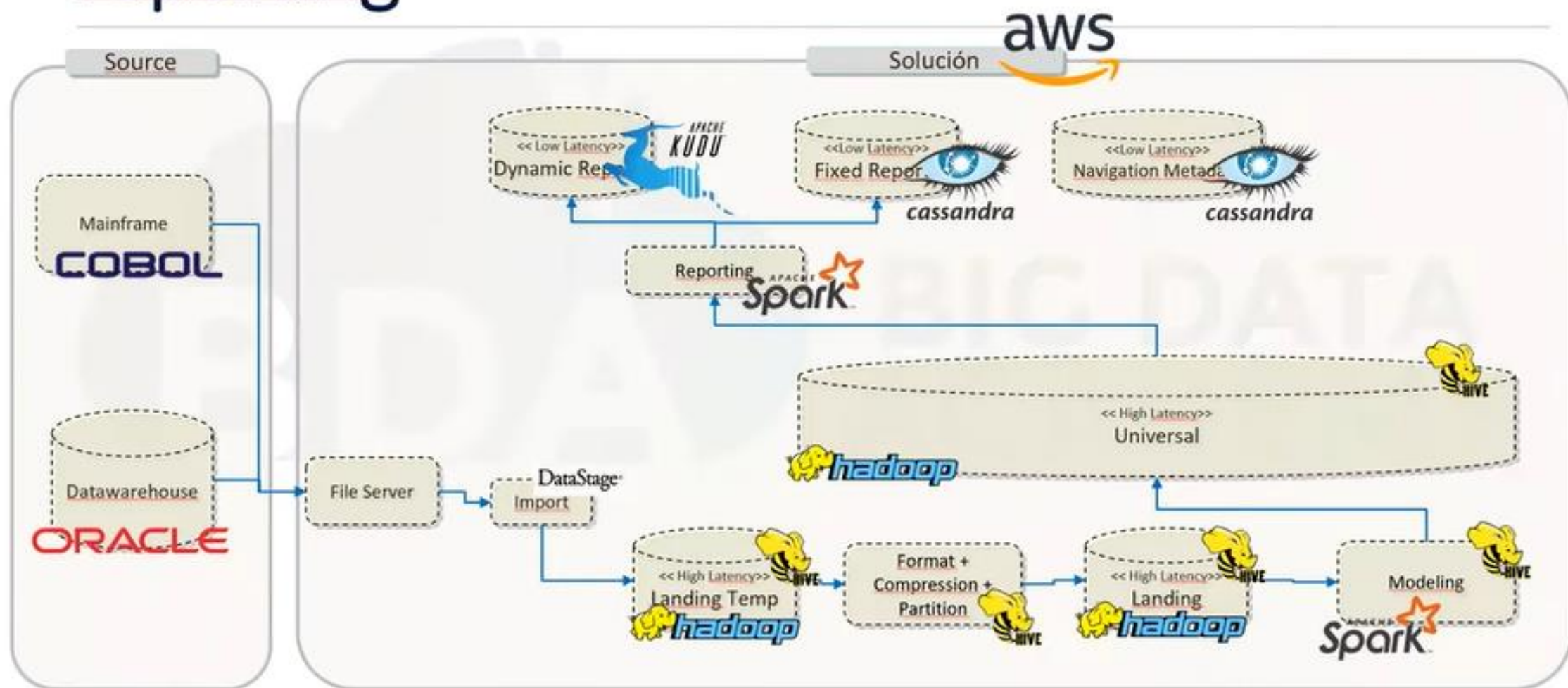




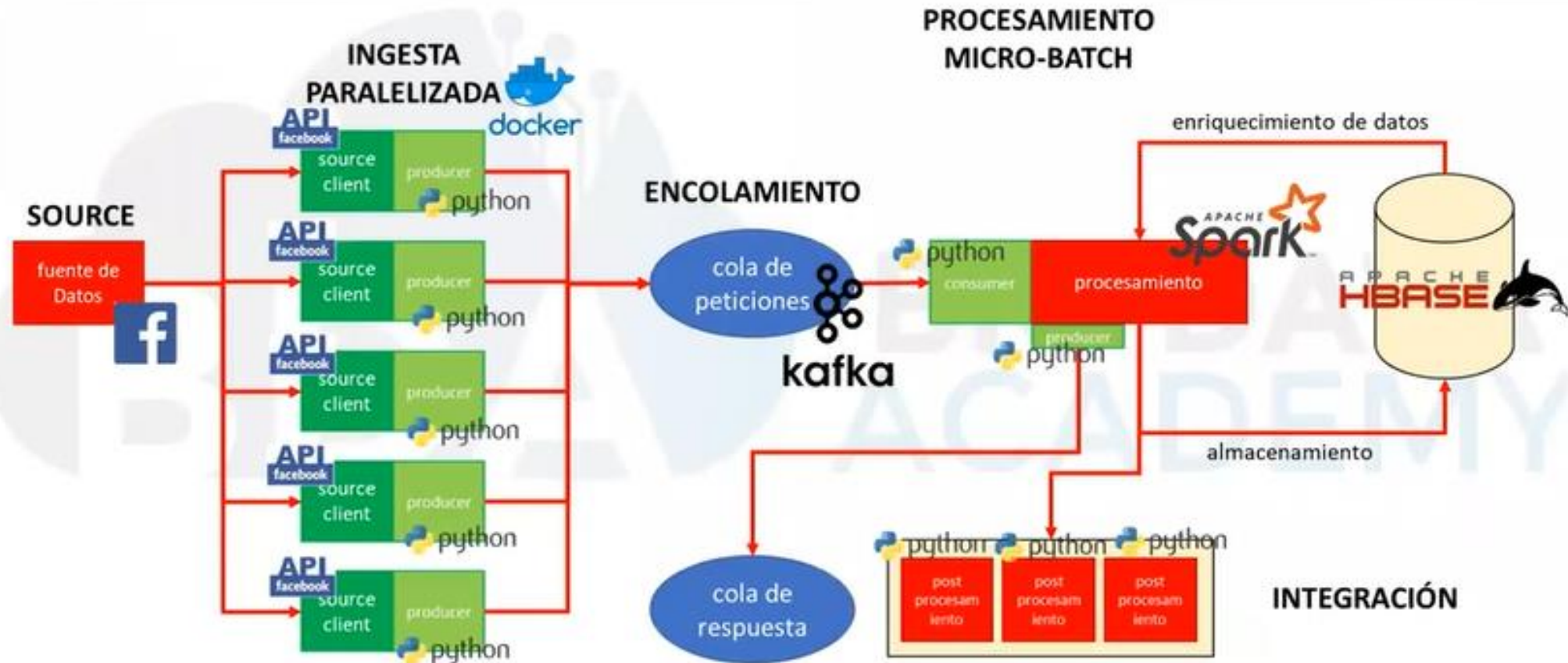
# Vista de patrones de diseño [ETL, Reporting, Real-Time]

