

Formación senior en Spark & Scala

PARTE DE LA FORMACIÓN DE BIG DATA ACADEMY

Sobre el profesor



Soy un profesional con 12 años de experiencia empresarial en el mundo de las TI y 8 años de experiencia en Big Data, Cloud & Machine Learning, en la parte académica soy investigador en temas de computación cuántica y asesor de tesis de Big Data & Machine Learning

ACTUALMENTE SOY

- Director de Big Data & Analytics en **Big Cloud Dev**
- Formador de cursos de Big Data en **Big Data Academy**



CONTACTO: alonsoraulmgs@gmail.com

LINKEDIN: <https://www.linkedin.com/in/alonsoraulmg/>

Concepto

BIG DATA
ACADEMY

Spark

Es un motor de procesamiento distribuido paralelo in-memory. Proporciona apis en Java, Scala, Python y R. Spark mantiene la escalabilidad lineal y la tolerancia a fallos de MapReduce, pero amplía sus bondades gracias a varias funcionalidades.



Objetivo fundamental

Objetivo fundamental de Spark

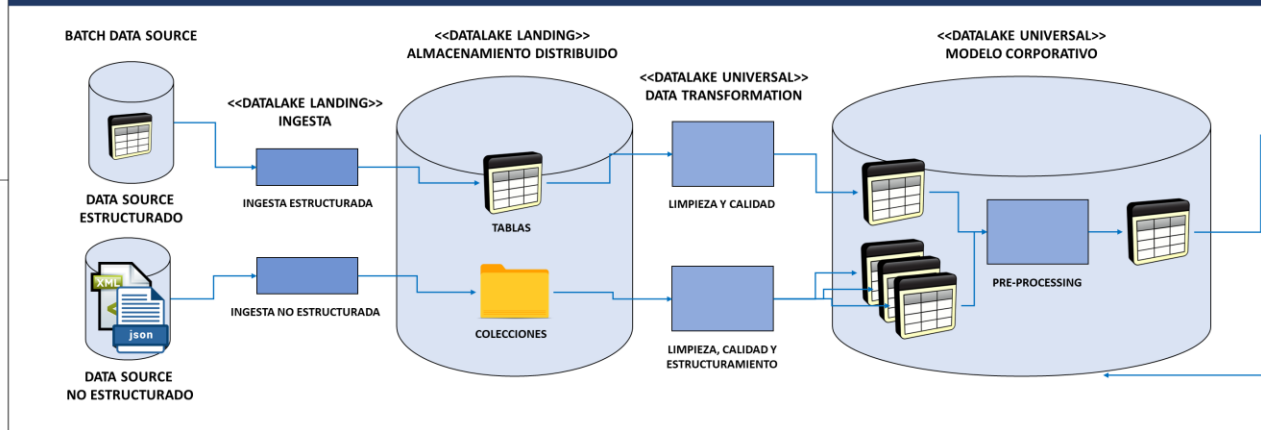
**Ejecutar procesos
lo más rápido
posible**

