

Curso Data Engineer: Creando un pipeline de datos

MÓDULO B - Clase 5

Agenda



- Jupyter Notebook
- GCP Ingest
- Transform
- Ejercicios



Jupyter Notebooks

Jupyter Notebook



127.0.0.1:8889/notebooks/Untitled.ipynb

jupyter Untitled Last Checkpoint: 3 minutes ago

File Edit View Run Kernel Settings Help

Trusted

JupyterLab Python 3 (ipykernel)

```
from pyspark.sql import SparkSession
spark = SparkSession.builder \
    .master("spark://localhost:7077") \
    .getOrCreate()

df = spark.read.option("header", "true").csv("/ingest/yellow_tripdata_2021-01.csv")

[3]: df.show(5)
```

	VendorID	tpep_pickup_datetime	tpep_dropoff_datetime	passenger_count	trip_distance	RatecodeID	store_and_fwd_flag	PULocationID	DOLocationID	payment_type	fare_amount	extra	mta_tax	tip_amount	tolls_amount	improvement_surcharge	total_amount	congestion_surcharge
8	3	2021-01-01 00:30:10	2021-01-01 00:36:12	0	1	2.10	1		142	43	2							
		0.5	0	0.3	11.8	.20	1		238	151	2							
3	0.5	2021-01-01 00:51:20	2021-01-01 00:52:19	0	0.3	4.3	0		132	165	1							
		0.5	0	0.3	1	14.70	1		138	132	1							
42	0.5	2021-01-01 00:43:30	2021-01-01 01:11:06	0	0.3	51.95	0		138	132	1							
		0.5	0	0.3	0	10.60	1		138	132	1							
29	0.5	2021-01-01 00:15:48	2021-01-01 00:31:01	0	0.3	36.35	0		68	33	1							
		0.5	0	0.3	1	4.94	1		68	33	1							
16.5	0.5	2021-01-01 00:31:49	2021-01-01 00:48:21	0	0.3	24.36	1		68	33	1							
		0.5	0	0.3	1	2.5												

only showing top 5 rows

A decorative graphic on the left side of the slide, consisting of a grid of small dots in various shades of blue and green, arranged in a pattern that tapers off to the right.

