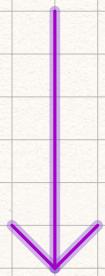


## # Data Warehouses.

- Poor fit, with mixed record of performance
- Still useful in many use cases.



## Following a Data Lake.

- The abundance of unstructured data (text, xml, json, logs, audio files, images, video...)
- Unpredicted data columns (social, IoT, machine-generated, etc.)
- The rise of Big Data technologies like HDFS, Spark, etc.
- New types of data analysis gaining momentum: graph analytics, predictive analytics...
- Emergence of new role Data Scientists.

## # Abundance of Unstructured Data:

- might be possible in the ETL process. Distill some elements of file and put them in tabular form
- But ↗■ Later we might decide we want it to be transformed differently → a particular transformation procedure → *Impose a strong commitment without enough knowledge.*



- Some data is hard to transform into tabular format → day marked XML
- Some data can be stored as "Blobs" → Remain swollen if not processed.



- HDFS reduces cost per TB compared to MPP databases
- Processing Frameworks work on the same mechanism employed by data lake.
- Possibility of performing analysis without requiring *previous schema definition*.
- Being data a new source of value; unstructured format applies no "pre-conception" to it

## # Big Data Technology Effects: On Warehouses

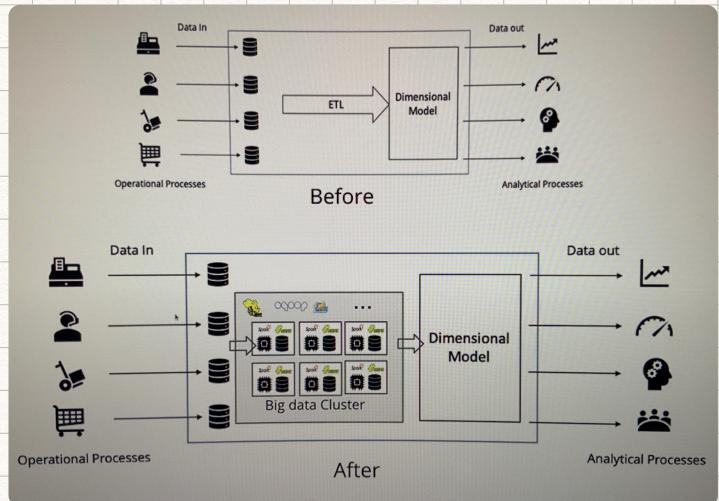
### # Lesser Cost & ETL Offloading

- ETL Offloading; not anymore a required step for data ingestion

- Same hardware for storage and processing  
→ no need for a special ETL box.

- Dimensional modelling for high known value data

- Low cost per TB → Room for storing  
= "Less known values data" (previously not available)



## # Schema-on-Read

- Traditionally, data in databases has been much easier to process than data in plain files

- Big Data tools from Hadoop ecosystem made it easy to work with unstructured data

↳ Creating a database

↳ Inserting data into the database

↳ Schema-on-read ↗ Inferred

↳ Specific! → when data is read, it is checked against the specified schema

## # Ch. 5: Structured Storage

- Read and Write: « files in many formats » variety of file system
  - a variety of schemas (SQL / No-SQL)
- All expand on a single common abstraction layer.

## # Data Lakes, maybe?

Before (Data Lake): data was extracted, then transformed in specific formats and finally loaded → ETL



Now (Data Lake): data is stored = as-is; transformation on demand.

