

# INFRAESTRUCTURA Y TECNOLOGÍAS PARA BIG DATA

Vicente Castelló Ferrer

# ÍNDICE

- Arquitectura de datos.
- Almacenamiento de archivos.
- Hadoop
- Spark
- Fabric

## arquitectura.

(Del lat. *architectūra*).

1. f. Arte de proyectar y construir edificios.

2. f. *Inform.* Estructura lógica y física de los componentes de un computador.



34.495.233 €

Importe Ventas

## ESTRUCTURA



ORIGEN



APLICACIONES



BASES DE DATOS



ARCHIVOS



EVENTOS / APIS



ESTRUCTURA

ACCESOS



ORIGEN

ACCESO



APLICACIONES



BASES DE DATOS



ARCHIVOS



EVENTOS / APIS



ETL / ELT



CONECTORES



MODELADO



STREAMING

ESTRUCTURA

ACCESOS

ALMACENAMIENTO





ORIGEN

ACCESO

ALMACENAMIENTO



APLICACIONES



BASES DE DATOS



ARCHIVOS



EVENTOS



ETL / ELT



CONECTORES



MODELADO



STREAMING



databricks



DELTA LAKE



DATA LAKES



amazon  
REDSHIFT



Google  
BigQuery



DATA WAREHOUSES



amazon  
S3

ARCHIVOS / OBJETOS

Microsoft Azure  
Blob Storage



FVMP  
Federació Valenciana  
de Municipis i Províncies

ESTRUCTURA

ACCESOS

ALMACENAMIENTO

DECORACIÓN



ORIGEN

ACCESO

ALMACENAMIENTO

CONSULTA



APLICACIONES



BASES DE DATOS



ARCHIVOS



EVENTOS



ETL / ELT



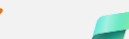
CONECTORES



MODELADO



STREAMING



DATA LAKES



DATA WAREHOUSES



ARCHIVOS / OBJETOS



CONSULTA HADOOP



CONSULTA



PREPARACIÓN



FVMP  
Federació Valenciana  
de Municipis i Províncies



ESTRUCTURA

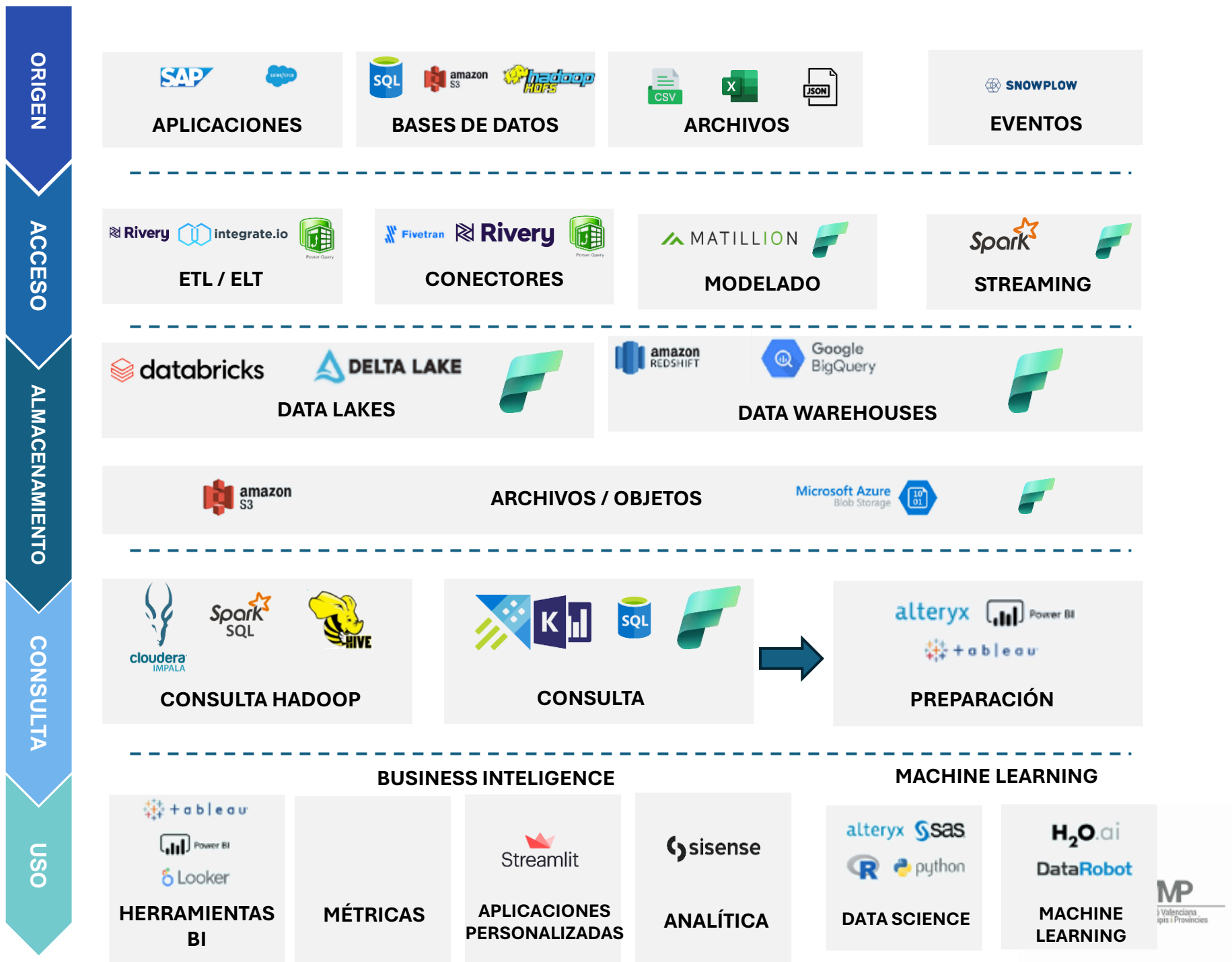
ACCESOS

ALMACENAMIENTO

DECORACIÓN

USO







GESTIÓN

CALIDAD

SEGURIDAD /  
GOBERNANZA

OPERACIONES

MANTENIMIENTO

COLABORACIÓN

ESTRUCTURA

ACCESOS

ALMACENAMIENTO

DECORACIÓN

USO



GESTIÓN DE METADATOS

CALIDAD

SEGURIDAD / GOBERNANZA

OPERACIONES

MANTENIMIENTO

COLABORACIÓN

OPERACIONES ML

ORIGEN

ACCESO

ALMACENAMIENTO

CONSULTA

USO



APLICACIONES



BASES DE DATOS



ARCHIVOS



SNOWFLOW

EVENTOS



ETL / ELT



CONECTORES



MATILLION



MODELADO



STREAMING



databricks



DELTA LAKE



DATA LAKES



amazon REDSHIFT



Google BigQuery



DATA WAREHOUSES



amazon S3

ARCHIVOS / OBJETOS

Microsoft Azure

Blob Storage



cloudera IMPALA



Spark SQL



HIVE

CONSULTA HADOOP



CONSULTA



Power BI



PREPARACIÓN

BUSINESS INTELIGENCIA

MACHINE LEARNING



+ a b l e a u



Power BI



Looker

HERRAMIENTAS BI

MÉTRICAS



Streamlit

APLICACIONES PERSONALIZADAS



sisense

ANALÍTICA



alteryx SAS



R python

DATA SCIENCE



H2O.ai



DataRobot

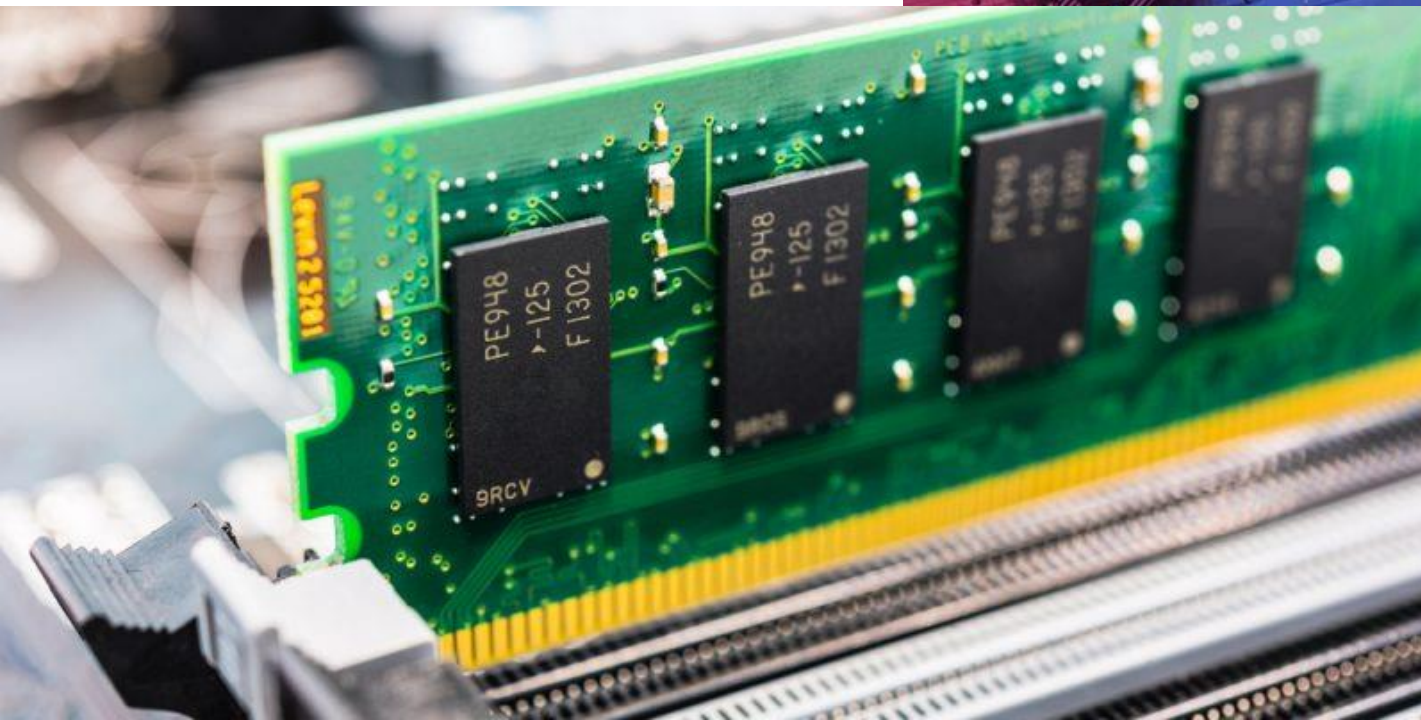
MACHINE LEARNING



MP  
Valenci  
pts i Provincies



# MEMORIA VS DISCO DURO





# ESTRUCTURA





# ESTRUCTURA DE ARCHIVO

	Ayuntamiento	Expediente	Área	Empresa	Importe
Fila 0	Ademuz	AD-0001	Urbanismo	Empresa 1	100.000 €
Fila 1	Alboraya	AL-0001	Urbanismo	Empresa 2	150.000 €
Fila 2	Valencia	VA-0901	Hacienda	Empresa 3	25.000 €
Fila 3	Manises	MA-0001	Urbanismo	Empresa 1	250.000 €
Fila 4	Alboraya	AL-3002	Contratación	Empresa 5	5.000 €
Fila 5	Alfajar	AL-5001	Comercio	Empresa 6	1.500 €