GCP Ingest

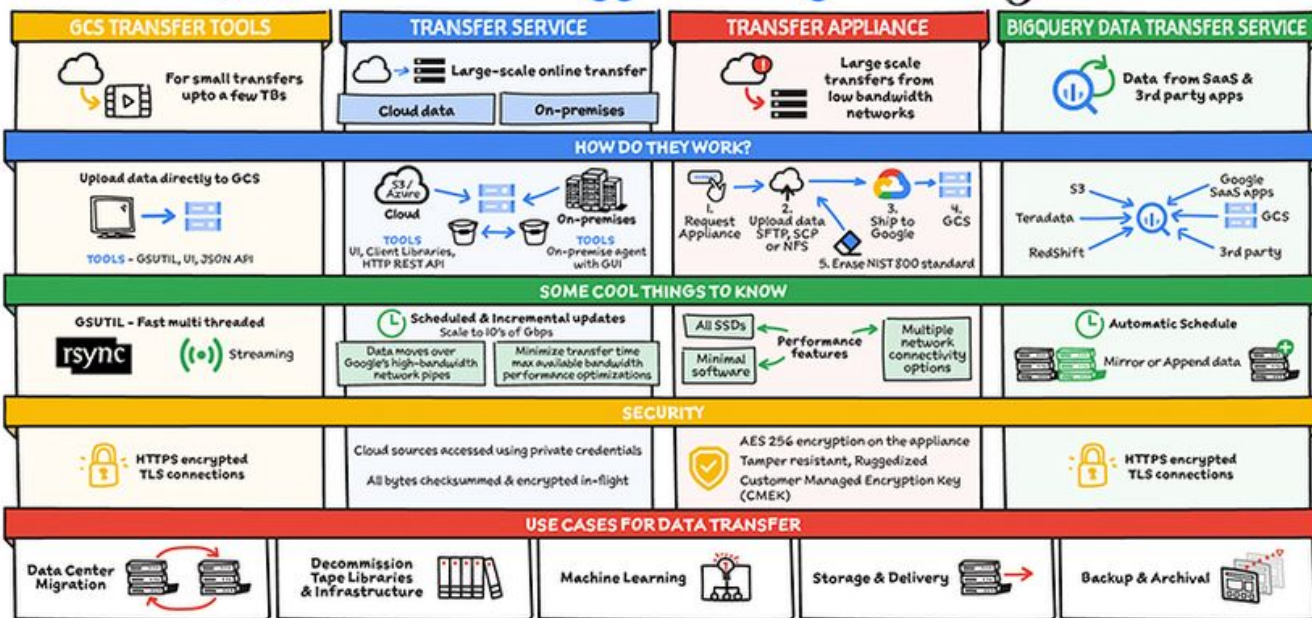
Ingest



#GCPSSketchnote
@PYERGADIA
THECLOUDGIRL.DEV
03.30.2021



Options to move data to Google Cloud



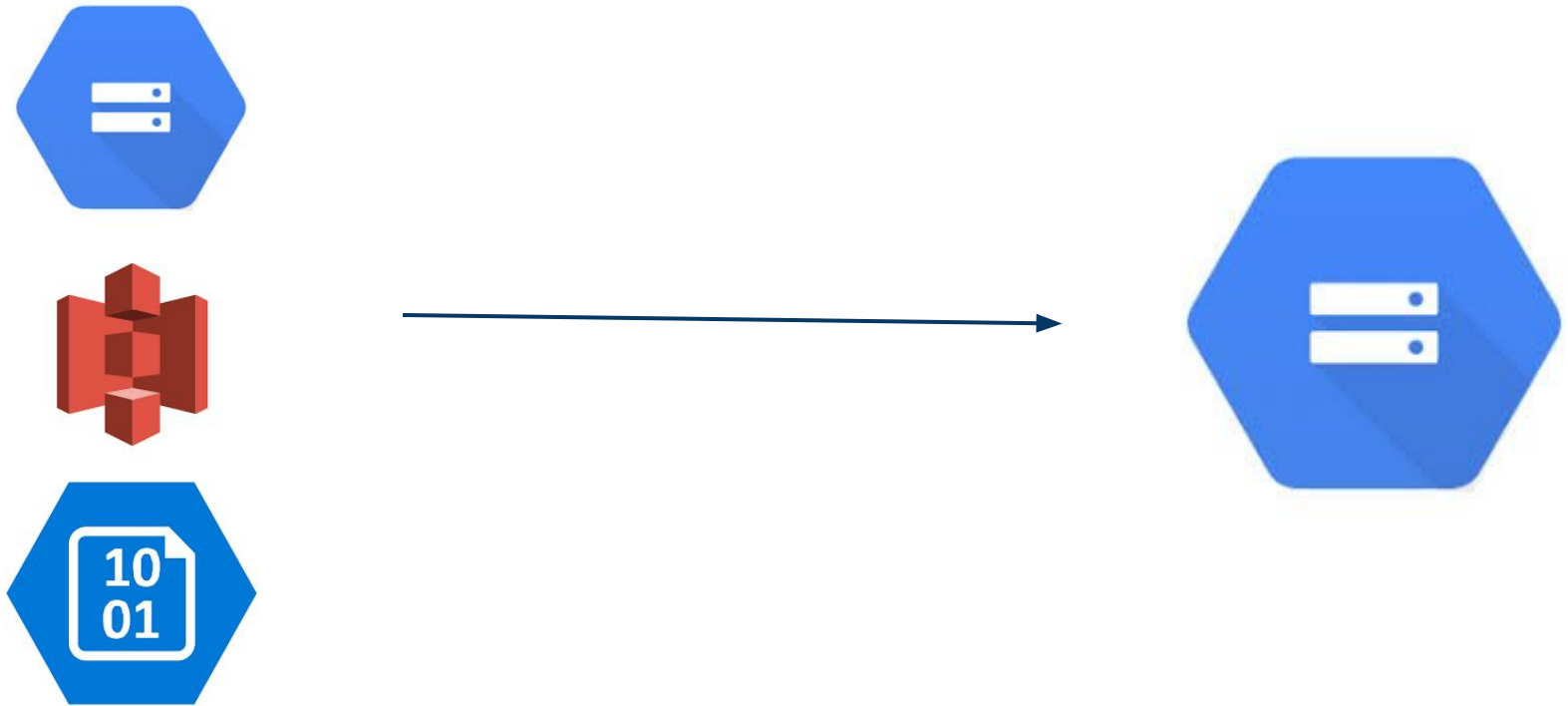
gsutil