Extract - Load - Transform

- ④ Bonus feature:
  - Possible parallelism & scalability comes out of the box
  - Column Storage using Parquet

## # Data Lake vs Data Warehouse:

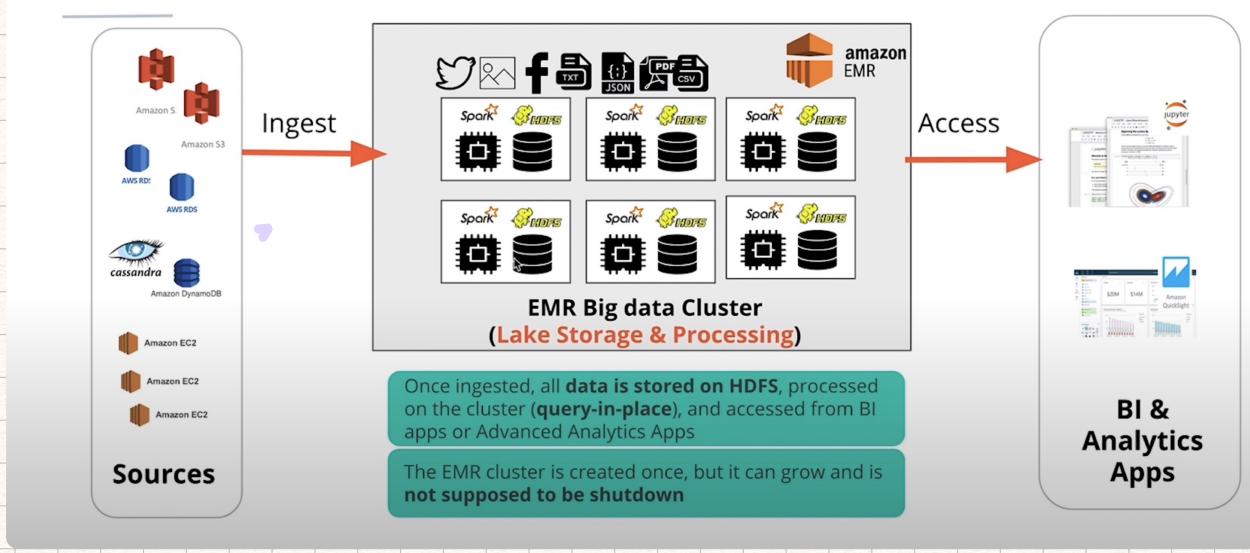
### Data Lake vs Data Warehouse

	Data Warehouse	Data Lake
Data form	Tabular format	All formats
Data value	High only	High-value, medium-value and to-be-discovered
Ingestion	ETL	ELT
Data model	Star & snowflake with conformed dimensions or data-marts and OLAP cubes	Star, snowflakes and OLAP are also possible but other ad-hoc representations are possible
Schema	Known before ingestion (schema-on-write)	On-the-fly at the time of analysis (schema-on-read)
Technology	Expensive MPP databases with expensive disks and connectivity	Commodity hardware with parallelism as first principle
Data Quality	High with effort for consistency and clear rules for accessibility	Mixed, some data remain in raw format, some data is transformed to higher quality
Users	Business analysts	Data scientists, Business analysts & ML engineers
Analytics	Reports and Business Intelligence visualizations	Machine Learning, graph analytics and data exploration

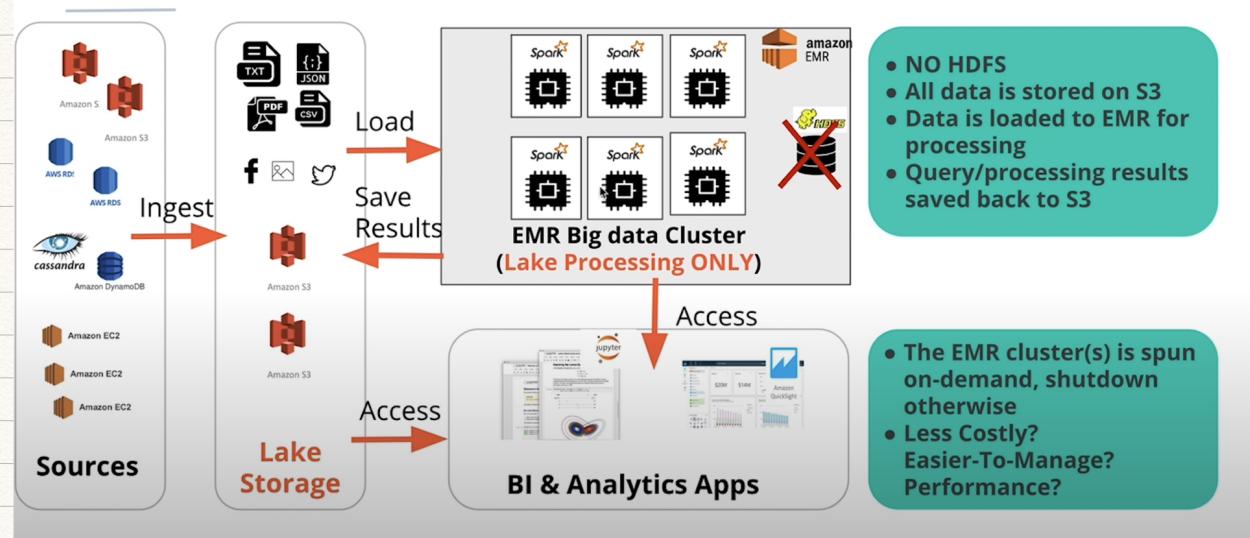
## # Data Lake options on AWS:

<del>HDFS</del>	Spark	AWS EMR (HDFS+Spark)	ECC + Vendor Solution
<del>S3</del>	Spark	AWS EMR (Spark)	ECC + Vendor Solution
<del>S3</del>	Serverless	AWS Athena	Lambda + Vendor Solution
= 'Storage'	'Processing'	'AWS-Managed Solution'	'Vendor Managed'

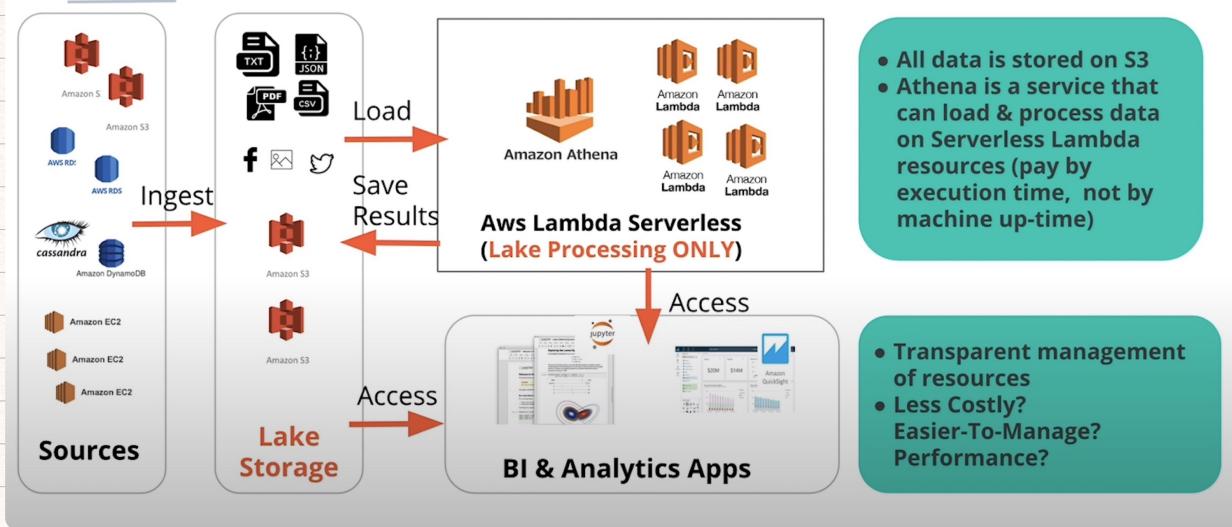
### Data Lake options: AWS EMR (HDFS+Spark)



### Data Lake options: AWS EMR (S3+Spark)



## Data Lake options: AWS Athena



### Data Lake Issues:

- Data Lakes are prone to be a *broader garbage dump*
  - Broader like detailed metadata to reduce risk.
- Since data lakes allow wide availability of cross-department data and external data.
  - Sometimes data governance is not easy to implement.
- Sometimes hard to know if a data lake should replace a data warehouse.