	Columna A	Columna B	Columna C	Columna D	Columna E
Fila 0	A0	B0	C0	D0	E0
Fila 1	A1	B1	C1	D1	E1
Fila 2	A2	B2	C2	D2	E2
Fila 3	A3	B3	C3	D3	E3
Fila 4	A4	B4	C4	D4	E4
Fila 5	A5	B5	C5	D5	E5

# ALMACENAMIENTO SECUENCIAL (BASE DE DATOS RELACIONAL)

	Ayuntamiento	Expediente	Área	Empresa	Importe
Fila 0	Ademuz	AD-0001	Urbanismo	Empresa 1	100.000 €
Fila 1	Alboraya	AL-0001	Urbanismo	Empresa 2	150.000 €
Fila 2	Valencia	VA-0901	Hacienda	Empresa 3	25.000 €
Fila 3	Manises	MA-0001	Urbanismo	Empresa 1	250.000 €
Fila 4	Alboraya	AL-3002	Contratación	Empresa 5	5.000 €
Fila 5	Alfajar	AL-5001	Comercio	Empresa 6	1.500 €

	Columna A	Columna B	Columna C	Columna D	Columna E
Fila 0	A0	B0	C0	D0	E0
Fila 1	A1	B1	C1	D1	E1
Fila 2	A2	B2	C2	D2	E2
Fila 3	A3	B3	C3	D3	E3
Fila 4	A4	B4	C4	D4	E4
Fila 5	A5	B5	C5	D5	E5

A0	B0	C0	D0	E0	A1	B1	C1	D1	E1	A2	B2	C2	D2	E2
A3	B3	C3	D3	E3	A4	B4	C4	D4	E4	A5	B5	C5	D5	E5

# DAME TODA LA INFORMACIÓN DEL EXPEDIENTE MA-0001

	Ayuntamiento	Expediente	Área	Empresa	Importe
Fila 0	Ademuz	AD-0001	Urbanismo	Empresa 1	100.000 €
Fila 1	Alboraya	AL-0001	Urbanismo	Empresa 2	150.000 €
Fila 2	Valencia	VA-0901	Hacienda	Empresa 3	25.000 €
Fila 3	Manises	MA-0001	Urbanismo	Empresa 1	250.000 €
Fila 4	Alboraya	AL-3002	Contratación	Empresa 5	5.000 €
Fila 5	Alfajar	AL-5001	Comercio	Empresa 6	1.500 €

	Columna A	Columna B	Columna C	Columna D	Columna E
Fila 0	A0	B0	C0	D0	E0
Fila 1	A1	B1	C1	D1	E1
Fila 2	A2	B2	C2	D2	E2
Fila 3	A3	B3	C3	D3	E3
Fila 4	A4	B4	C4	D4	E4
Fila 5	A5	B5	C5	D5	E5

A0	B0	C0	D0	E0	A1	B1	C1	D1	E1	A2	B2	C2	D2	E2
A3	B3	C3	D3	E3	A4	B4	C4	D4	E4	A5	B5	C5	D5	E5

# DAME TODA LA INFORMACIÓN DEL EXPEDIENTE MA-0001

	Ayuntamiento	Expediente	Área	Empresa	Importe
Fila 0	Ademuz	AD-0001	Urbanismo	Empresa 1	100.000 €
Fila 1	Alboraya	AL-0001	Urbanismo	Empresa 2	150.000 €
Fila 2	Valencia	VA-0901	Hacienda	Empresa 3	25.000 €
Fila 3	Manises	MA-0001	Urbanismo	Empresa 1	250.000 €
Fila 4	Alboraya	AL-3002	Contratación	Empresa 5	5.000 €
Fila 5	Alfajar	AL-5001	Comercio	Empresa 6	1.500 €

	Columna A	Columna B	Columna C	Columna D	Columna E
Fila 0	A0	B0	C0	D0	E0
Fila 1	A1	B1	C1	D1	E1
Fila 2	A2	B2	C2	D2	E2
Fila 3	A3	B3	C3	D3	E3
Fila 4	A4	B4	C4	D4	E4
Fila 5	A5	B5	C5	D5	E5

A0	B0	C0	D0	E0	A1	B1	C1	D1	E1	A2	B2	C2	D2	E2
A3	B3	C3	D3	E3	A4	B4	C4	D4	E4	A5	B5	C5	D5	E5

DAME TODA LA INFORMACIÓN DEL EXPEDIENTE MA-0001

	Ayuntamiento	Expediente	Área	Empresa	Importe
Fila 0	Ademuz	AD-0001	Urbanismo	Empresa 1	100.000 €
Fila 1	Alboraya	AL-0001	Urbanismo	Empresa 2	150.000 €
Fila 2	Valencia	VA-0901	Hacienda	Empresa 3	25.000 €
Fila 3	Manises	MA-0001	Urbanismo	Empresa 1	250.000 €
Fila 4	Alboraya	AL-3002	Contratación	Empresa 5	5.000 €
Fila 5	Alfajar	AL-5001	Comercio	Empresa 6	1.500 €

	Columna A	Columna B	Columna C	Columna D	Columna E
Fila 0	A0	B0	C0	D0	E0
Fila 1	A1	B1	C1	D1	E1
Fila 2	A2	B2	C2	D2	E2
Fila 3	A3	B3	C3	D3	E3
Fila 4	A4	B4	C4	D4	E4
Fila 5	A5	B5	C5	D5	E5

A0	B0	C0	D0	E0	A1	B1	C1	D1	E1	A2	B2	C2	D2	E2
A3	B3	C3	D3	E3	A4	B4	C4	D4	E4	A5	B5	C5	D5	E5

# DAME LA SUMA DEL IMPORTE DE TODOS LOS EXPEDIENTES

	Ayuntamiento	Expediente	Área	Empresa	Importe
Fila 0	Ademuz	AD-0001	Urbanismo	Empresa 1	100.000 €
Fila 1	Alboraya	AL-0001	Urbanismo	Empresa 2	150.000 €
Fila 2	Valencia	VA-0901	Hacienda	Empresa 3	25.000 €
Fila 3	Manises	MA-0001	Urbanismo	Empresa 1	250.000 €
Fila 4	Alboraya	AL-3002	Contratación	Empresa 5	5.000 €
Fila 5	Alfajar	AL-5001	Comercio	Empresa 6	1.500 €

	Columna A	Columna B	Columna C	Columna D	Columna E
Fila 0	A0	B0	C0	D0	E0
Fila 1	A1	B1	C1	D1	E1
Fila 2	A2	B2	C2	D2	E2
Fila 3	A3	B3	C3	D3	E3
Fila 4	A4	B4	C4	D4	E4
Fila 5	A5	B5	C5	D5	E5

A0	B0	C0	D0	E0	A1	B1	C1	D1	E1	A2	B2	C2	D2	E2
A3	B3	C3	D3	E3	A4	B4	C4	D4	E4	A5	B5	C5	D5	E5

# DAME LA SUMA DEL IMPORTE DE TODOS LOS EXPEDIENTES

	Ayuntamiento	Expediente	Área	Empresa	Importe
Fila 0	Ademuz	AD-0001	Urbanismo	Empresa 1	100.000 €
Fila 1	Alboraya	AL-0001	Urbanismo	Empresa 2	150.000 €
Fila 2	Valencia	VA-0901	Hacienda	Empresa 3	25.000 €
Fila 3	Manises	MA-0001	Urbanismo	Empresa 1	250.000 €
Fila 4	Alboraya	AL-3002	Contratación	Empresa 5	5.000 €
Fila 5	Alfajar	AL-5001	Comercio	Empresa 6	1.500 €

	Columna A	Columna B	Columna C	Columna D	Columna E
Fila 0	A0	B0	C0	D0	E0
Fila 1	A1	B1	C1	D1	E1
Fila 2	A2	B2	C2	D2	E2
Fila 3	A3	B3	C3	D3	E3
Fila 4	A4	B4	C4	D4	E4
Fila 5	A5	B5	C5	D5	E5

A0	B0	C0	D0	E0	A1	B1	C1	D1	E1	A2	B2	C2	D2	E2
A3	B3	C3	D3	E3	A4	B4	C4	D4	E4	A5	B5	C5	D5	E5



# DAME LA SUMA DEL IMPORTE DE TODOS LOS EXPEDIENTES

	Ayuntamiento	Expediente	Área	Empresa	Importe
Fila 0	Ademuz	AD-0001	Urbanismo	Empresa 1	100.000 €
Fila 1	Alboraya	AL-0001	Urbanismo	Empresa 2	150.000 €
Fila 2	Valencia	VA-0901	Hacienda	Empresa 3	25.000 €
Fila 3	Manises	MA-0001	Urbanismo	Empresa 1	250.000 €
Fila 4	Alboraya	AL-3002	Contratación	Empresa 5	5.000 €
Fila 5	Alfajar	AL-5001	Comercio	Empresa 6	1.500 €

	Columna A	Columna B	Columna C	Columna D	Columna E
Fila 0	A0	B0	C0	D0	E0
Fila 1	A1	B1	C1	D1	E1
Fila 2	A2	B2	C2	D2	E2
Fila 3	A3	B3	C3	D3	E3
Fila 4	A4	B4	C4	D4	E4
Fila 5	A5	B5	C5	D5	E5

A0	B0	C0	D0	E0	A1	B1	C1	D1	E1	A2	B2	C2	D2	E2
A3	B3	C3	D3	E3	A4	B4	C4	D4	E4	A5	B5	C5	D5	E5

# DAME LA SUMA DEL IMPORTE DE TODOS LOS EXPEDIENTES

	Ayuntamiento	Expediente	Área	Empresa	Importe
Fila 0	Ademuz	AD-0001	Urbanismo	Empresa 1	100.000 €
Fila 1	Alboraya	AL-0001	Urbanismo	Empresa 2	150.000 €
Fila 2	Valencia	VA-0901	Hacienda	Empresa 3	25.000 €
Fila 3	Manises	MA-0001	Urbanismo	Empresa 1	250.000 €
Fila 4	Alboraya	AL-3002	Contratación	Empresa 5	5.000 €
Fila 5	Alfajar	AL-5001	Comercio	Empresa 6	1.500 €

	Columna A	Columna B	Columna C	Columna D	Columna E
Fila 0	A0	B0	C0	D0	E0
Fila 1	A1	B1	C1	D1	E1
Fila 2	A2	B2	C2	D2	E2
Fila 3	A3	B3	C3	D3	E3
Fila 4	A4	B4	C4	D4	E4
Fila 5	A5	B5	C5	D5	E5

A0	B0	C0	D0	E0	A1	B1	C1	D1	E1	A2	B2	C2	D2	E2
A3	B3	C3	D3	E3	A4	B4	C4	D4	E4	A5	B5	C5	D5	E5

# DAME LA SUMA DEL IMPORTE DE TODOS LOS EXPEDIENTES

	Ayuntamiento	Expediente	Área	Empresa	Importe
Fila 0	Ademuz	AD-0001	Urbanismo	Empresa 1	100.000 €
Fila 1	Alboraya	AL-0001	Urbanismo	Empresa 2	150.000 €
Fila 2	Valencia	VA-0901	Hacienda	Empresa 3	25.000 €
Fila 3	Manises	MA-0001	Urbanismo	Empresa 1	250.000 €
Fila 4	Alboraya	AL-3002	Contratación	Empresa 5	5.000 €
Fila 5	Alfajar	AL-5001	Comercio	Empresa 6	1.500 €