 Parquet



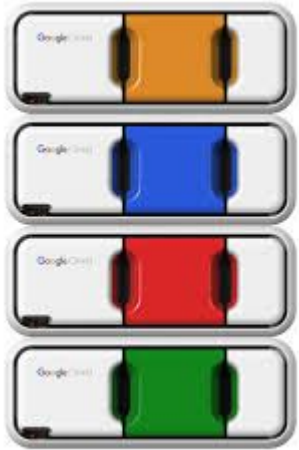
transfer service - data transfer



transfer appliance



transfer appliance



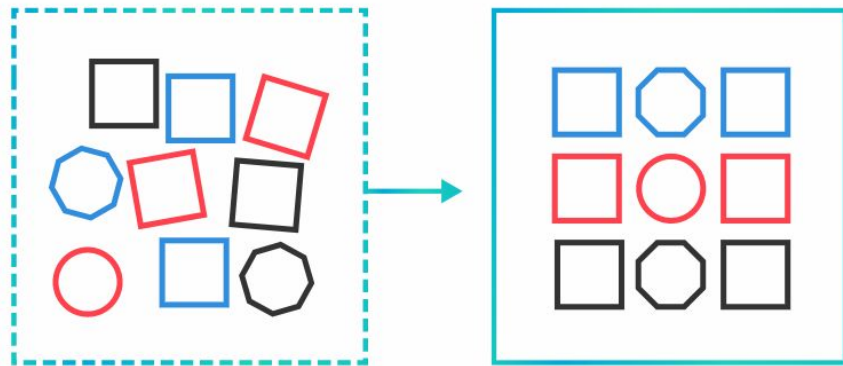


Transform

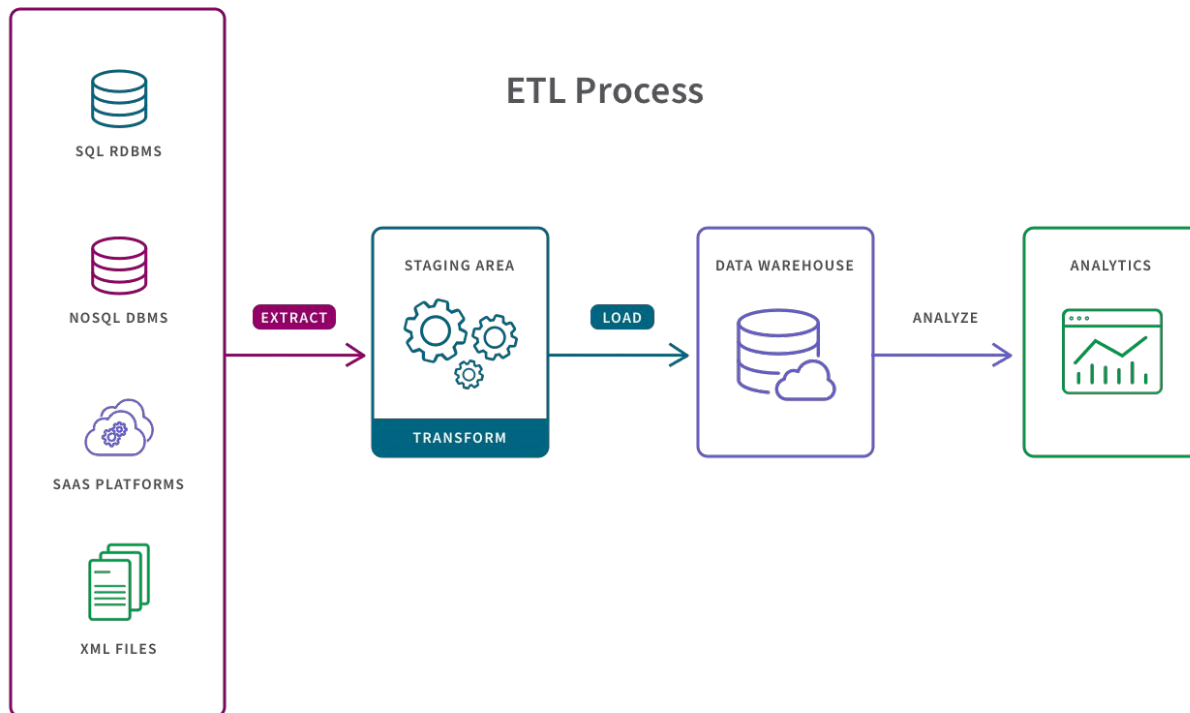