Naturaleza de funcionamiento

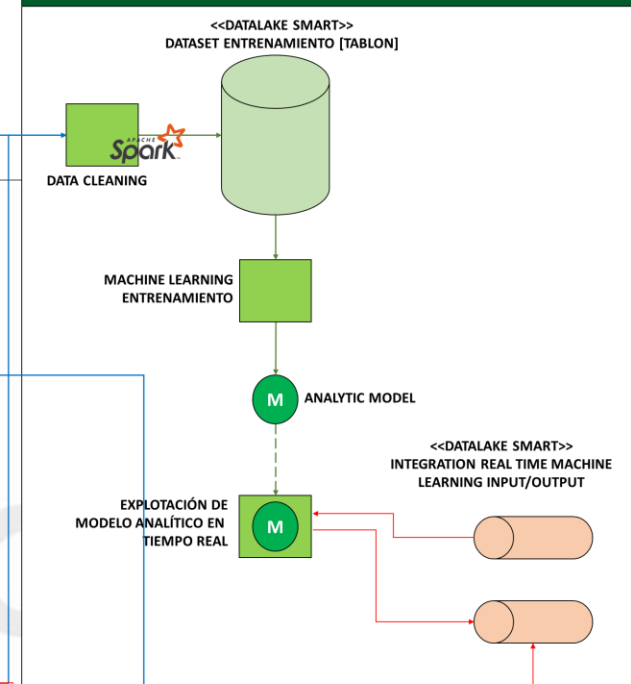


Vista Conceptual con Patrones de Diseño de Big Data

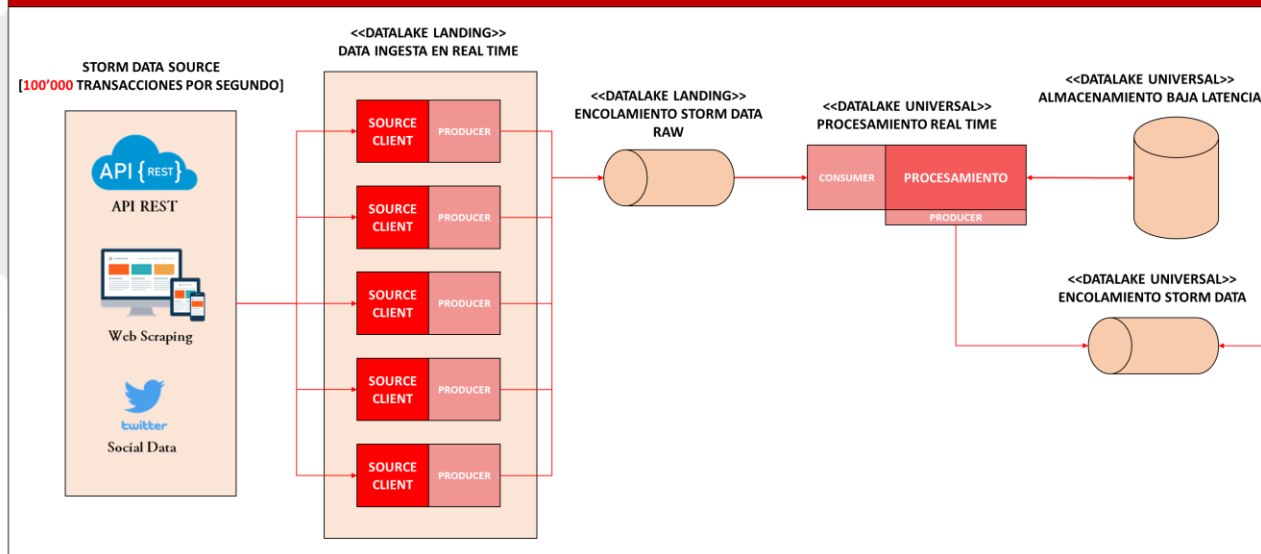
ETL ENRICHMENT



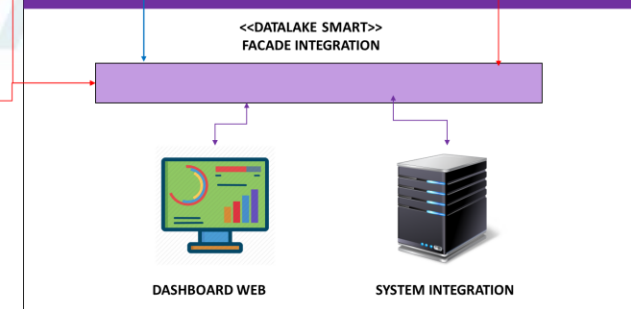
TEXT ANALYTICS



STORM DATA CAPTURE

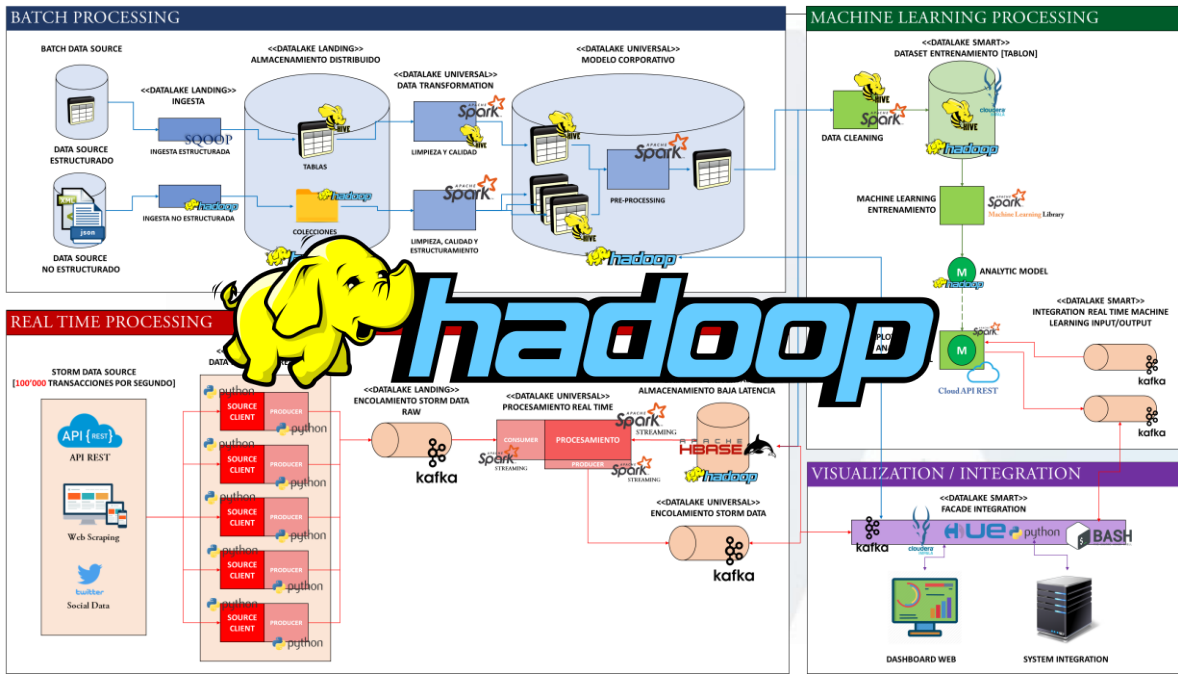


VISUALIZATION



¡Este diagrama sí lo podemos codificar!

El ecosistema estándar basado en Hadoop



Basado en el Ecosistema **Estándar Hadoop**

Los patrones de diseño son agnósticos al
ecosistema tecnológico

This section displays the Hadoop ecosystem architecture mapped onto four major cloud providers, demonstrating its agnostic design:

- AWS:** The architecture is mapped onto Amazon Web Services, showing how Hadoop components like HDFS, MapReduce, and Spark can be implemented using AWS services like S3, EMR, and Lambda.
- Azure:** The architecture is mapped onto Microsoft Azure, showing how Hadoop components can be implemented using Azure services like Data Lake Storage, HDInsight, and Databricks.
- Google Cloud:** The architecture is mapped onto Google Cloud Platform, showing how Hadoop components can be implemented using Google Cloud services like Big Data Storage, Dataproc, and Dataflow.
- HUAWEI CLOUD:** The architecture is mapped onto Huawei Cloud, showing how Hadoop components can be implemented using Huawei Cloud services like OBS, EMR, and Spark.

