

# DFS and File Formats

# Agenda



- Changes in the exercise hours
- Discord!
- Some terms
- Hadoop's HDFS
- Parquet
- Avro

# Discord!

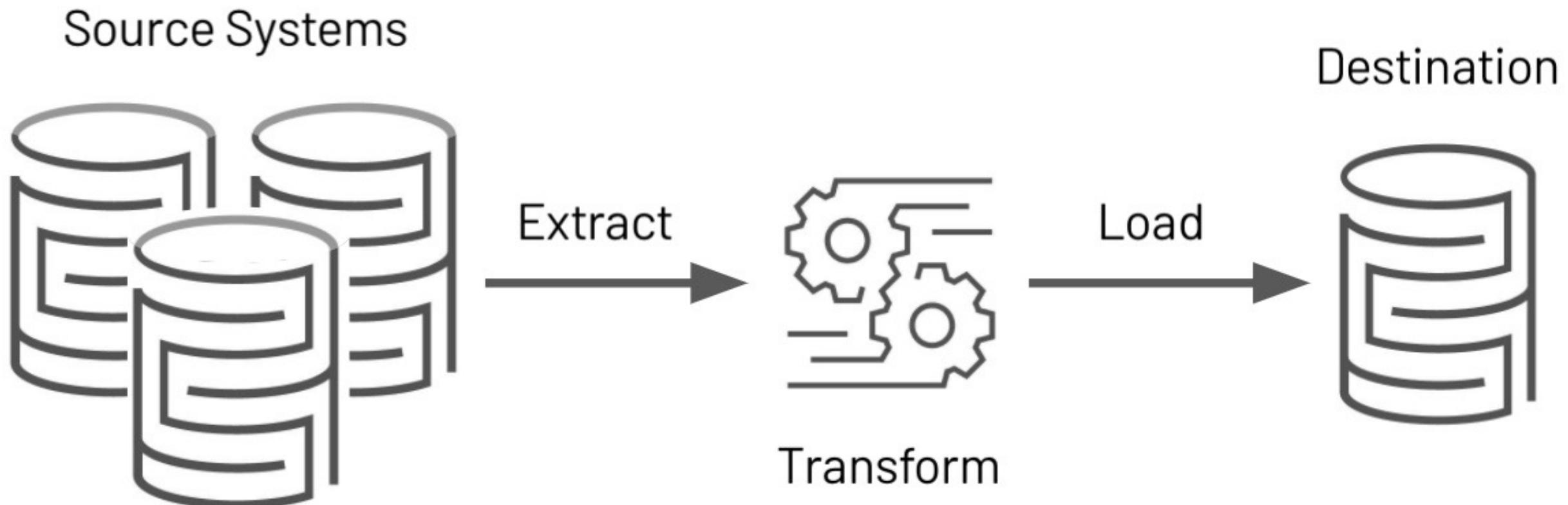


**DISCORD**

- Our common chatroom!
- Announcements!
- Ask the teacher, instructors, alumni, or other students!
- Join here: <https://discord.gg/mHvjRdn28>
- Change your name to your real name to get access!