Transform



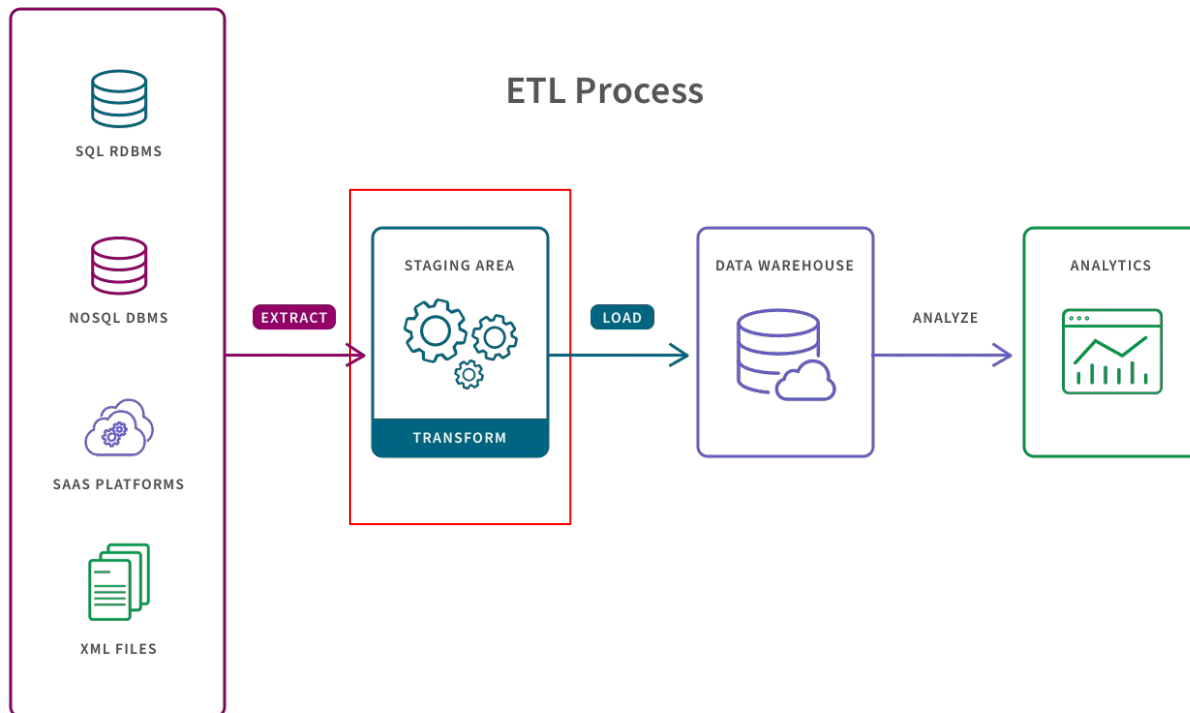
Es el proceso de convertir, limpiar y estructurar datos en un formato utilizable que se pueda analizar para respaldar los procesos de toma de decisiones.