Herramientas cloud equivalentes

	ECOSISTEMA ESTÁNDAR	aws	Google Cloud Platform	Azure
Almacenamiento				
Interfaz SQL				
ETL				
Soluciones				
Workflow				

Spark vs Hadoop

Apache Spark™ is a unified analytics engine for large-scale data processing.

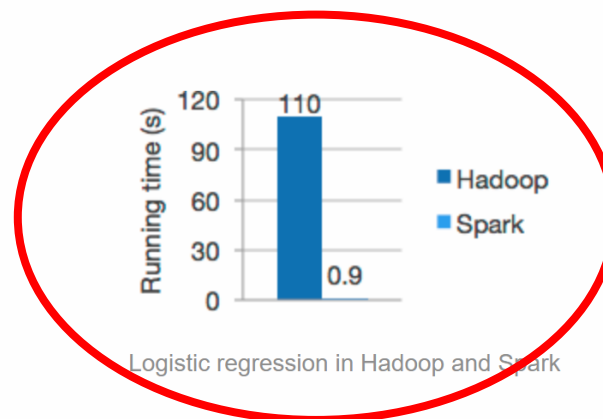
Speed

Run workloads 100x faster.

Apache Spark achieves high performance for both batch and streaming data, using a state-of-the-art DAG scheduler, a query optimizer, and a physical execution engine.

Ease of Use

Write applications quickly in Java, Scala, Python, R,



Logistic regression in Hadoop and Spark

```
df = spark.read.json("logs.json")
df.where("age > 21")
  .select("name.first").show()
```

Latest News

Spark 2.3.3 released
Spark 2.2.3 released
Spark+AI Summit 2019, San Francisco (Dec 19, 2018)
Spark 2.4.0 released



Download

Built-in Libraries

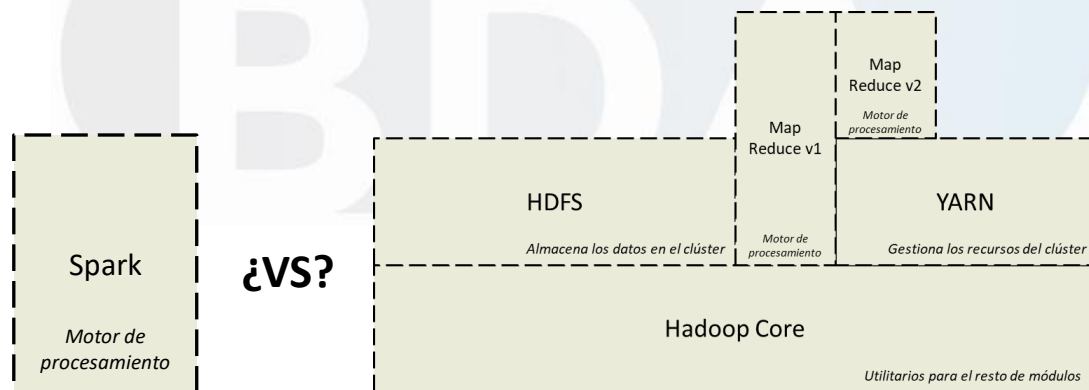
SQL and DataFrames

Spark puede llegar a ser hasta **100 veces más rápido que Hadoop**

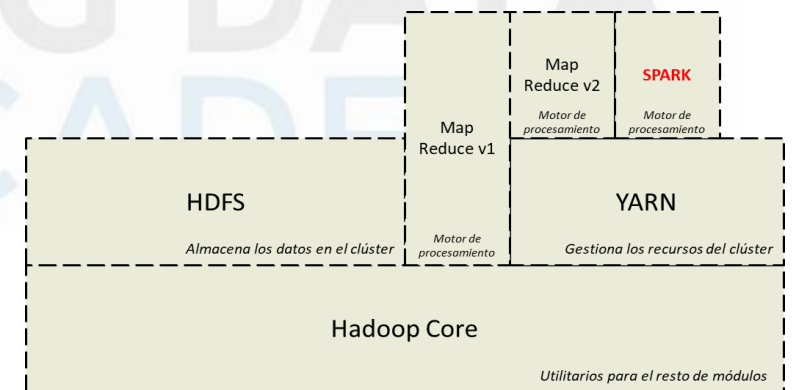
Si Spark es “mejor” entonces, ¿reemplaza a Hadoop?

No, Spark es un motor de procesamiento, Hadoop es un ecosistema que incluye un motor de procesamiento llamado MapReduce v2 y otros componentes (HDFS para almacenar en disco duro, YARN para gestionar recursos). **Spark no reemplaza a Hadoop, sino que lo potencia.**