	Columna A	Columna B	Columna C	Columna D	Columna E
Fila 0	A0	B0	C0	D0	E0
Fila 1	A1	B1	C1	D1	E1
Fila 2	A2	B2	C2	D2	E2
Fila 3	A3	B3	C3	D3	E3
Fila 4	A4	B4	C4	D4	E4
Fila 5	A5	B5	C5	D5	E5

A0	B0	C0	D0	E0	A1	B1	C1	D1	E1	A2	B2	C2	D2	E2
A3	B3	C3	D3	E3	A4	B4	C4	D4	E4	A5	B5	C5	D5	E5

# DAME LA SUMA DEL IMPORTE DE LOS EXPEDIENTES DE URBANISMO

	Ayuntamiento	Expediente	Área	Empresa	Importe
Fila 0	Ademuz	AD-0001	Urbanismo	Empresa 1	100.000 €
Fila 1	Alboraya	AL-0001	Urbanismo	Empresa 2	150.000 €
Fila 2	Valencia	VA-0901	Hacienda	Empresa 3	25.000 €
Fila 3	Manises	MA-0001	Urbanismo	Empresa 1	250.000 €
Fila 4	Alboraya	AL-3002	Contratación	Empresa 5	5.000 €
Fila 5	Alfajar	AL-5001	Comercio	Empresa 6	1.500 €

	Columna A	Columna B	Columna C	Columna D	Columna E
Fila 0	A0	B0	C0	D0	E0
Fila 1	A1	B1	C1	D1	E1
Fila 2	A2	B2	C2	D2	E2
Fila 3	A3	B3	C3	D3	E3
Fila 4	A4	B4	C4	D4	E4
Fila 5	A5	B5	C5	D5	E5

A0	B0	C0	D0	E0	A1	B1	C1	D1	E1	A2	B2	C2	D2	E2
A3	B3	C3	D3	E3	A4	B4	C4	D4	E4	A5	B5	C5	D5	E5

# DAME LA SUMA DEL IMPORTE DE LOS EXPEDIENTES DE URBANISMO

	Ayuntamiento	Expediente	Área	Empresa	Importe
Fila 0	Ademuz	AD-0001	Urbanismo	Empresa 1	100.000 €
Fila 1	Alboraya	AL-0001	Urbanismo	Empresa 2	150.000 €
Fila 2	Valencia	VA-0901	Hacienda	Empresa 3	25.000 €
Fila 3	Manises	MA-0001	Urbanismo	Empresa 1	250.000 €
Fila 4	Alboraya	AL-3002	Contratación	Empresa 5	5.000 €
Fila 5	Alfajar	AL-5001	Comercio	Empresa 6	1.500 €

	Columna A	Columna B	Columna C	Columna D	Columna E
Fila 0	A0	B0	C0	D0	E0
Fila 1	A1	B1	C1	D1	E1
Fila 2	A2	B2	C2	D2	E2
Fila 3	A3	B3	C3	D3	E3
Fila 4	A4	B4	C4	D4	E4
Fila 5	A5	B5	C5	D5	E5

A0	B0	C0	D0	E0	A1	B1	C1	D1	E1	A2	B2	C2	D2	E2
A3	B3	C3	D3	E3	A4	B4	C4	D4	E4	A5	B5	C5	D5	E5

# DAME LA SUMA DEL IMPORTE DE LOS EXPEDIENTES DE URBANISMO

	Ayuntamiento	Expediente	Área	Empresa	Importe
Fila 0	Ademuz	AD-0001	Urbanismo	Empresa 1	100.000 €
Fila 1	Alboraya	AL-0001	Urbanismo	Empresa 2	150.000 €
Fila 2	Valencia	VA-0901	Hacienda	Empresa 3	25.000 €
Fila 3	Manises	MA-0001	Urbanismo	Empresa 1	250.000 €
Fila 4	Alboraya	AL-3002	Contratación	Empresa 5	5.000 €
Fila 5	Alfajar	AL-5001	Comercio	Empresa 6	1.500 €

	Columna A	Columna B	Columna C	Columna D	Columna E
Fila 0	A0	B0	C0	D0	E0
Fila 1	A1	B1	C1	D1	E1
Fila 2	A2	B2	C2	D2	E2
Fila 3	A3	B3	C3	D3	E3
Fila 4	A4	B4	C4	D4	E4
Fila 5	A5	B5	C5	D5	E5

A0	B0	C0	D0	E0	A1	B1	C1	D1	E1	A2	B2	C2	D2	E2
A3	B3	C3	D3	E3	A4	B4	C4	D4	E4	A5	B5	C5	D5	E5

# ESTRUCTURA DE ARCHIVO

	Ayuntamiento	Expediente	Área	Empresa	Importe
Fila 0	Ademuz	AD-0001	Urbanismo	Empresa 1	100.000 €
Fila 1	Alboraya	AL-0001	Urbanismo	Empresa 2	150.000 €
Fila 2	Valencia	VA-0901	Hacienda	Empresa 3	25.000 €
Fila 3	Manises	MA-0001	Urbanismo	Empresa 1	250.000 €
Fila 4	Alboraya	AL-3002	Contratación	Empresa 5	5.000 €
Fila 5	Alfajar	AL-5001	Comercio	Empresa 6	1.500 €

	Columna A	Columna B	Columna C	Columna D	Columna E
Fila 0	A0	B0	C0	D0	E0
Fila 1	A1	B1	C1	D1	E1
Fila 2	A2	B2	C2	D2	E2
Fila 3	A3	B3	C3	D3	E3
Fila 4	A4	B4	C4	D4	E4
Fila 5	A5	B5	C5	D5	E5

# ESTRUCTURA TABULAR - COLUMNAS

	Ayuntamiento	Expediente	Área	Empresa	Importe
Fila 0	Ademuz	AD-0001	Urbanismo	Empresa 1	100.000 €
Fila 1	Alboraya	AL-0001	Urbanismo	Empresa 2	150.000 €
Fila 2	Valencia	VA-0901	Hacienda	Empresa 3	25.000 €
Fila 3	Manises	MA-0001	Urbanismo	Empresa 1	250.000 €
Fila 4	Alboraya	AL-3002	Contratación	Empresa 5	5.000 €
Fila 5	Alfajar	AL-5001	Comercio	Empresa 6	1.500 €

	Columna A	Columna B	Columna C	Columna D	Columna E
Fila 0	A0	B0	C0	D0	E0
Fila 1	A1	B1	C1	D1	E1
Fila 2	A2	B2	C2	D2	E2
Fila 3	A3	B3	C3	D3	E3
Fila 4	A4	B4	C4	D4	E4
Fila 5	A5	B5	C5	D5	E5

A0	A1	A2	A3	A4	A5	B0	B1	B2	B3	B4	B5	C0	C1	C2
C3	C4	C5	D0	D1	D2	D3	D4	D5	E0	E1	E2	E3	E4	E5



# ESTRUCTURA TABULAR - COLUMNAS

	Ayuntamiento	Expediente	Área	Empresa	Importe
Fila 0	Ademuz	AD-0001	Urbanismo	Empresa 1	100.000 €
Fila 1	Alboraya	AL-0001	Urbanismo	Empresa 2	150.000 €
Fila 2	Valencia	VA-0901	Hacienda	Empresa 3	25.000 €
Fila 3	Manises	MA-0001	Urbanismo	Empresa 1	250.000 €
Fila 4	Alboraya	AL-3002	Contratación	Empresa 5	5.000 €
Fila 5	Alfajar	AL-5001	Comercio	Empresa 6	1.500 €

	Columna A	Columna B	Columna C	Columna D	Columna E
Fila 0	A0	B0	C0	D0	E0
Fila 1	A1	B1	C1	D1	E1
Fila 2	A2	B2	C2	D2	E2
Fila 3	A3	B3	C3	D3	E3
Fila 4	A4	B4	C4	D4	E4
Fila 5	A5	B5	C5	D5	E5

A0	A1	A2	A3	A4	A5	B0	B1	B2	B3	B4	B5	C0	C1	C2
C3	C4	C5	D0	D1	D2	D3	D4	D5	E0	E1	E2	E3	E4	E5

# ESTRUCTURA DE BASE DE DATOS - ¿QUÉ ES LO QUE QUIERO HACER?

A

Realizar MUCHAS operaciones PEQUEÑAS que me devuelven una FILA COMPLETA



**OLTP – Online Transaction Processing**

B

Realizar POCAS operaciones GRANDES que me devuelven TODAS las columnas



**OLAP – Online Analytical Processing**

**HECHO** – Los accesos I/O cuestan – Memoria, disco, red, ...

# ALMACENAMIENTO SECUENCIAL

	Ayuntamiento	Expediente	Área	Empresa	Importe
Fila 0	Ademuz	AD-0001	Urbanismo	Empresa 1	100.000 €
Fila 1	Alboraya	AL-0001	Urbanismo	Empresa 2	150.000 €
Fila 2	Valencia	VA-0901	Hacienda	Empresa 3	25.000 €
Fila 3	Manises	MA-0001	Urbanismo	Empresa 1	250.000 €
Fila 4	Alboraya	AL-3002	Contratación	Empresa 5	5.000 €
Fila 5	Alfajar	AL-5001	Comercio	Empresa 6	1.500 €



OLTP

	Columna A	Columna B	Columna C	Columna D	Columna E
Fila 0	A0	B0	C0	D0	E0
Fila 1	A1	B1	C1	D1	E1
Fila 2	A2	B2	C2	D2	E2
Fila 3	A3	B3	C3	D3	E3
Fila 4	A4	B4	C4	D4	E4
Fila 5	A5	B5	C5	D5	E5



OLAP

A0	B0	C0	D0	E0	A1	B1	C1	D1	E1	A2	B2	C2	D2	E2
A3	B3	C3	D3	E3	A4	B4	C4	D4	E4	A5	B5	C5	D5	E5

# ESTRUCTURA TABULAR - COLUMNAS

	Ayuntamiento	Expediente	Área	Empresa	Importe
Fila 0	Ademuz	AD-0001	Urbanismo	Empresa 1	100.000 €
Fila 1	Alboraya	AL-0001	Urbanismo	Empresa 2	150.000 €
Fila 2	Valencia	VA-0901	Hacienda	Empresa 3	25.000 €
Fila 3	Manises	MA-0001	Urbanismo	Empresa 1	250.000 €
Fila 4	Alboraya	AL-3002	Contratación	Empresa 5	5.000 €
Fila 5	Alfajar	AL-5001	Comercio	Empresa 6	1.500 €



OLAP

	Columna A	Columna B	Columna C	Columna D	Columna E
Fila 0	A0	B0	C0	D0	E0
Fila 1	A1	B1	C1	D1	E1
Fila 2	A2	B2	C2	D2	E2
Fila 3	A3	B3	C3	D3	E3
Fila 4	A4	B4	C4	D4	E4
Fila 5	A5	B5	C5	D5	E5



OLTP

A0	A1	A2	A3	A4	A5	B0	B1	B2	B3	B4	B5	C0	C1	C2
C3	C4	C5	D0	D1	D2	D3	D4	D5	E0	E1	E2	E3	E4	E5

	Expediente	Ayuntamiento	Área	Empresa	Importe
Fila 0	AD-0001	Ademuz	Urbanismo	Empresa 1	100.000 €
Fila 1	AL-0001	Alboraya	Urbanismo	Empresa 2	150.000 €
Fila 2	VA-0901	Valencia	Hacienda	Empresa 3	25.000 €
Fila 3	MA-0001	Manises	Urbanismo	Empresa 1	250.000 €
Fila 4	AL-3002	Alboraya	Contratación	Empresa 5	5.000 €
Fila 5	AL-5001	Alfafar	Comercio	Empresa 6	1.500 €