# First some terms

# ETL vs ELT vs EtLT

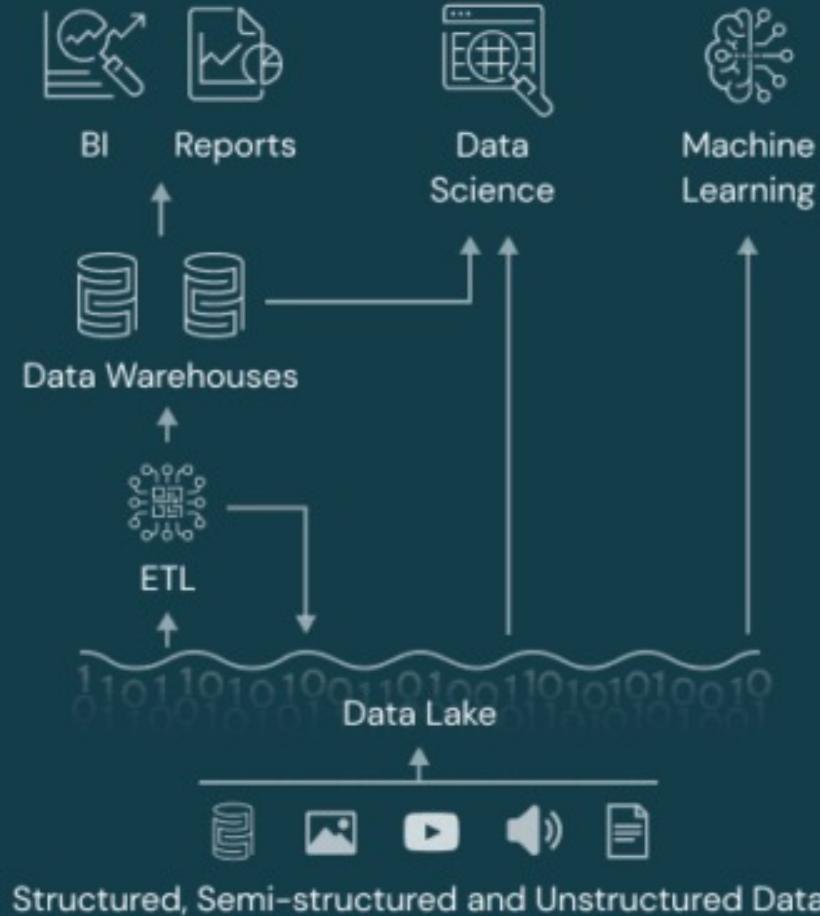


# Data Warehouse/Lake/Lakehouse (and Data Mesh)

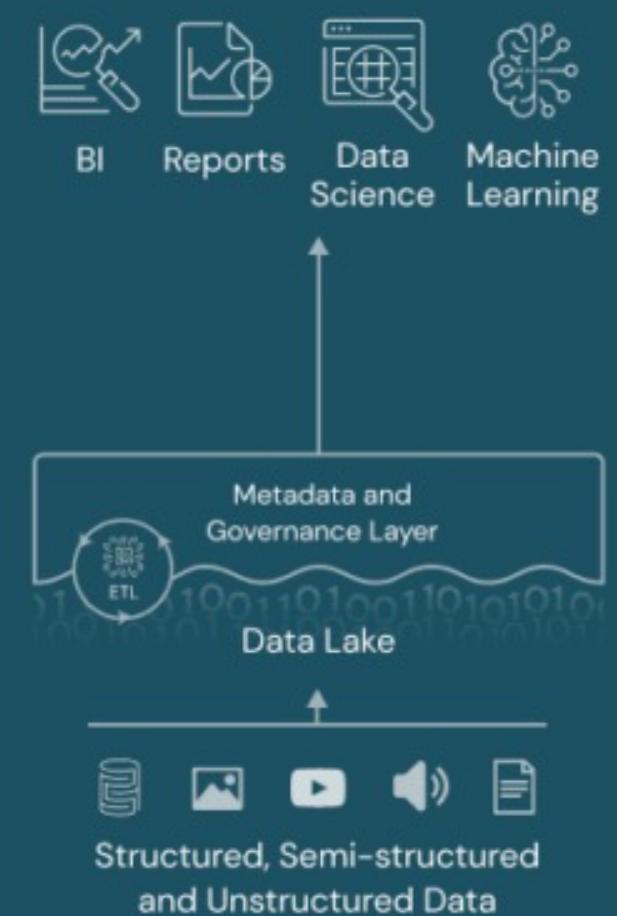