¿SPARK VS HADOOP?: **¡NO!**



SPARK ON HADOOP

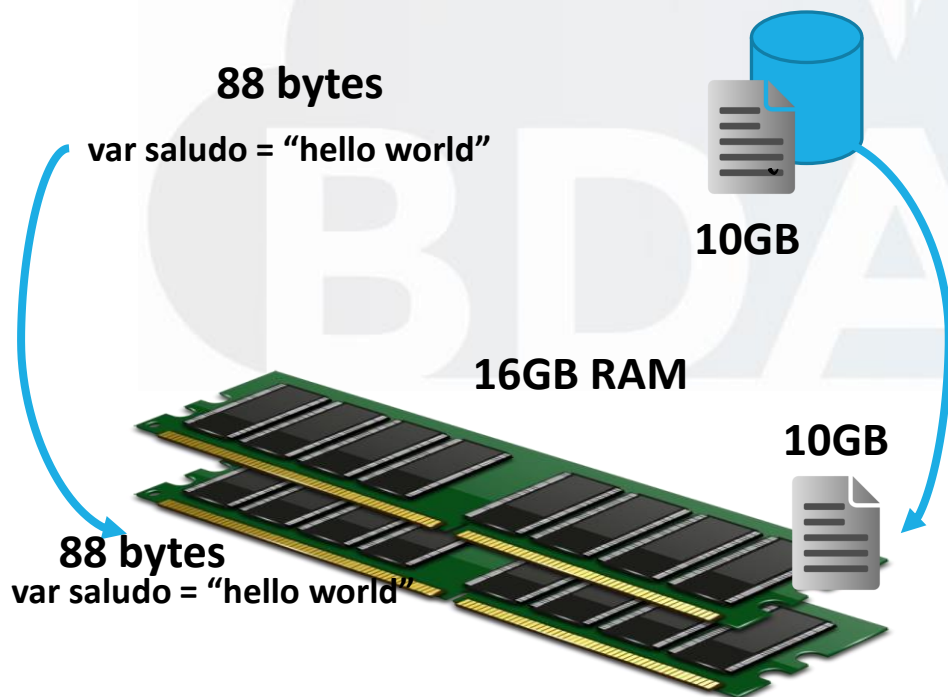


Variables in-memory

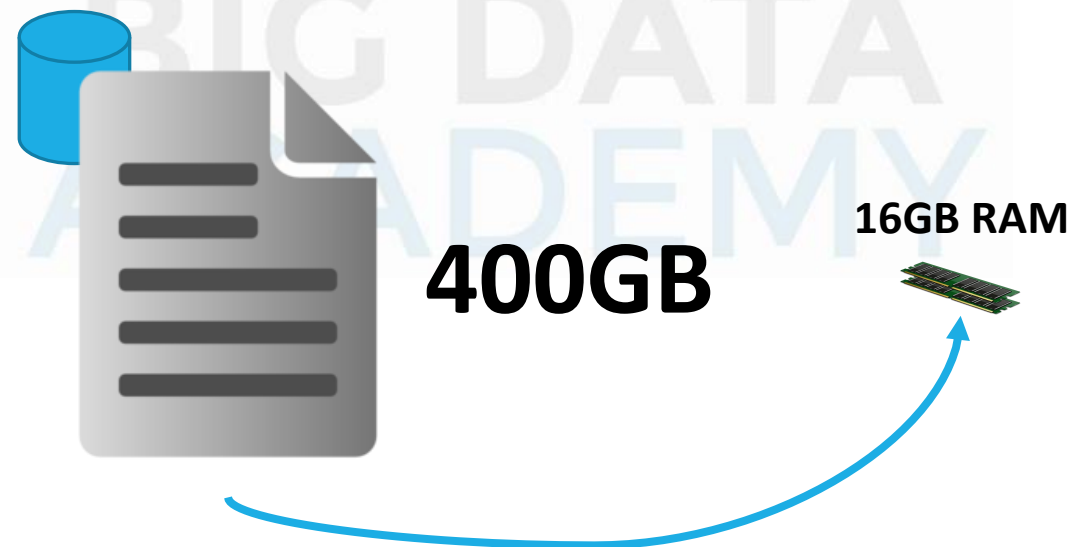
A large, faint, light blue watermark is centered in the background of the slide. It consists of a cloud shape on the left and the text "BDA BIG DATA ACADEMY" on the right, matching the logo in the top right corner.

Variables en memoria

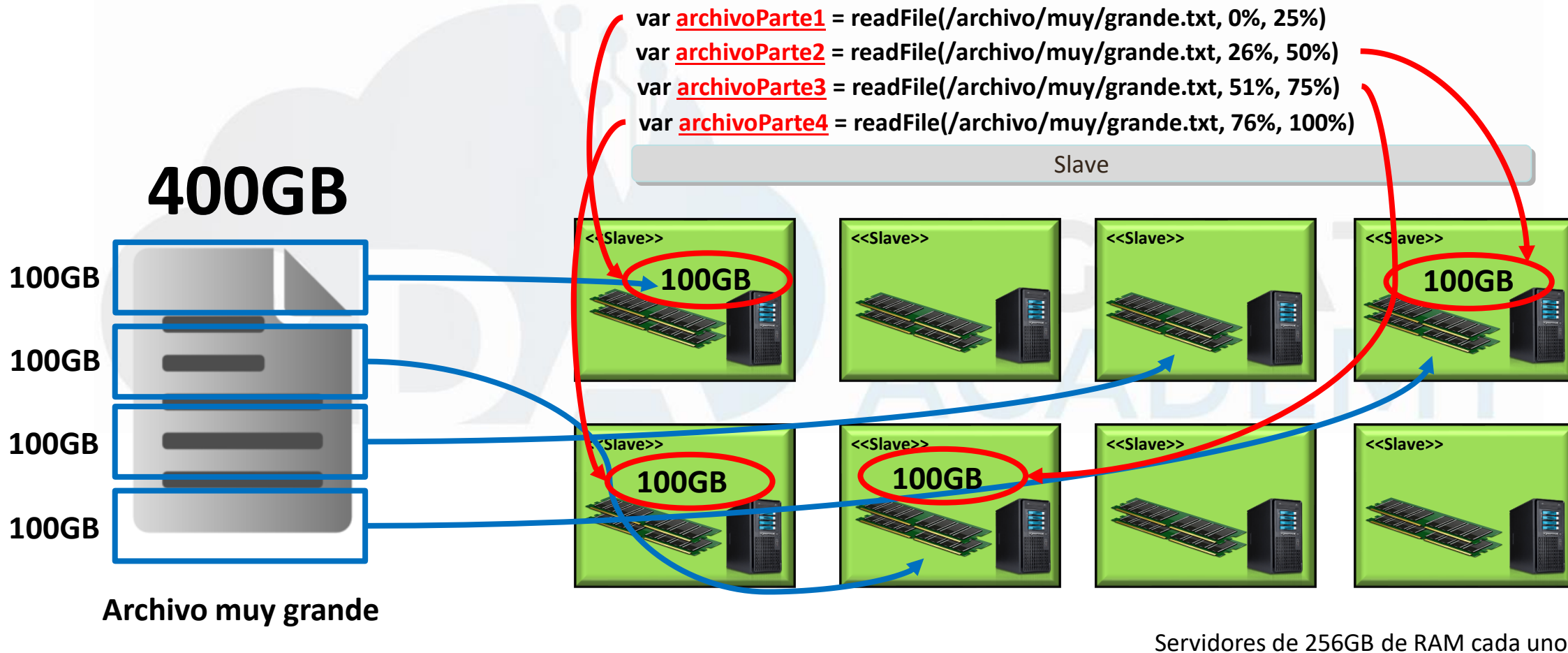
¿Cómo se crea una variable en memoria?