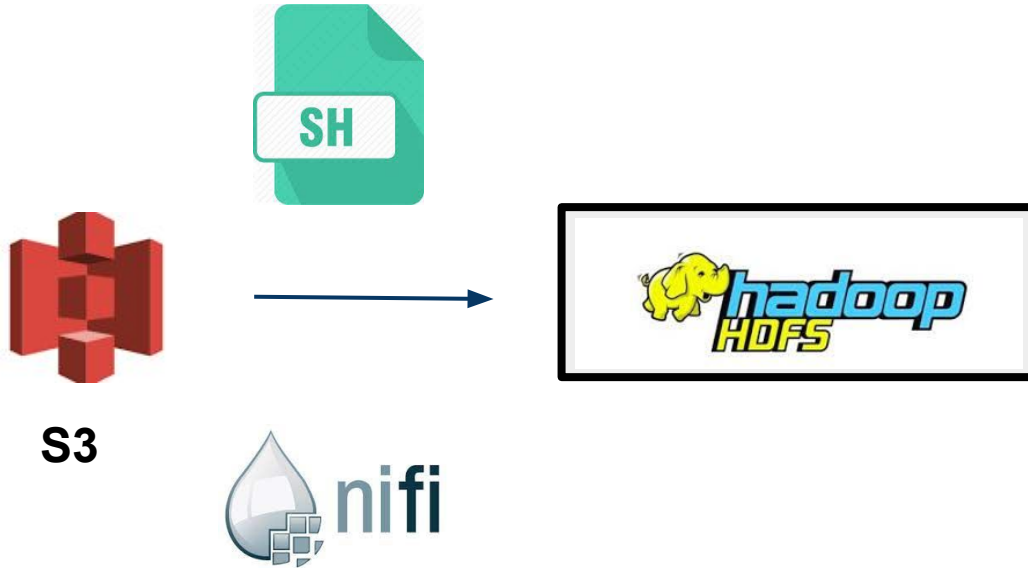
Transform



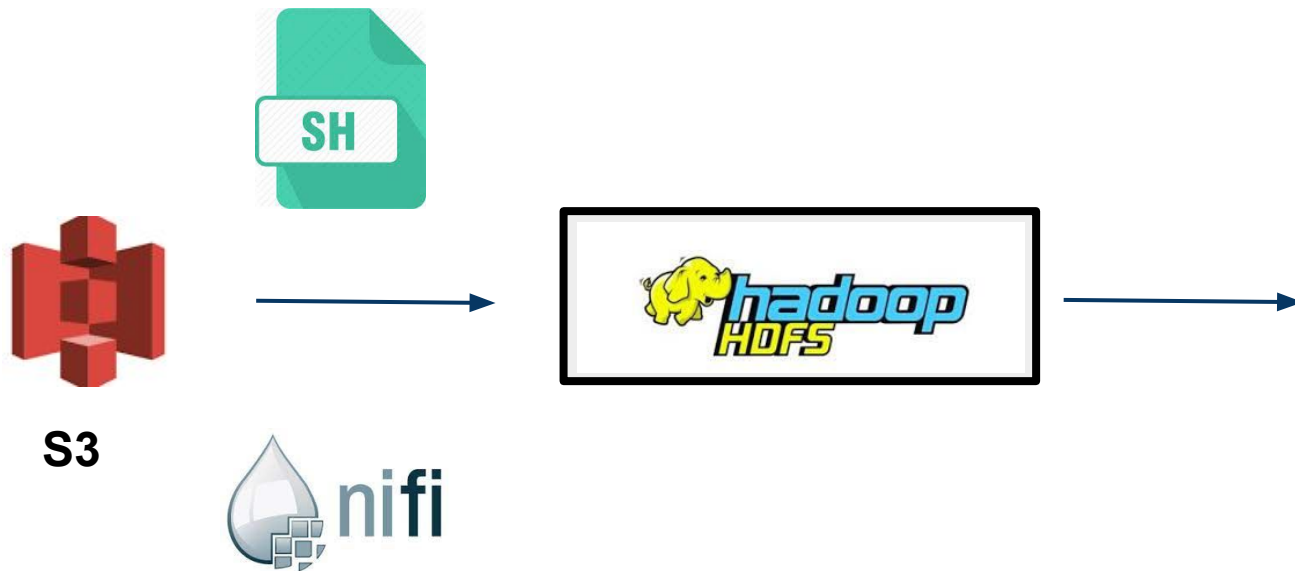
Transform



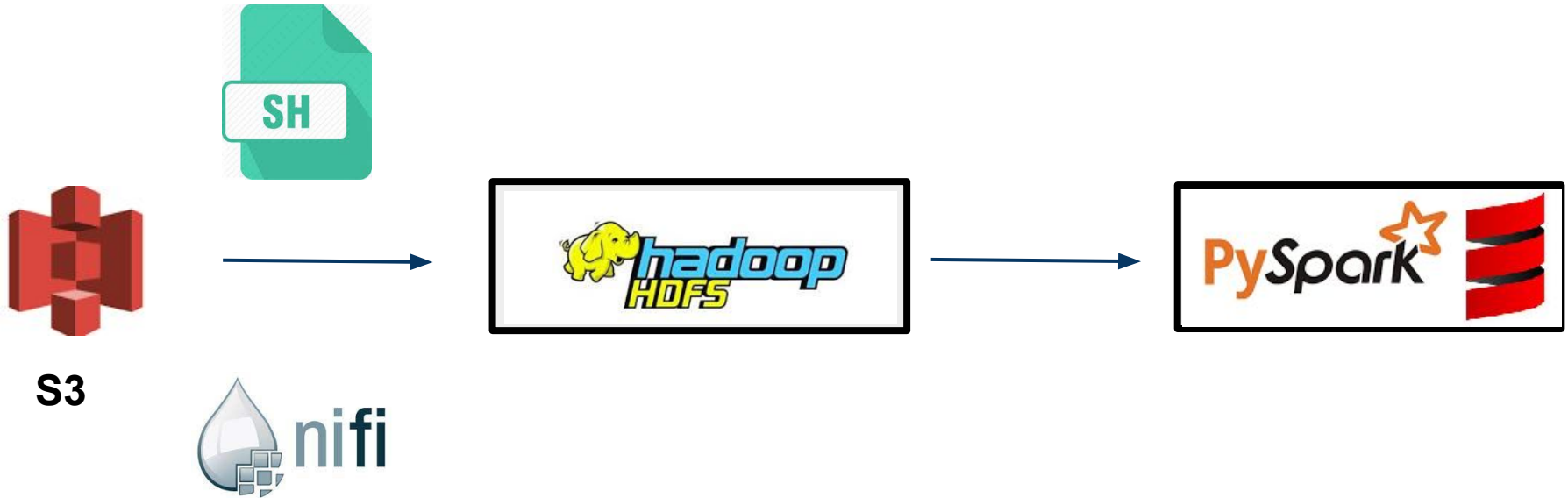
Transform



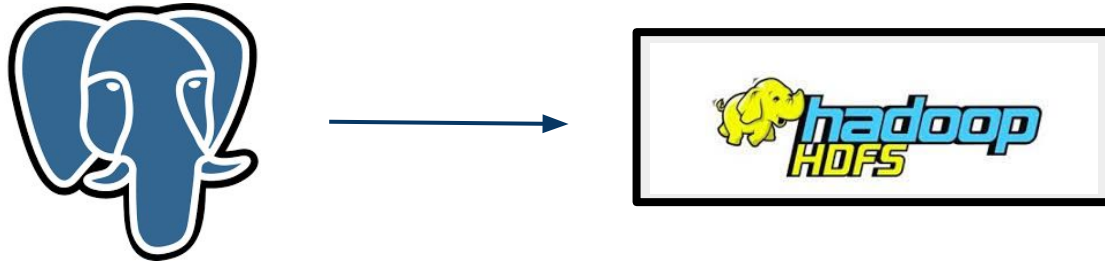
Transform



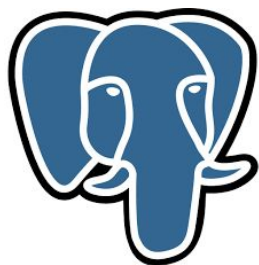
Transform



Transform



Transform



Transform



ETL VS ELT



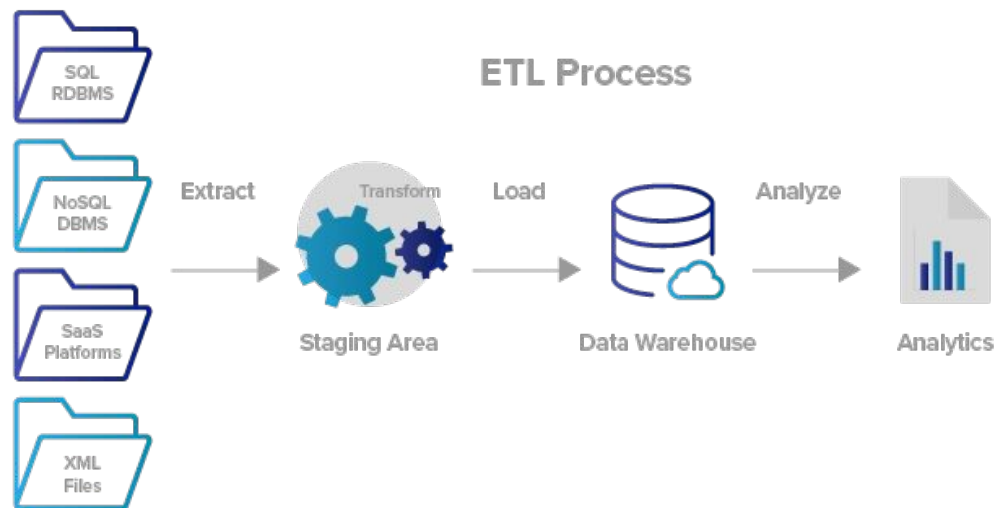
ETL