# ESTRUCTURA HÍBRIDA - PARQUET



	Expediente	Ayuntamiento	Área	Empresa	Importe	
A	Fila 0	AD-0001	Ademuz	Urbanismo	Empresa 1	100.000 €
	Fila 1	AL-0001	Alboraya	Urbanismo	Empresa 2	150.000 €
B	Fila 2	VA-0901	Valencia	Hacienda	Empresa 3	25.000 €
	Fila 3	MA-0001	Manises	Urbanismo	Empresa 1	250.000 €
C	Fila 4	AL-3002	Alboraya	Contratación	Empresa 5	5.000 €
	Fila 5	AL-5001	Alfafar	Comercio	Empresa 6	1.500 €

Almacenamiento **COLUMNAR** con **PARTICIONADO** de filas

# ESTRUCTURA HÍBRIDA - PARQUET



A

B

C

	Expediente	Ayuntamiento	Área	Empresa	Importe
Fila 0	AD-0001	Ademuz	Urbanismo	Empresa 1	100.000 €
Fila 1	AL-0001	Alboraya	Urbanismo	Empresa 2	150.000 €
Fila 2	VA-0901	Valencia	Hacienda	Empresa 3	25.000 €
Fila 3	MA-0001	Manises	Urbanismo	Empresa 1	250.000 €
Fila 4	AL-3002	Alboraya	Contratación	Empresa 5	5.000 €
Fila 5	AL-5001	Alfafar	Comercio	Empresa 6	1.500 €

A0	A1	A2	B0	B1	B2	C0	C1	C2	D0	D1	D2	E0	E1	E2
A3	A4	A5	B3	B4	B5	C3	C4	C5	D3	D4	D5	E3	E4	E5

DAME LA SUMA DEL IMPORTE DE LOS EXPEDIENTES DE URBANISMO



	Expediente	Ayuntamiento	Área	Empresa	Importe	
A	Fila 0	AD-0001	Ademuz	Urbanismo	Empresa 1	100.000 €
	Fila 1	AL-0001	Alboraya	Urbanismo	Empresa 2	150.000 €
B	Fila 2	VA-0901	Valencia	Hacienda	Empresa 3	25.000 €
	Fila 3	MA-0001	Manises	Urbanismo	Empresa 1	250.000 €
C	Fila 4	AL-3002	Alboraya	Contratación	Empresa 5	5.000 €
	Fila 5	AL-5001	Alfafar	Comercio	Empresa 6	1.500 €



ESTE GRUPO DE FILAS NO SE TIENE EN CUENTA



# DAME LA SUMA DEL IMPORTE DE LOS EXPEDIENTES DE URBANISMO

	Ayuntamiento	Expediente	Área	Empresa	Importe
Fila 0	Ademuz	AD-0001	Urbanismo	Empresa 1	100.000 €
Fila 1	Alboraya	AL-0001	Urbanismo	Empresa 2	150.000 €
Fila 2	Valencia	VA-0901	Hacienda	Empresa 3	25.000 €
Fila 3	Manises	MA-0001	Urbanismo	Empresa 1	250.000 €
Fila 4	Alboraya	AL-3002	Contratación	Empresa 5	5.000 €
Fila 5	Alfafar	AL-5001	Comercio	Empresa 6	1.500 €

TIPO ALMACENAMIENTO	ACCIÓN
SECUENCIAL (FILAS)	Buscamos en 5 columnas con 6 filas

# DAME LA SUMA DEL IMPORTE DE LOS EXPEDIENTES DE URBANISMO

	Ayuntamiento	Expediente	Área	Empresa	Importe
Fila 0	Ademuz	AD-0001	Urbanismo	Empresa 1	100.000 €
Fila 1	Alboraya	AL-0001	Urbanismo	Empresa 2	150.000 €
Fila 2	Valencia	VA-0901	Hacienda	Empresa 3	25.000 €
Fila 3	Manises	MA-0001	Urbanismo	Empresa 1	250.000 €
Fila 4	Alboraya	AL-3002	Contratación	Empresa 5	5.000 €
Fila 5	Alfajar	AL-5001	Comercio	Empresa 6	1.500 €

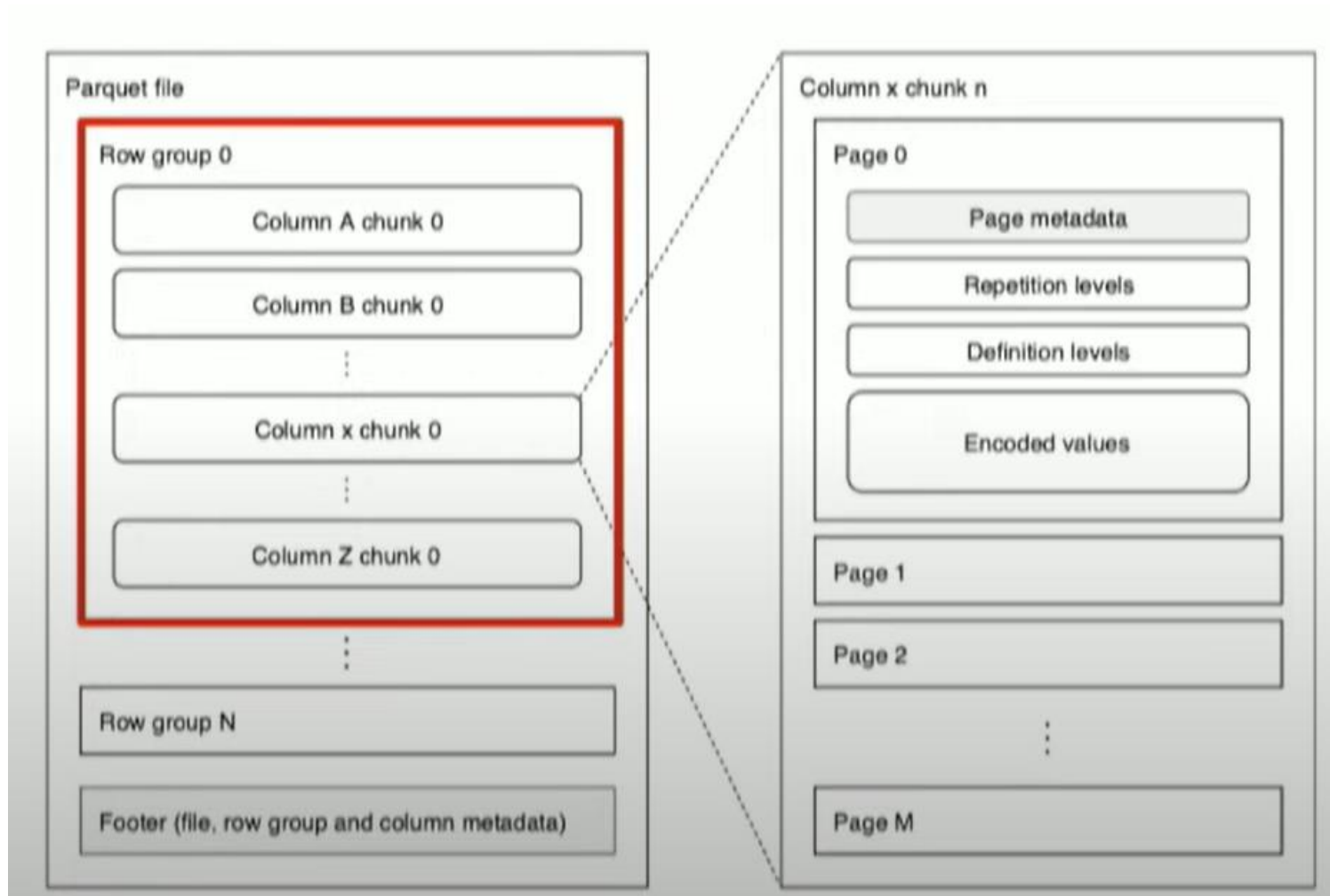
TIPO ALMACENAMIENTO	ACCIÓN
SECUENCIAL (FILAS)	Buscamos en 5 columnas con 6 filas
TABULAR (COLUMNAS)	Buscamos en 2 columnas con 6 filas

# DAME LA SUMA DEL IMPORTE DE LOS EXPEDIENTES DE URBANISMO

	Ayuntamiento	Expediente	Área	Empresa	Importe
Fila 0	Ademuz	AD-0001	Urbanismo	Empresa 1	100.000 €
Fila 1	Alboraya	AL-0001	Urbanismo	Empresa 2	150.000 €
Fila 2	Valencia	VA-0901	Hacienda	Empresa 3	25.000 €
Fila 3	Manises	MA-0001	Urbanismo	Empresa 1	250.000 €
Fila 4	Alboraya	AL-3002	Contratación	Empresa 5	5.000 €
Fila 5	Alfajar	AL-5001	Comercio	Empresa 6	1.500 €

TIPO ALMACENAMIENTO	ACCIÓN
SECUENCIAL (FILAS)	Buscamos en 5 columnas con 6 filas
TABULAR (COLUMNAS)	Buscamos en 2 columnas con 6 filas
HÍBRIDO (PARQUET)	Buscamos en 2 columnas con 4 filas

# ESTRUCTURA METADATOS



Row groups: Por defecto de 128 MB  
Trozos de columnas

Páginas:

- Metadata:
- Min
- Max
- Count

# ESTRUCTURA METADATOS



```
##### file meta data #####

created_by: parquet-cpp version 1.5.1-SNAPSHOT

num_columns: 3

num_rows: 3

num_row_groups: 1

format_version: 1.0

serialized_size: 2226

##### Columns #####

one

two

three

##### Column(one) #####

name: one

path: one

max_definition_level: 1

max_repetition_level: 0

physical_type: DOUBLE

logical_type: None
```

```
##### Column(two) #####

name: two

path: two

max_definition_level: 1

max_repetition_level: 0

physical_type: BYTE_ARRAY

logical_type: String

converted_type (legacy): UTF8

##### Column(three) #####

name: three

path: three

max_definition_level: 1

max_repetition_level: 0

physical_type: BOOLEAN

logical_type: None

converted_type (legacy): NONE
```

[¿Por qué deberías de usar ficheros Parquet si procesas muchos datos? | datos.gob.es](https://datos.gob.es/es/blog/por-que-deberias-de-usar-ficheros-parquet-si-procesas-muchos-datos)

<https://datos.gob.es/es/blog/por-que-deberias-de-usar-ficheros-parquet-si-procesas-muchos-datos>

# ESTRUCTURA METADATOS



Dataset	Size on Amazon S3	Query Run time	Data Scanned	Cost
Data stored as CSV files	1 TB	236 seconds	1.15 TB	\$5.75
Data stored in Apache Parquet format*	130 GB	6.78 seconds	2.51 GB	\$0.01
Savings / Speedup	87% less with Parquet	34x faster	99% less data scanned	99.7% savings

Dataset	Columns	Size on Amazon S3	Data Scanned	Cost
Data stored as CSV file	4	4TB	4TB	\$20 (4TB x \$5/TB)
Data stored as GZIP CSV file	4	1TB	1TB	\$5 (1TB x \$5/TB)
Data stored as Parquet file	4	1TB	.25TB	\$1.25 (.25TB x \$5/TB)

[¿Por qué deberías de usar ficheros Parquet si procesas muchos datos? | datos.gob.es](https://datos.gob.es/es/blog/por-que-deberias-de-usar-ficheros-parquet-si-procesas-muchos-datos)

<https://datos.gob.es/es/blog/por-que-deberias-de-usar-ficheros-parquet-si-procesas-muchos-datos>

## CODIFICACIÓN

1	Ademuz
2	Ador
3	Atzeneta d'Albaida
4	Agullent
5	Alaquàs
6	Albaida
7	Albal
8	Albalat de la Ribera
9	Albalat dels Sorells
10	Albalat dels Tarongers
11	Alberic
12	Alborache
.....	.....
.....	.....
261	Yátova
262	Yesa, La
263	Zarra

## RUN-LENGTH ENCODING

Ademuz	(1,345)
Ador	(2,243)
Atzeneta d'Albaida	(3,65)
Agullent	(4,235)
Alaquàs	(5,876)

	Código	Inicio	Fin
Ademuz	1	1	345
Ador	2	346	589
Atzeneta d'Albaida	3	590	655
Agullent	4	656	891
Alaquàs	5	892	1768





**Hadoop es un entorno de licencia libre para almacenar  
datos y ejecutar aplicaciones en clusters de hardware  
“estándar”**

# HISTORIA



- A principios de siglo Google empezó a tener problemas debido al aumento de usuarios de Internet.
- El problema principal era cómo almacenar los datos de sus usuarios en servidores tradicionales.
  - 4,1 millones de búsquedas.
  - Cada búsqueda implicaba la consulta de cientos de MB de datos (alto coste computacional).
  - Tolerancia a fallos.