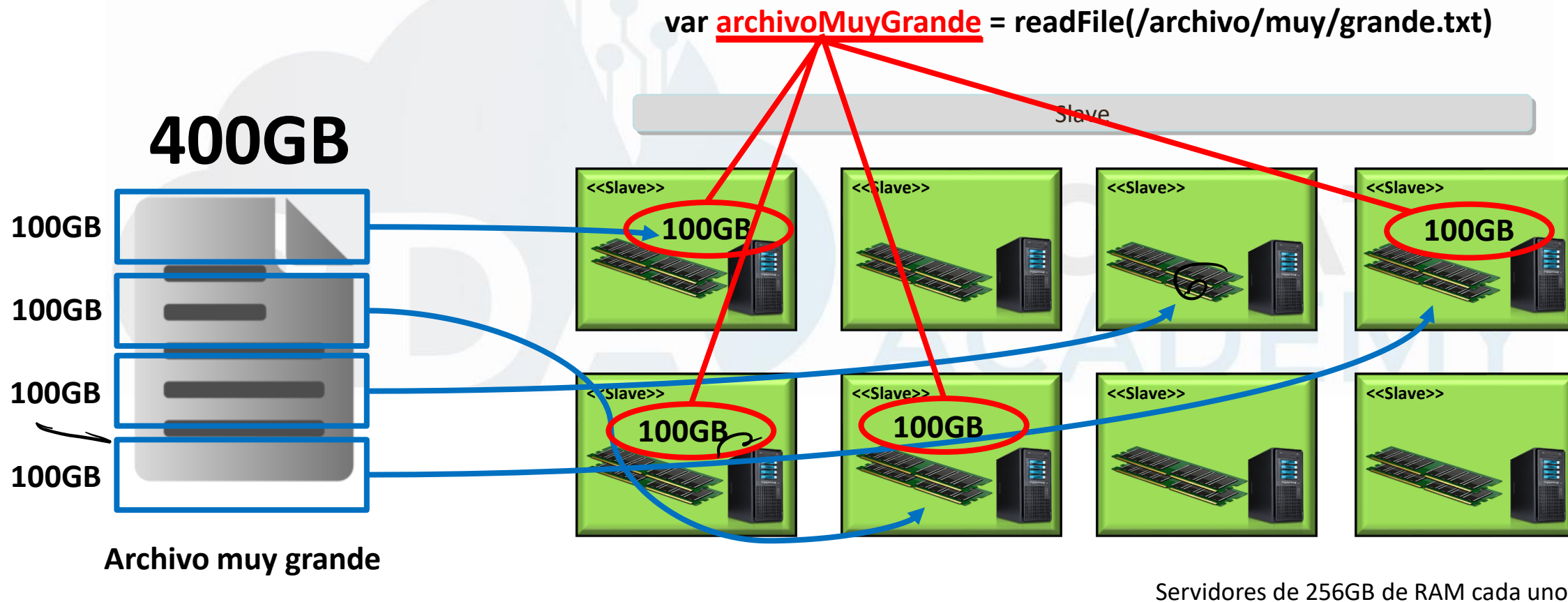
¿Y si tengo un archivo muy grande?