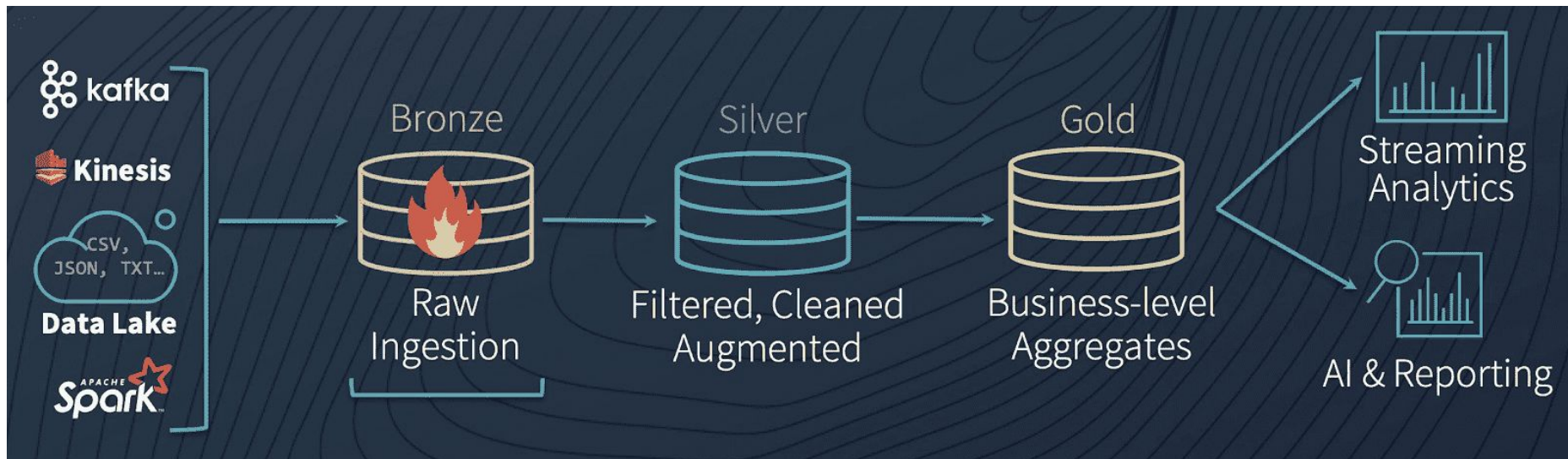
ELT



ETL



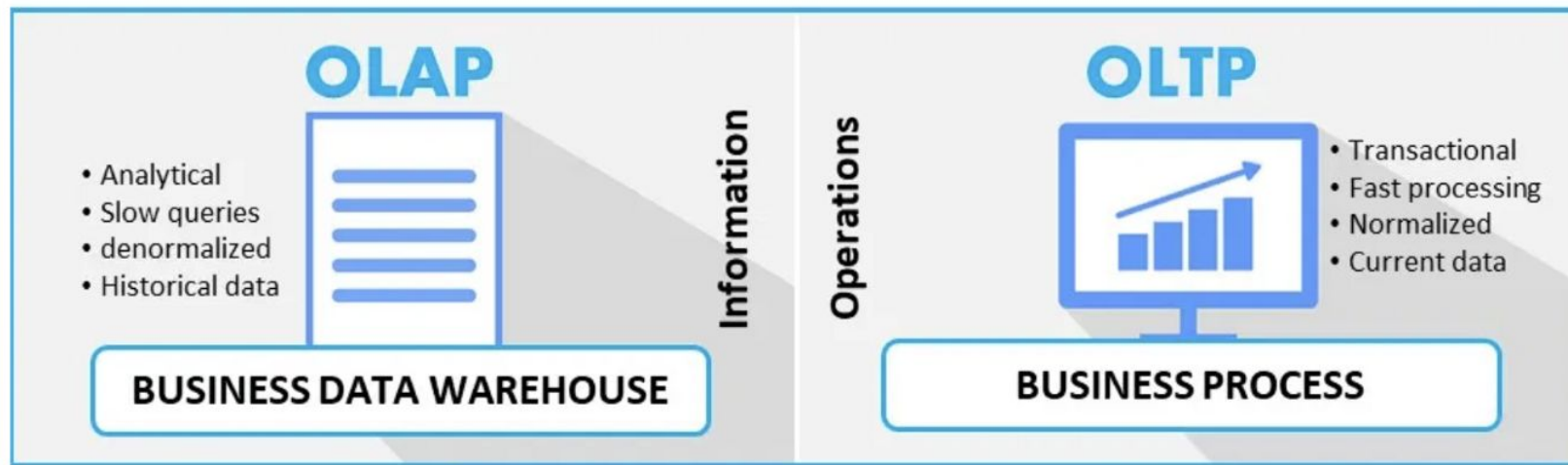
ELT



Online Analytical Processing Vs Online Transaction Processing



OLAP Vs OLTP



On line Analytical Processing Vs On line Transaction Processing



Students	
id	student_name
1	Juan García
2	José Perez
3	Alberto Quiroga

Courses	
id	course_name
1	SQL
2	Python
3	R
4	Java