# ESCALABILIDAD



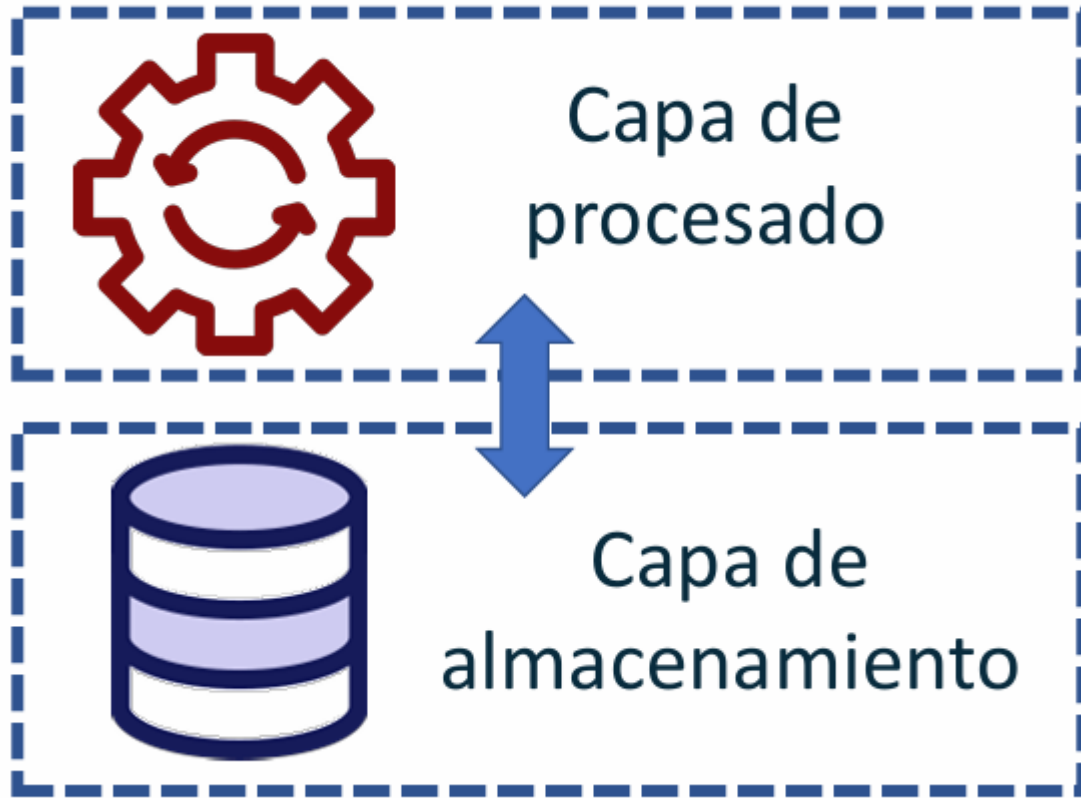
- A la hora de diseñar un servicio en la nube, hay que predecir que la demanda puede crecer (o decrecer).
- **Escalabilidad:** Añadir recursos a la infraestructura.



**ESCALABILIDAD HORIZONTAL**



**ESCALABILIDAD VERTICAL**



# ELEMENTOS CLAVE



## NameNode

1. Es el master en HDFS.
2. Distribuye los datos en bloques entre los diferentes DataNodes
3. Almacena los metadatos para saber dónde están los archivos en cada momento, qué DataNodes están disponibles etc.
4. Es el encargado de crear, renombrar o eliminar archivos o carpetas etc.

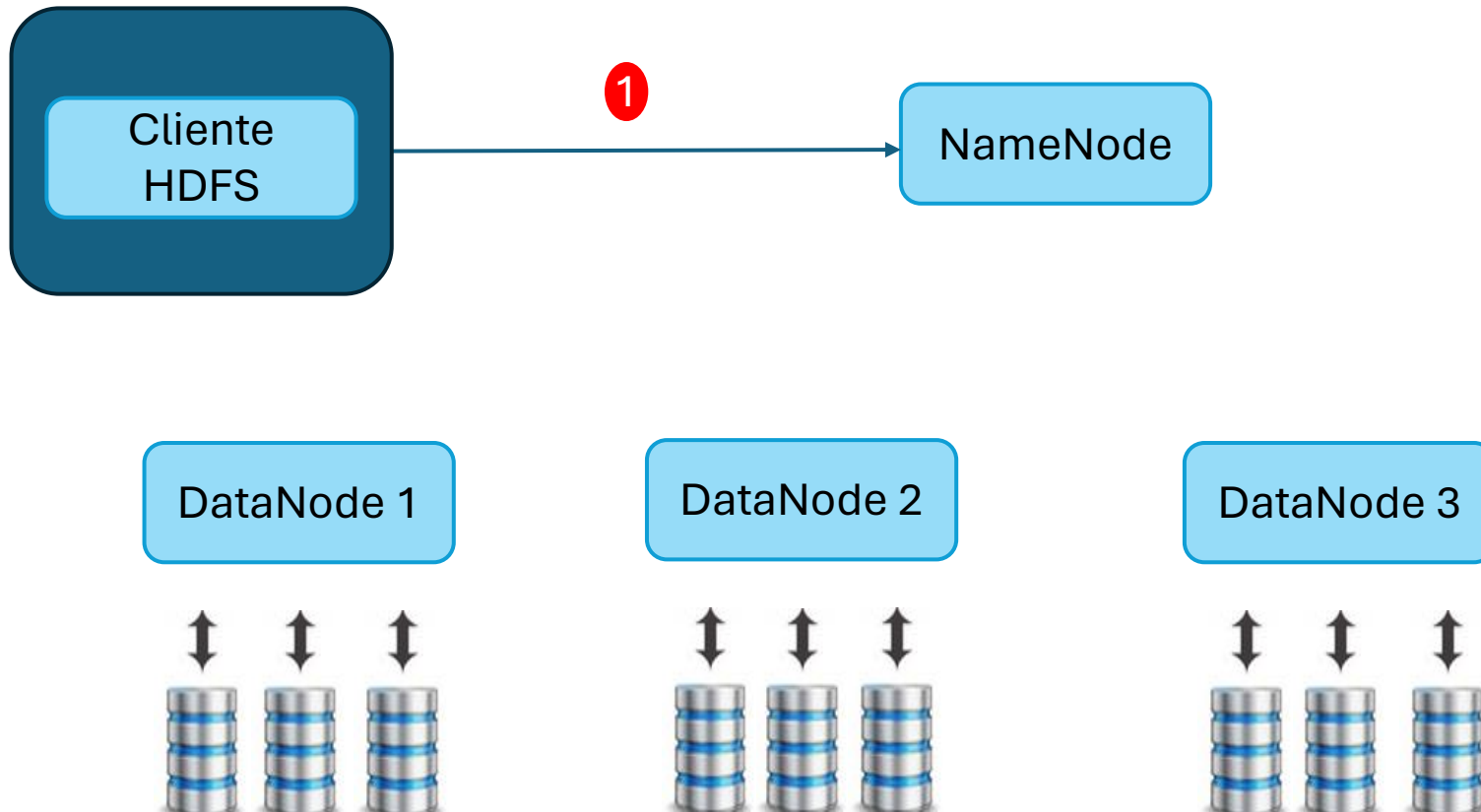
## DataNode

1. Es el esclavo en HDFS.
2. Es el encargado de almacenar los datos y en el caso de que un cliente contacte enviarle los datos directamente a éste.
3. Se comunican periódicamente con el NameNode (heartbeat cada 3s), de manera que éste conoce su estado en todo momento.
4. Crea, borra o replica bloques cuando el NameNode se lo encarga.

---

**El NameNode es el servidor primario que maneja los metadatos y dirige a los clientes hacia los DataNodes. Los DataNodes se encargan de almacenar los datos**