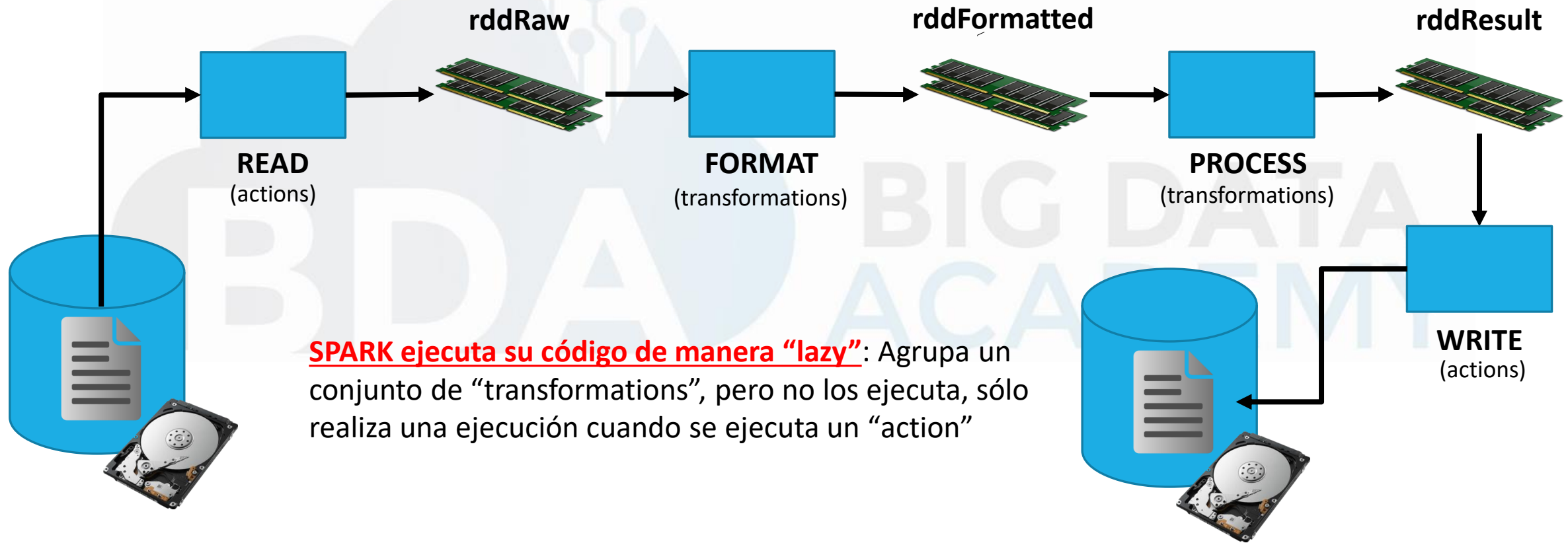
En un clúster clásico



RDD: Resilient Distributed Dataset



Arquetipo de procesamiento en SPARK



Equivalente de un GROUP BY

```
rddPersonaGroupBy = sc.textFile("/dataset/persona.data"). \
map(lambda line : line.split("|")). \
filter(lambda register : register[0] != "ID"). \
map(lambda register : ( \
    int(register[7]), \
    ( \
        1, \
        float(register[6]), \
        float(register[6]), \
        float(register[6]), \
        int(register[5]) \
    ) \
) \
). \
reduceByKey(lambda v1, v2 : ( \
    v1[0] + v2[0], \
    v1[1] + v2[1], \
    max(v1[2], v2[2]), \
    min(v1[3], v2[3]), \
    v1[4] + v2[4] \
)). \
map(lambda register : (( \
    register[0], \
    ( \
        register[1][0], \
        register[1][1], \
        register[1][2], \
        register[1][3], \
        register[1][4]/register[1][0] \
    ) \
))). \
sortBy(lambda register : register[0])
```

```
df5 = dfData.\
groupBy("EDAD").\
agg(\
    f.count("EDAD"), \
    f.min("FECHA_INGRESO"), \
    f.sum("SALARIO"), \
    f.max("SALARIO")\
)
```

Agregando estructura a los RDD: Los Dataframes

RDD

+

METADATA

=

DATAFRAME



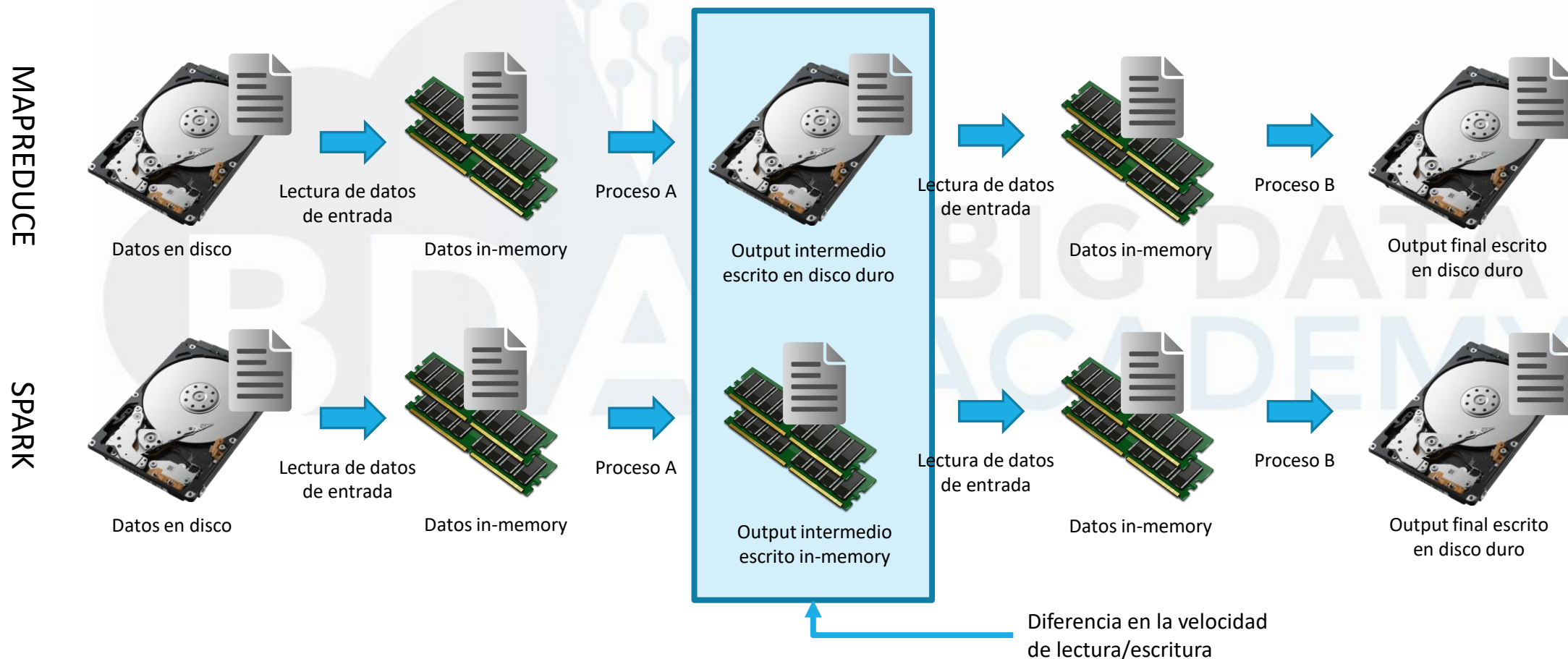
(CAMPO1 STRING,
CAMPO2 INT,
CAMPO3 DOUBLE,
...)