HAGA CLIC PARA AGREGAR TEXTO

# Patrón de procesamiento batch aplicado a reporting

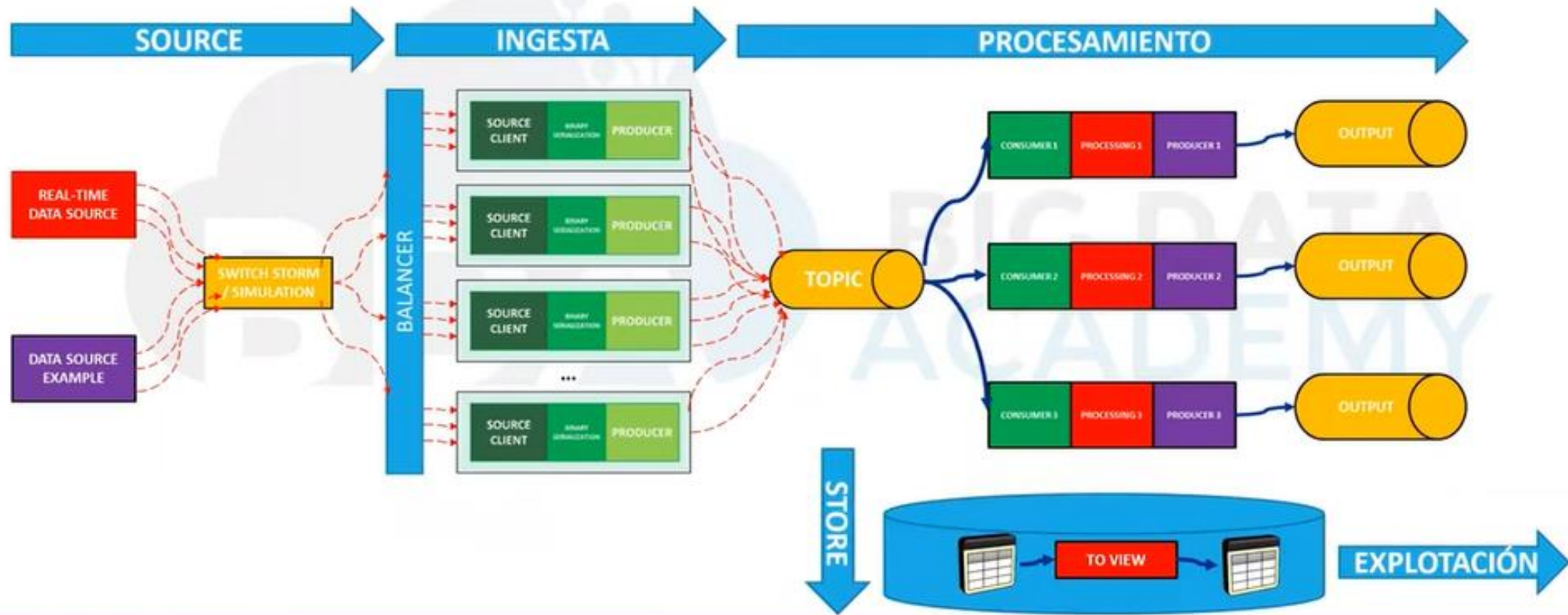


# Patrón de procesamiento en real time aplicado a redes sociales

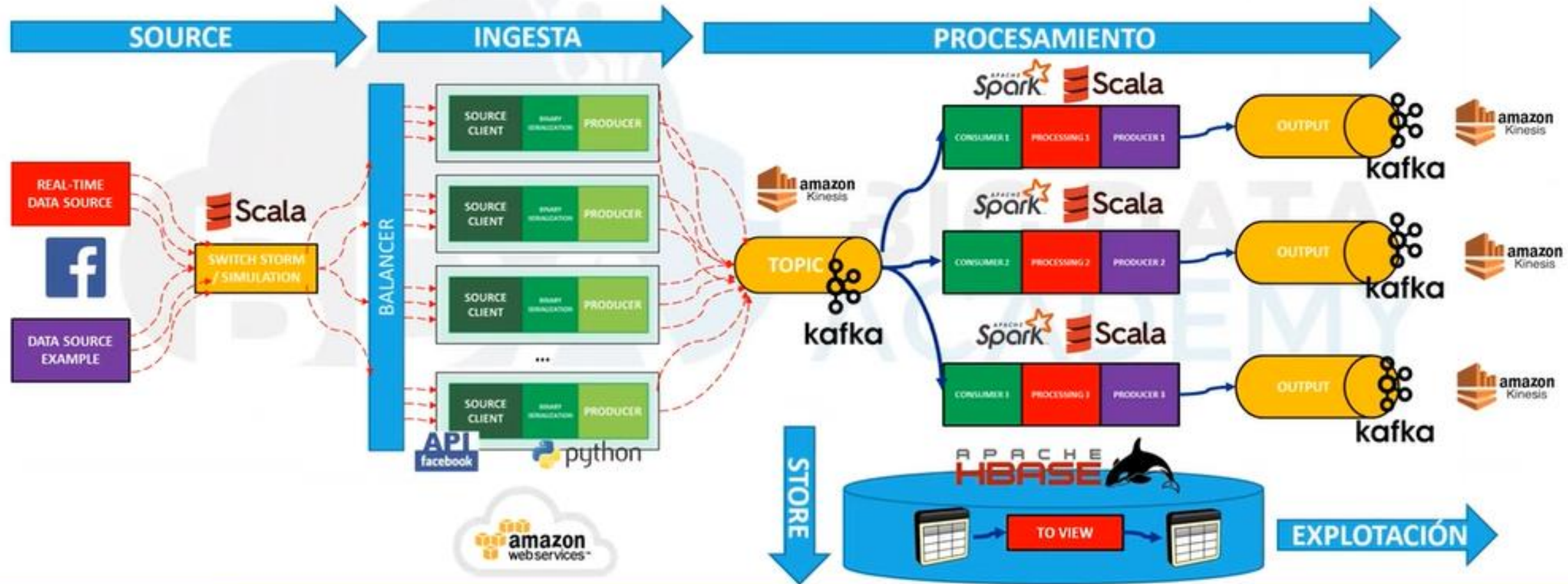




# Arquitectura estándar extendida para procesamiento en Real-Time

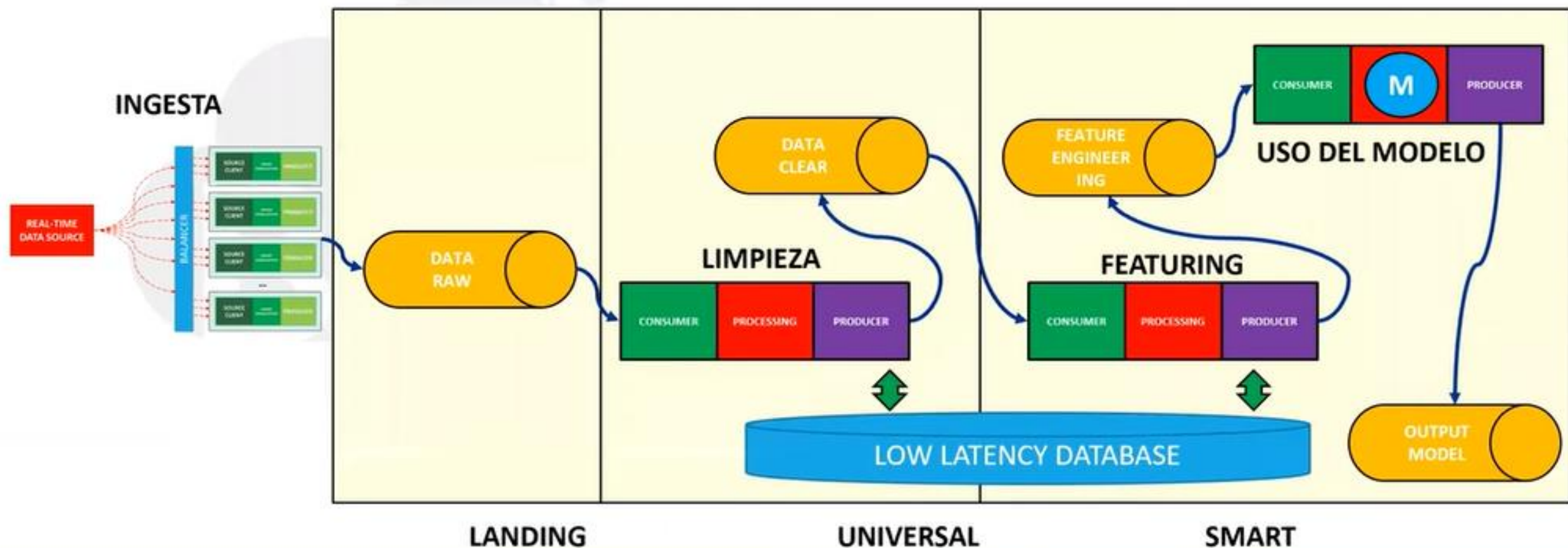


# Arquitectura estándar extendida para procesamiento en Real-Time





# Machine Learning on Datalake para explotación en Real-Time



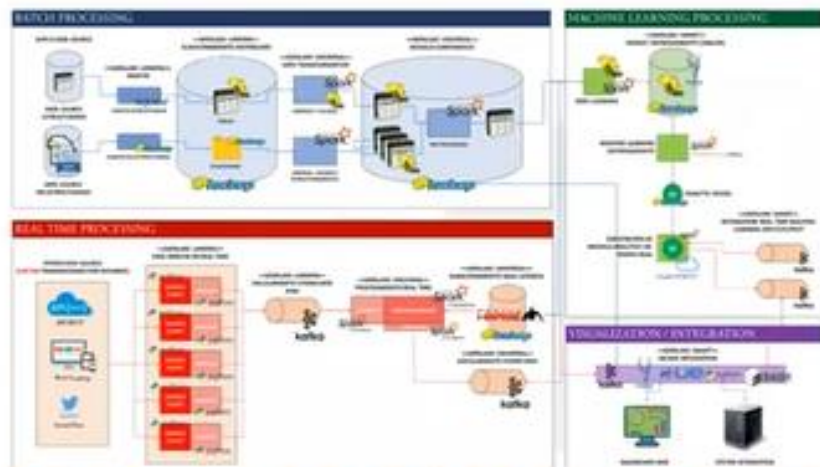


# Conclusiones

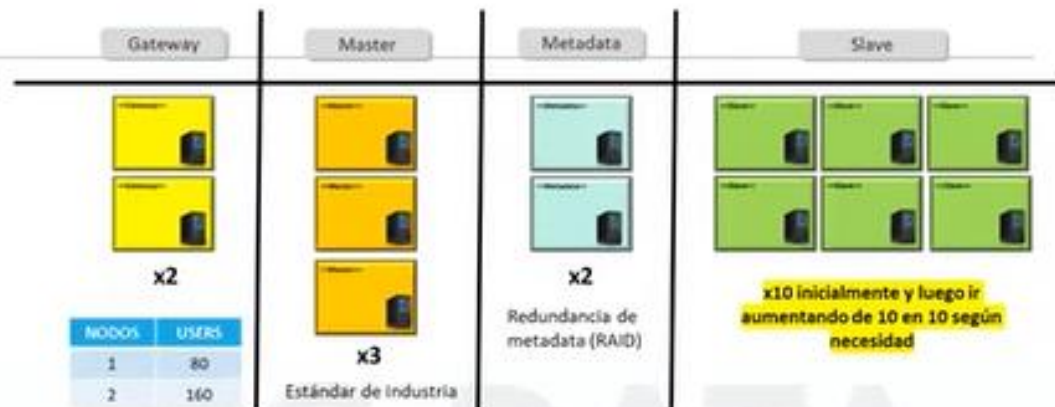
HAGA CLIC PARA AGREGAR TEXTO