# ¿CÓMO ESCRIBE HADOOP UN FICHERO?



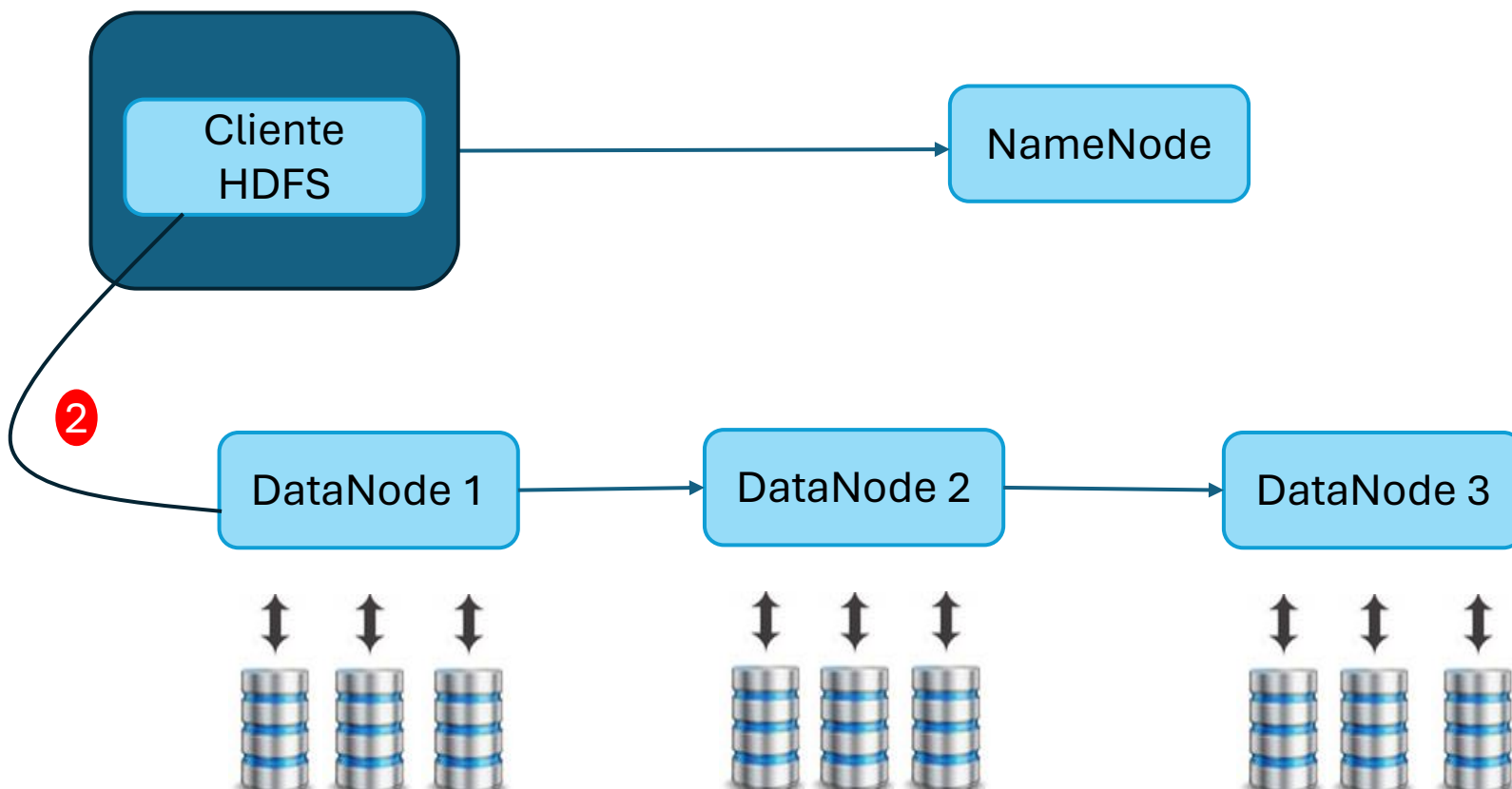
1

El cliente HDFS escribirá donde le indica el NameNode

El NameNode se asegura de que al menos una réplica se encuentre en otro DataNode



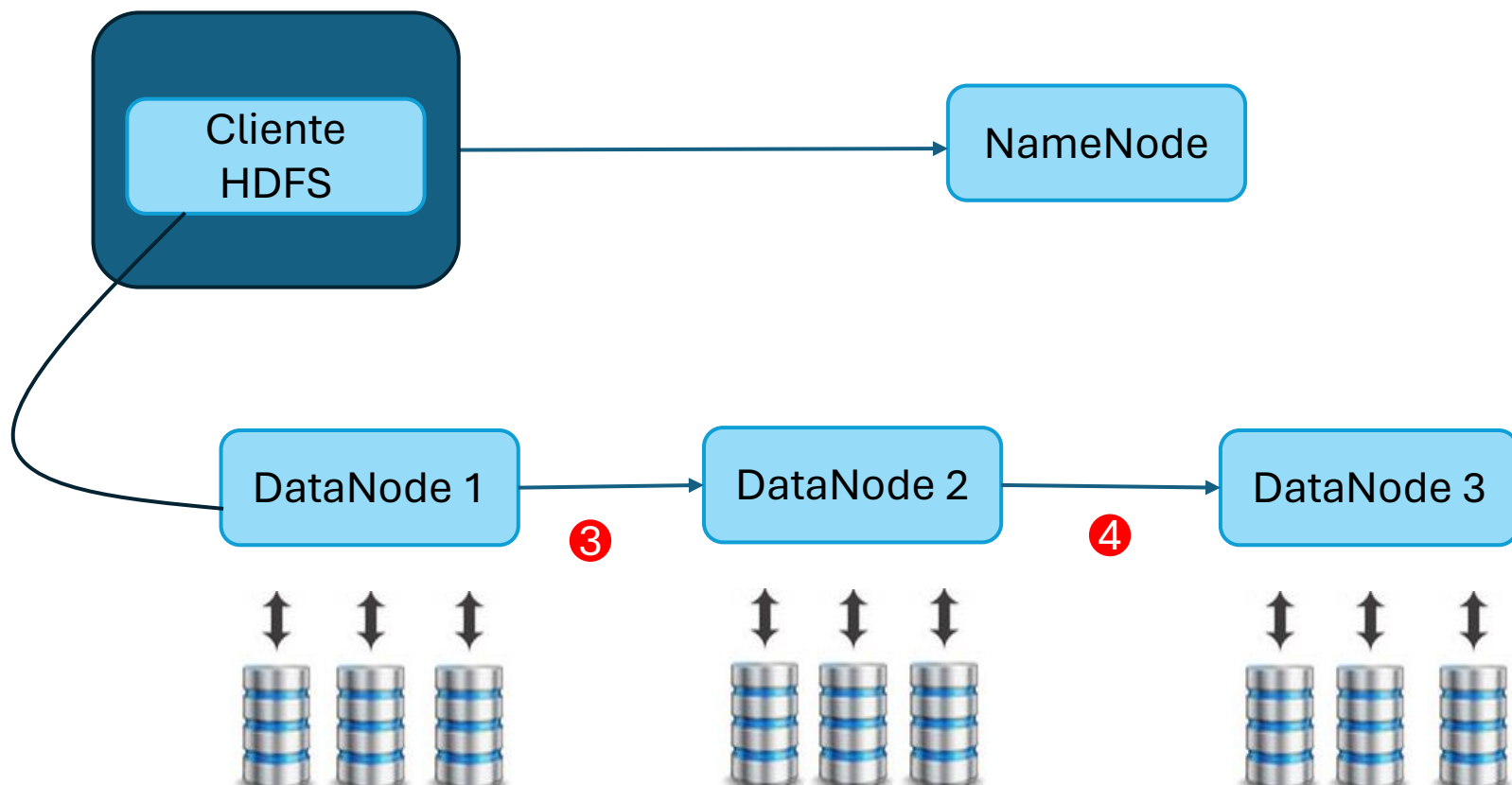
# ¿CÓMO ESCRIBE HADOOP UN FICHERO?



2

El cliente envía los datos al primer DataNode

# ¿CÓMO ESCRIBE HADOOP UN FICHERO?

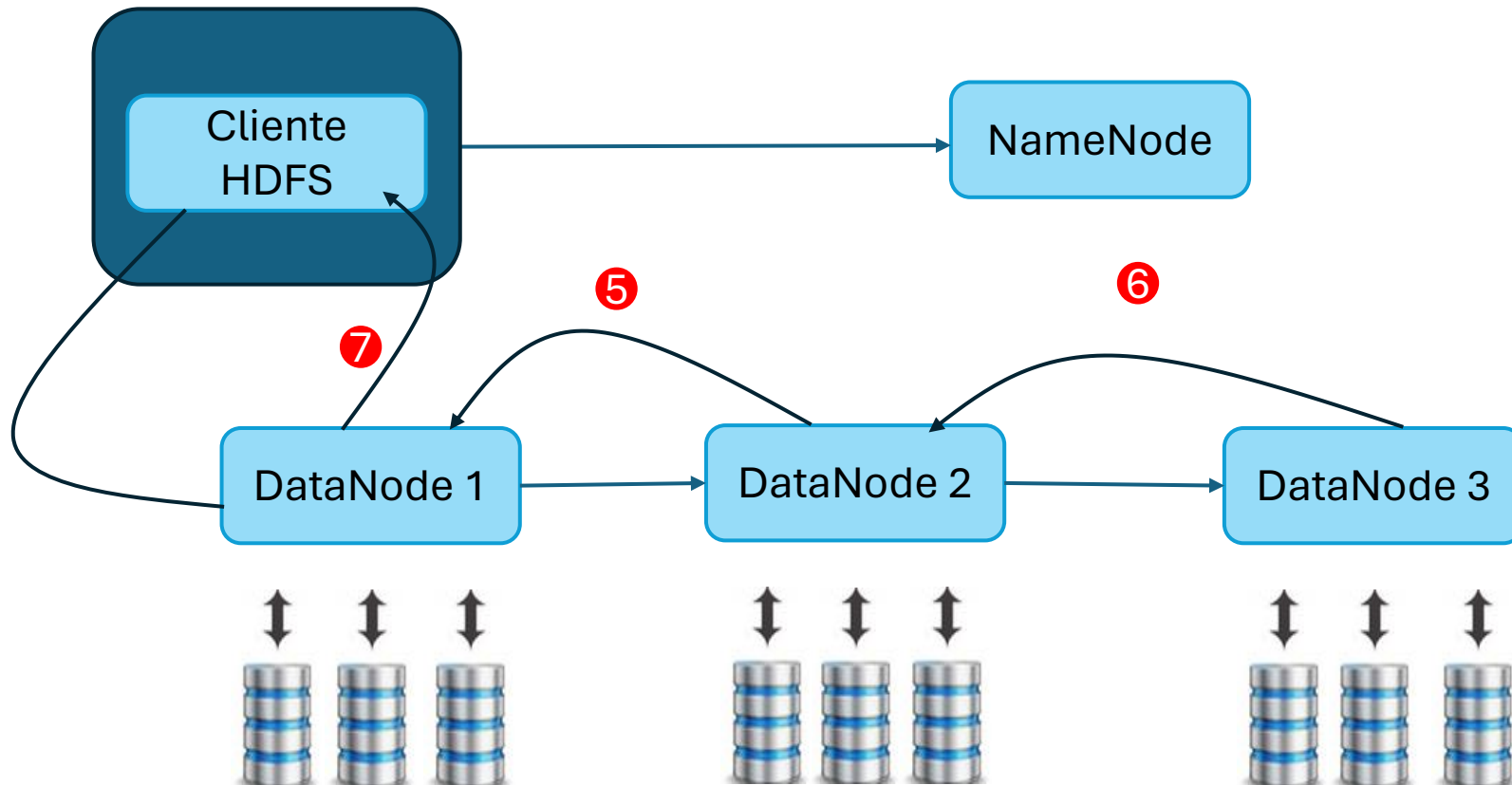


3

4

Se hace una copia de los datos.