Procesamiento in-memory

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

¿Por qué SPARK es más rápido que MapReduce?



Procesamiento in-memory

Lectura

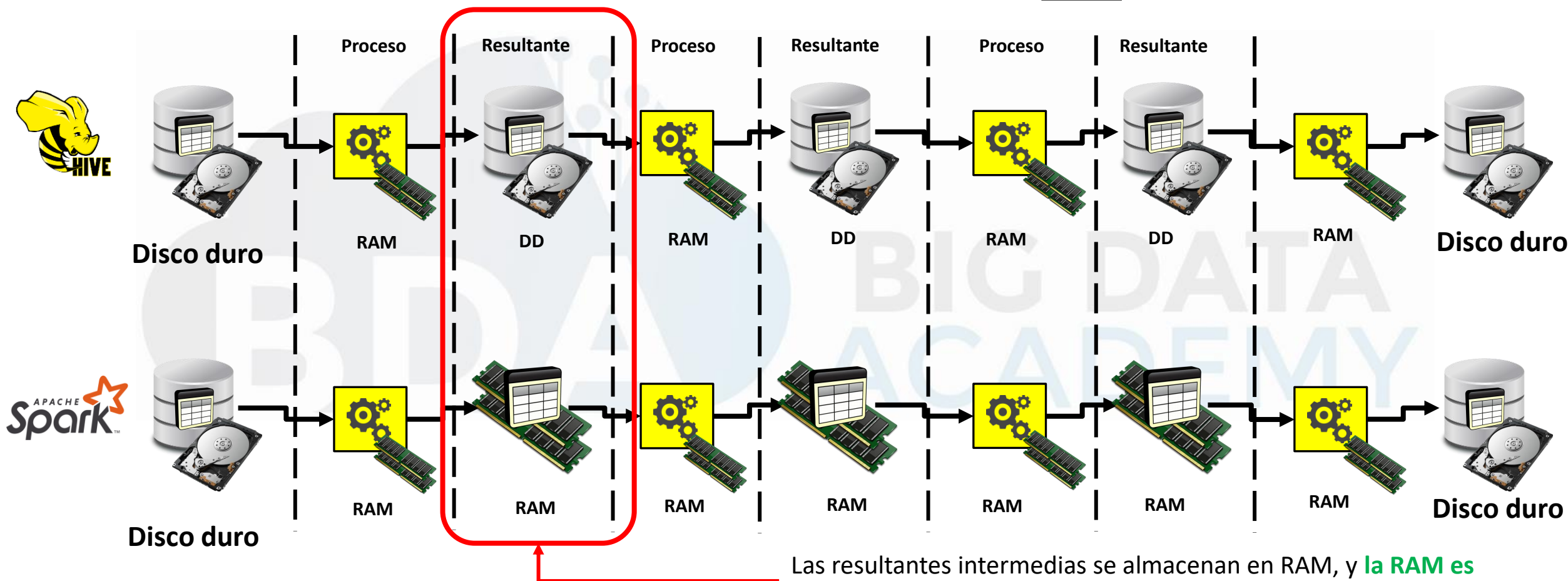
Procesamiento

Escritura

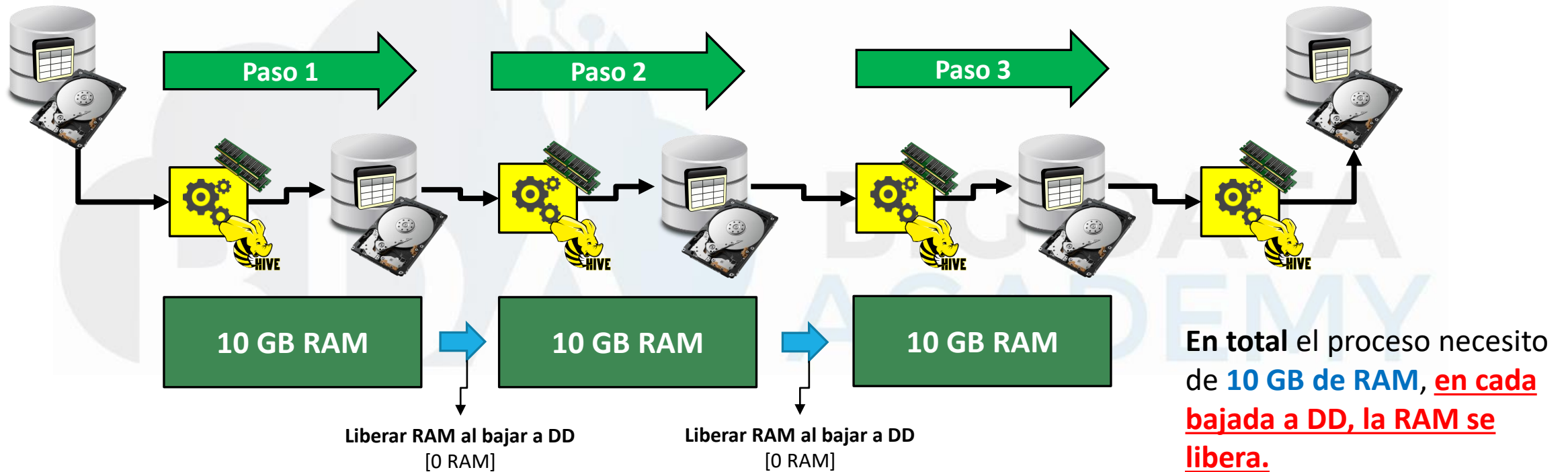
Paso 1

Paso 2

Paso 3

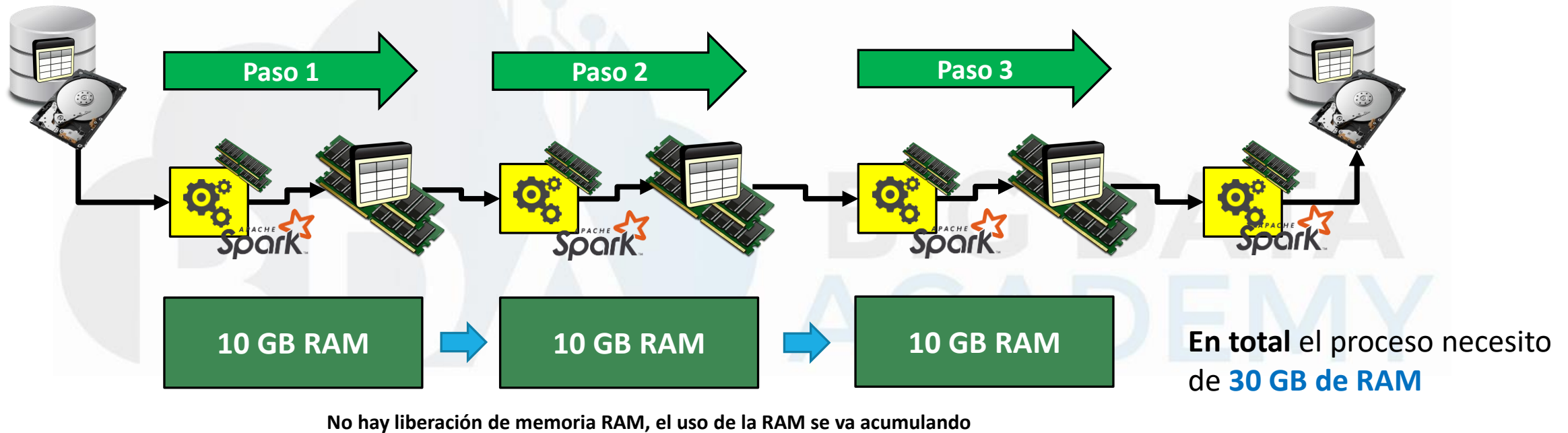


¿El procesamiento in-memory es mejor?



Un proceso que intercala procesamiento entra RAM y DD, no consume mucha memoria RAM

¿El procesamiento in-memory es mejor?



Un proceso que usa intensivamente la RAM, consume mucha memoria RAM

Mientras más pasos, más memoria RAM

```
1 #PASO 1
2 df1 = dfData.groupby(dfData["EDAD"]).agg(
3     f.count(dfData["EDAD"]).alias("CANTIDAD"),
4     f.min(dfData["FECHA_INGRESO"]).alias("FECHA_CONTRATO_MAS_RECIENTE"),
5     f.sum(dfData["SALARIO"]).alias("SUMA_SALARIOS"),
6     f.max(dfData["SALARIO"]).alias("SALARIO_MAYOR")
7 )
8
9 #PASO 2
10 df2 = df1.filter(df1["EDAD"] > 35)
11
12 #...