# Próximos pasos

## Arquitecturas y patrones

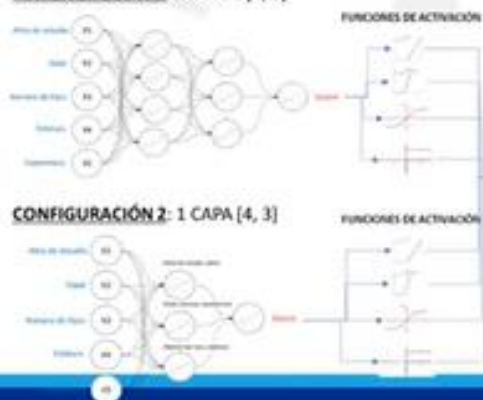


## Infraestructura Cloud y On-Premise

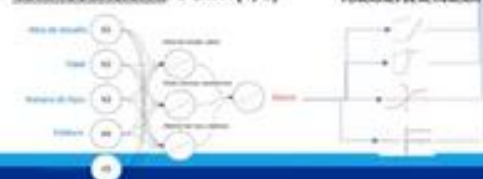


## Machine Learning Escalable

CONFIGURACIÓN 1: 2 CAPAS [4, 3]



CONFIGURACIÓN 2: 1 CAPA [4, 3]



MODELOS ENTRENADOS



## Arquetipos de código

```

def formatData(register):
    registerFormatted = []
    try:
        # Cada registro recibido es procesado por una regla de clasificación
        # Esta regla es un array de dos registros, el campo "id" es la clave y el campo "v"
        # es el valor, la "clave" siempre es None, solo nos interesa el registro, así que
        registerValue = register[1]
        # Convertimos el registro a un array de JSON
        registerFormatted = json.loads(registerValue)
    except Exception as e:
        # Si el string no contiene un formato JSON, devolvemos un array vacío
        registerFormatted = []
    return registerFormatted

def enrichMetadata(register):
    # Se agregan los campos de fecha y hora
    dateTime = datetime.datetime.now()
    register["FECHA_TRANSACCION"] = str(dateTime.date())
    register["HORA_TRANSACCION"] = str(dateTime.time())
    
```