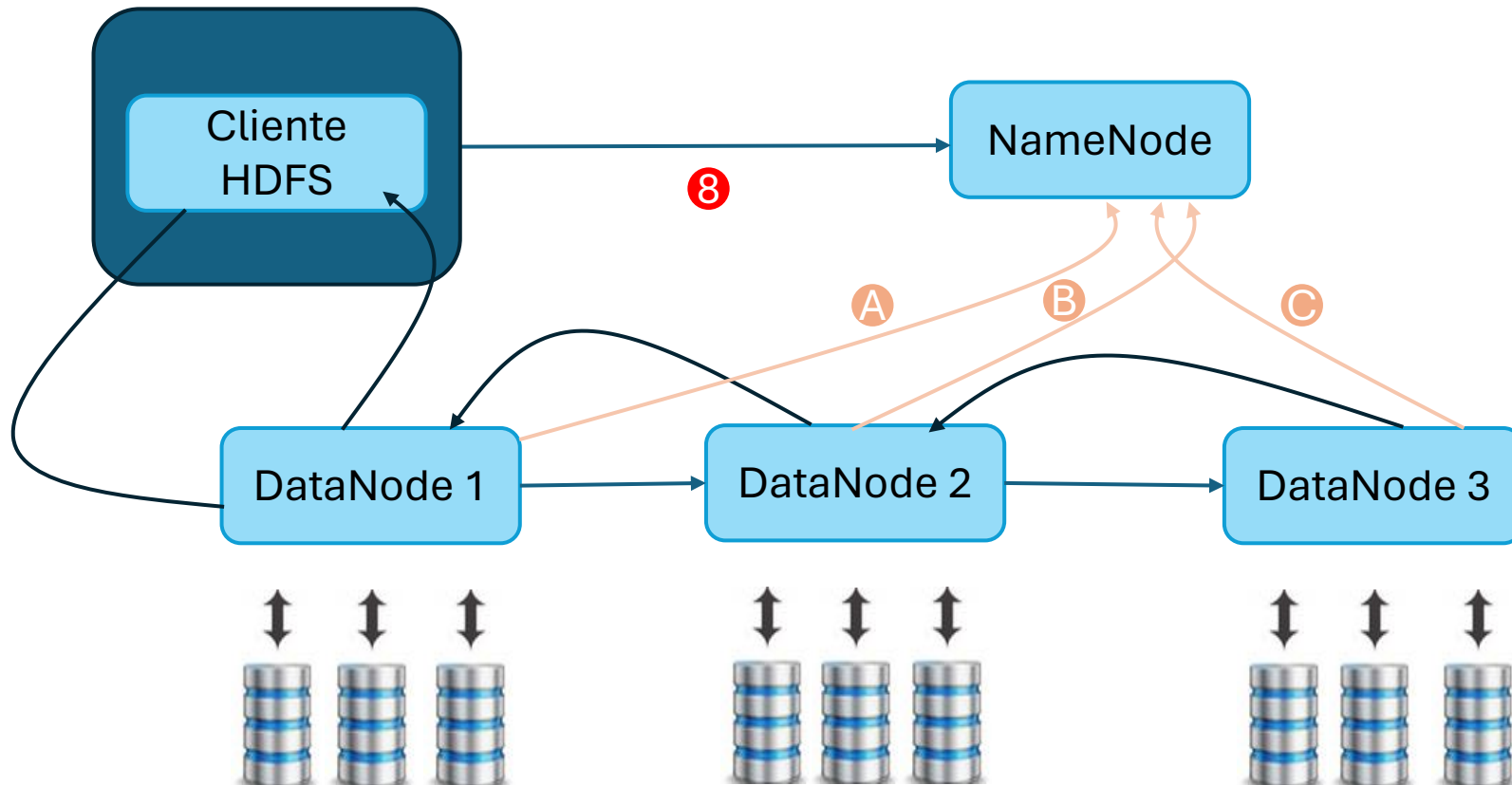
# ¿CÓMO ESCRIBE HADOOP UN FICHERO?



5 6 7

Se confirma las copias en sentido inverso.

# ¿CÓMO ESCRIBE HADOOP UN FICHERO?



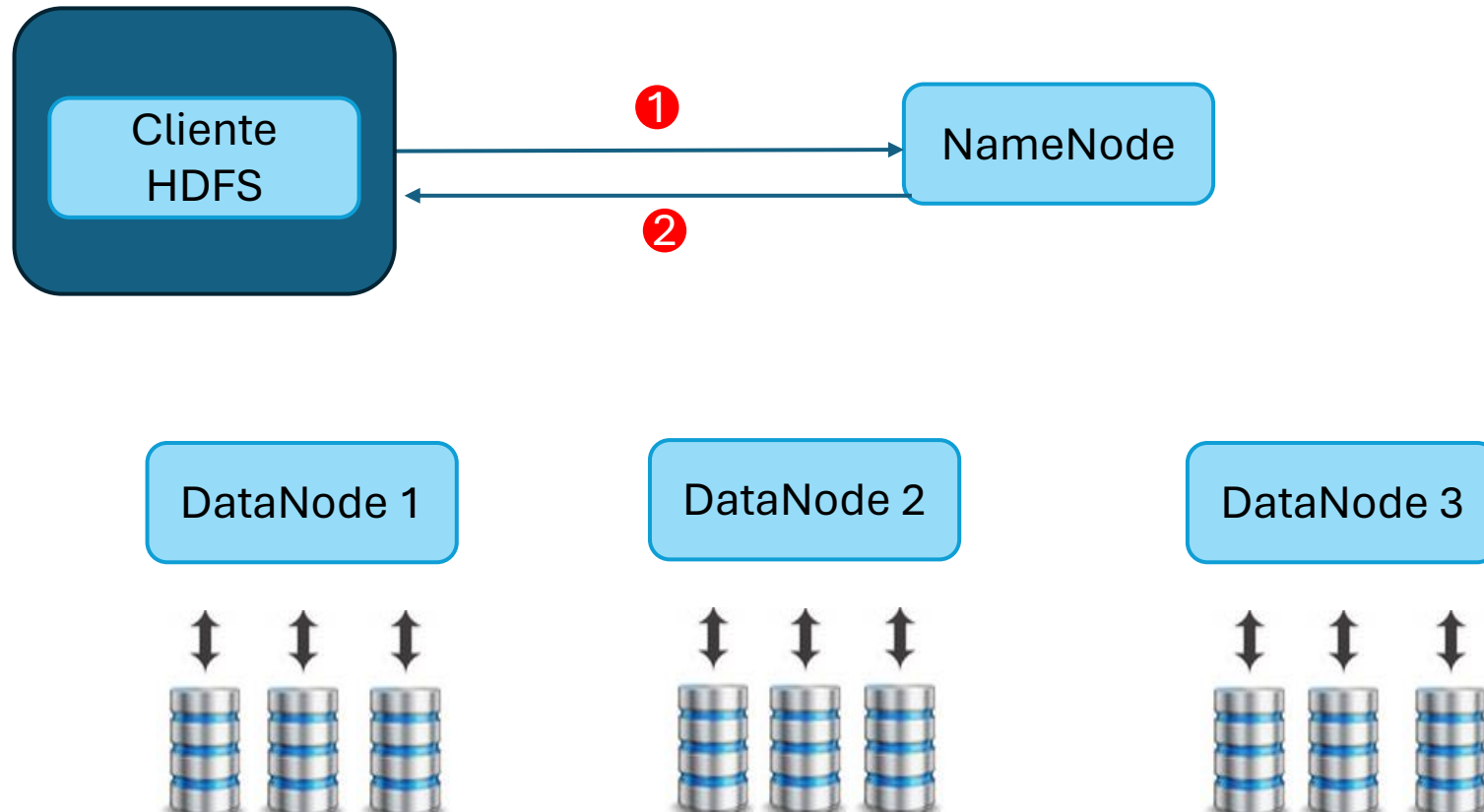
8

El cliente informa al NameNode que ha terminado.

A B C

Los DataNodes informan al NameNode de que ha terminado.

# ¿CÓMO LEE HADOOP UN FICHERO?

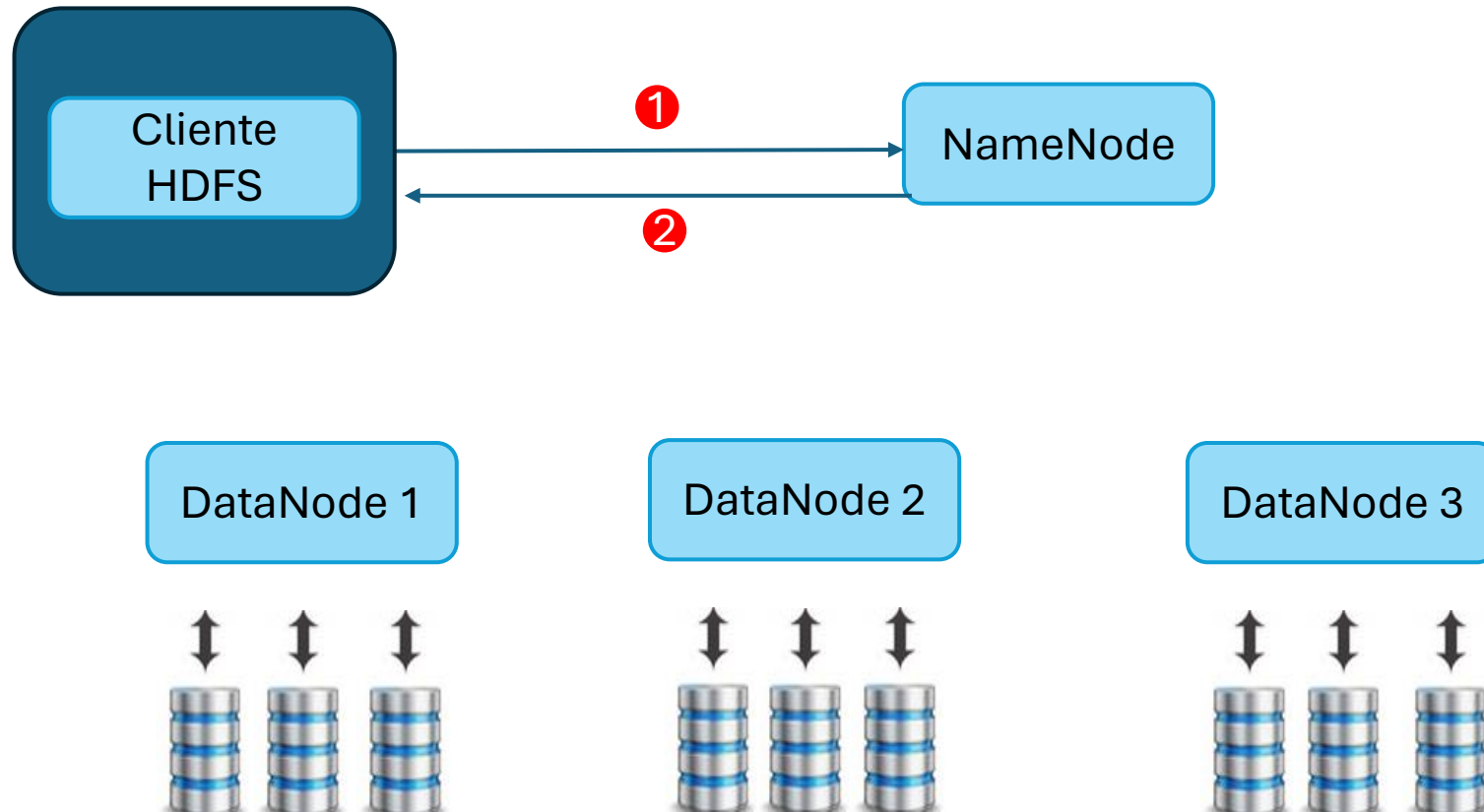


1

El cliente contacta con el NameNode para saber dónde están los bloques de las copias

Los DataNodes informan al NameNode de que ha terminado.