13
14 #PASO 10
15 df10 = df9.filter(df2["SUMA_SALARIOS"] > 500000)
16
17 #RESULTANTE
18 dfResultado = df3.filter(df3["SALARIO_MAYOR"] > 1000)
19 dfResultado.show()
```



¿Encadenar todo el proceso en un único gran paso lo soluciona?

¿Y si tengo **CIEN** pasos, o **MIL** pasos?

¿Y si encadeno procesos, optimizo?

```
1 dfResultado = dfData.groupby(dfData["EDAD"]).agg( #PASO 1
2     f.count(dfData["EDAD"]).alias("CANTIDAD"),
3     f.min(dfData["FECHA_INGRESO"]).alias("FECHA_CONTRATO_I"),
4     f.sum(dfData["SALARIO"]).alias("SUMA_SALARIOS"),
5     f.max(dfData["SALARIO"]).alias("SALARIO_MAYOR")
6 ).alias("D").\
7 filter(f.col("D.EDAD") > 35).\ #PASO 2
8 ....
9 filter(f.col("D.SUMA_SALARIOS") > 5000).\ #PASO 10
10 filter(f.col("D.SALARIO_MAYOR") > 1000) #RESULTANTE
```



Ejecutas un único gran paso, que seguirá usando la misma cantidad de memoria RAM, **no estamos optimizando.**

¿Nuestro script tiene la suficiente memoria RAM (100 GB) para funcionar?