Students_Courses			
Date	id_student	id_courses	mark
11/01/2022	1	1	7
10/25/2022	1	2	6
10/28/2022	2	4	6
10/03/2022	3	3	5

ETL



1. Con Sqoop hago un export en archivos parquet/text/avro, etc.
2. Ingesto esos archivos en HDFS (hdfs dfs -put origen destino)
3. Transformo esos archivos con pyspark/scala/sql
4. Cargo esos datos en el Data Warehouse (Hive)

OLAP

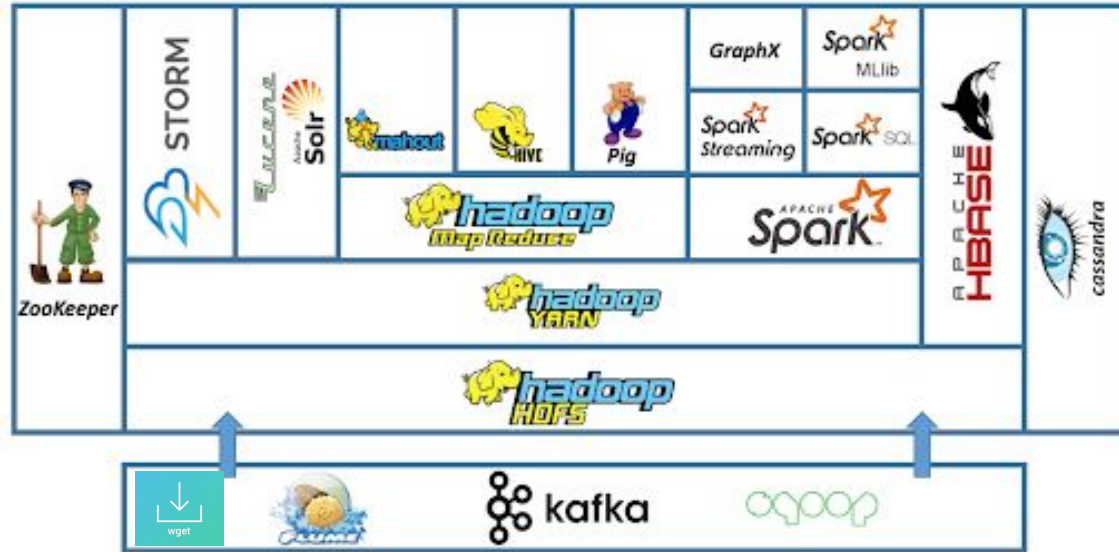


Marks				
Date	student	courses	mark	assist
11/01/2022	Juan García	SQL	7	100
10/25/2022	Juan García	Python	6	93
10/28/2022	José Perez	Java	6	95
10/03/2022	Alberto Quiroga	R	5	80



Ejercicio

Ecosistema Hadoop



Ejercicios



- GCP
 - Creación de cuenta y configuración de consola
 - Ingesta de archivos con Gsutil
 - Storage transfer