# ¿CÓMO LEE HADOOP UN FICHERO?



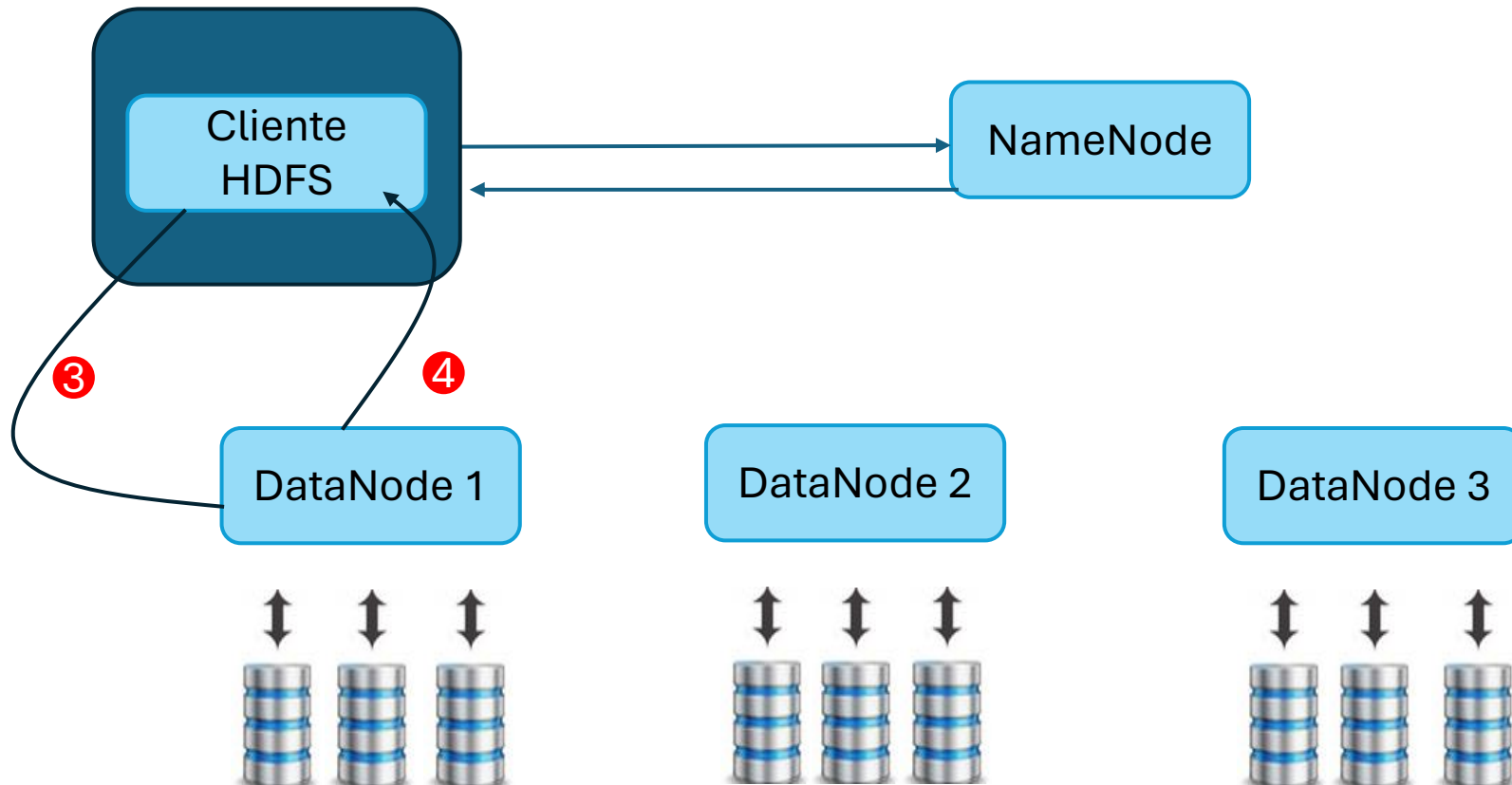
1

El cliente contacta con el NameNode para saber dónde están los bloques de las copias

Los DataNodes informan al NameNode de que ha terminado.



# ¿CÓMO LEE HADOOP UN FICHERO?



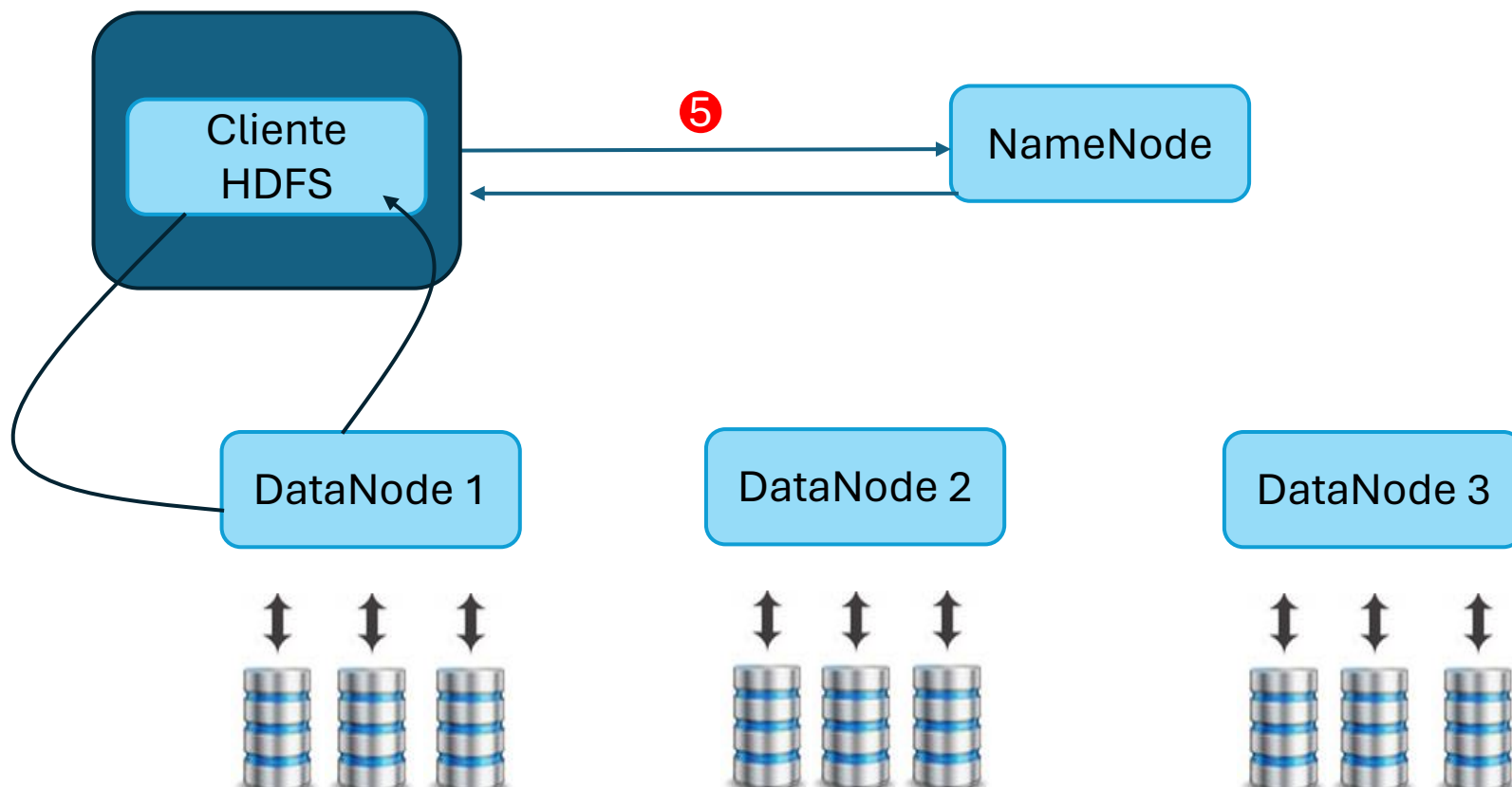
3

El cliente abre la comunicación con el DataNode correspondiente.

4

El DataNode responde con los datos

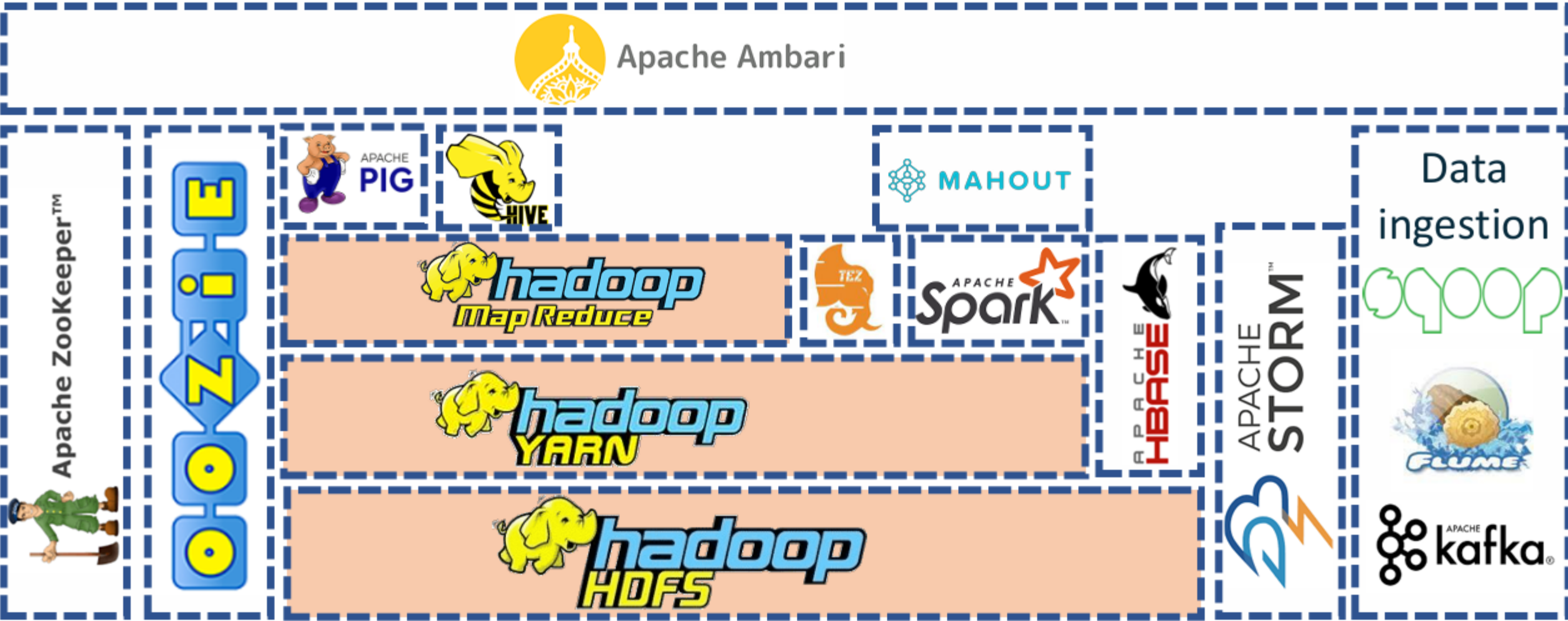
# ¿CÓMO LEE HADOOP UN FICHERO?



5

Se indica que se ha finalizado el proceso

# ECOSISTEMA DE HADOOP





**Spark es una plataforma de proceso de datos que utiliza una capa de memoria por encima del almacenamiento de datos desde el que los datos pueden ser cargados y procesados en paralelo a lo largo de todo el cluster**





Google Cloud



databricks



Microsoft Azure







Puede ser utilizado para procesar archivos en lote (batch processing) y procesamiento en streaming

Usado para procesamiento en lote

Más rápido por usar la memoria

Más lento por usar disco. Latencia de lectura/escritura

Tiene muchas APIs para soportar procesamiento de Big Data

Tiene un menor número de APIs

Utiliza RDD para la tolerancia a fallos

Utiliza replicación para la tolerancia a fallos

Utiliza menos líneas de código

Tiene niveles superiores de seguridad



**Microsoft Fabric es una plataforma de datos y análisis integral diseñada para empresas que requieren una solución unificada. Abarca el movimiento de datos, el procesamiento, la ingesta, la transformación, el enrutamiento de eventos en tiempo real y la creación de informes.**



Synapse Data  
Warehousing



Synapse Data  
Engineering



Data  
Factory



Synapse Data  
Science



Synapse Real  
Time Analytics



Power BI

T-SQL

Spark

Serverless  
compute

KQL

Analysis  
Services

OneSecurity

Warehouse

Lakehouse

Delta –  
Parquet  
Format

Delta –  
Parquet  
Format



OneLake

Kusto DB

Dataset

Delta –  
Parquet  
Format

Delta –  
Parquet  
Format