## Data Warehouse



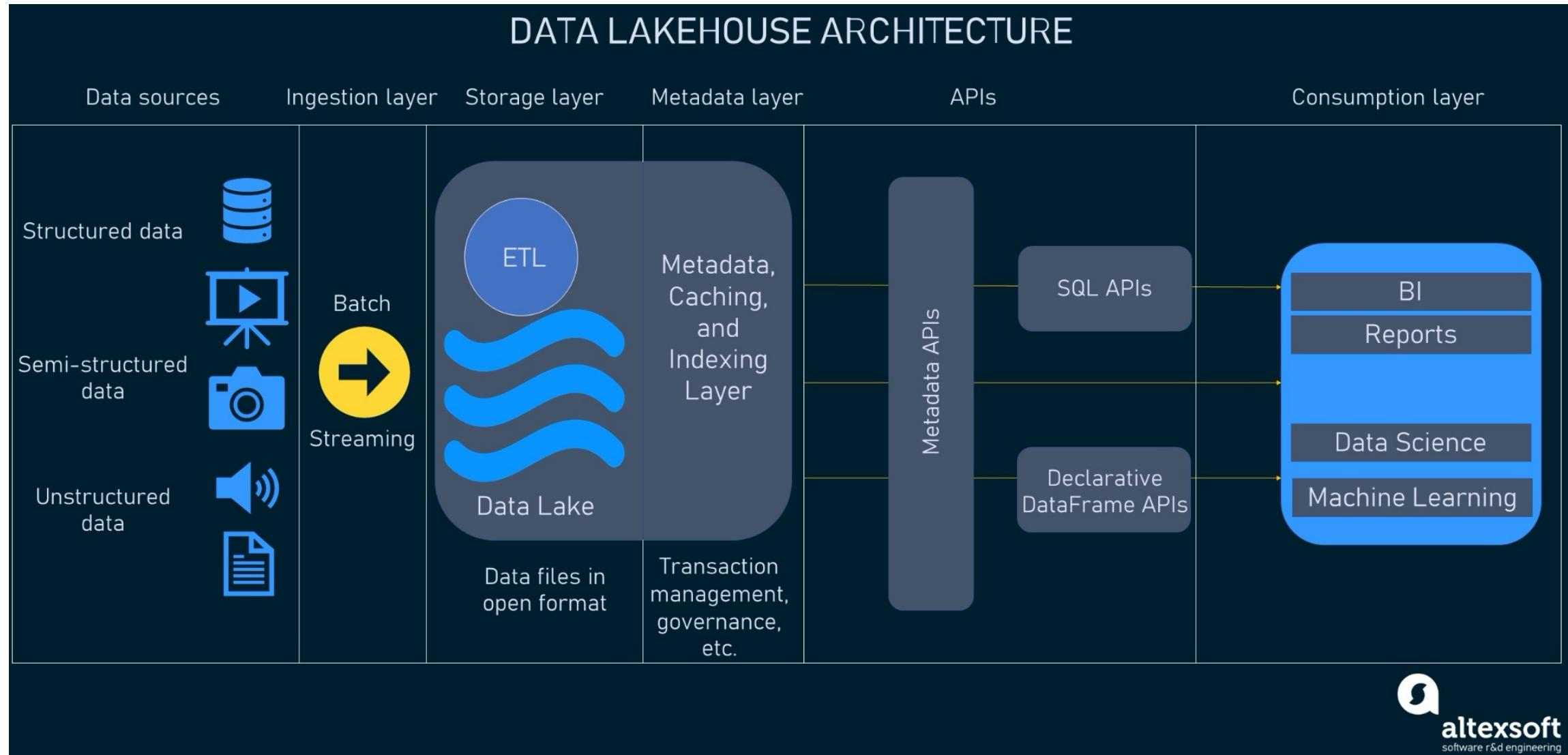
## Data Lake



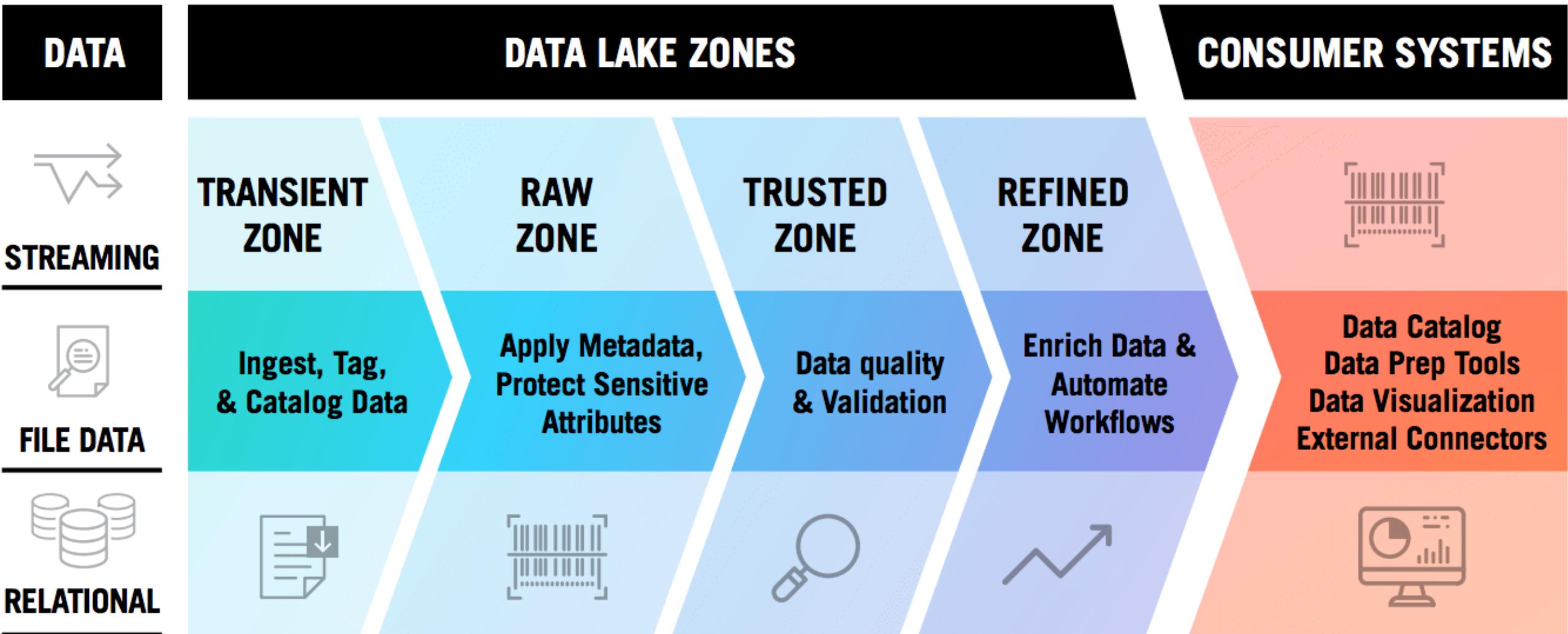
## Data Lakehouse



# Typical structure and missing use cases

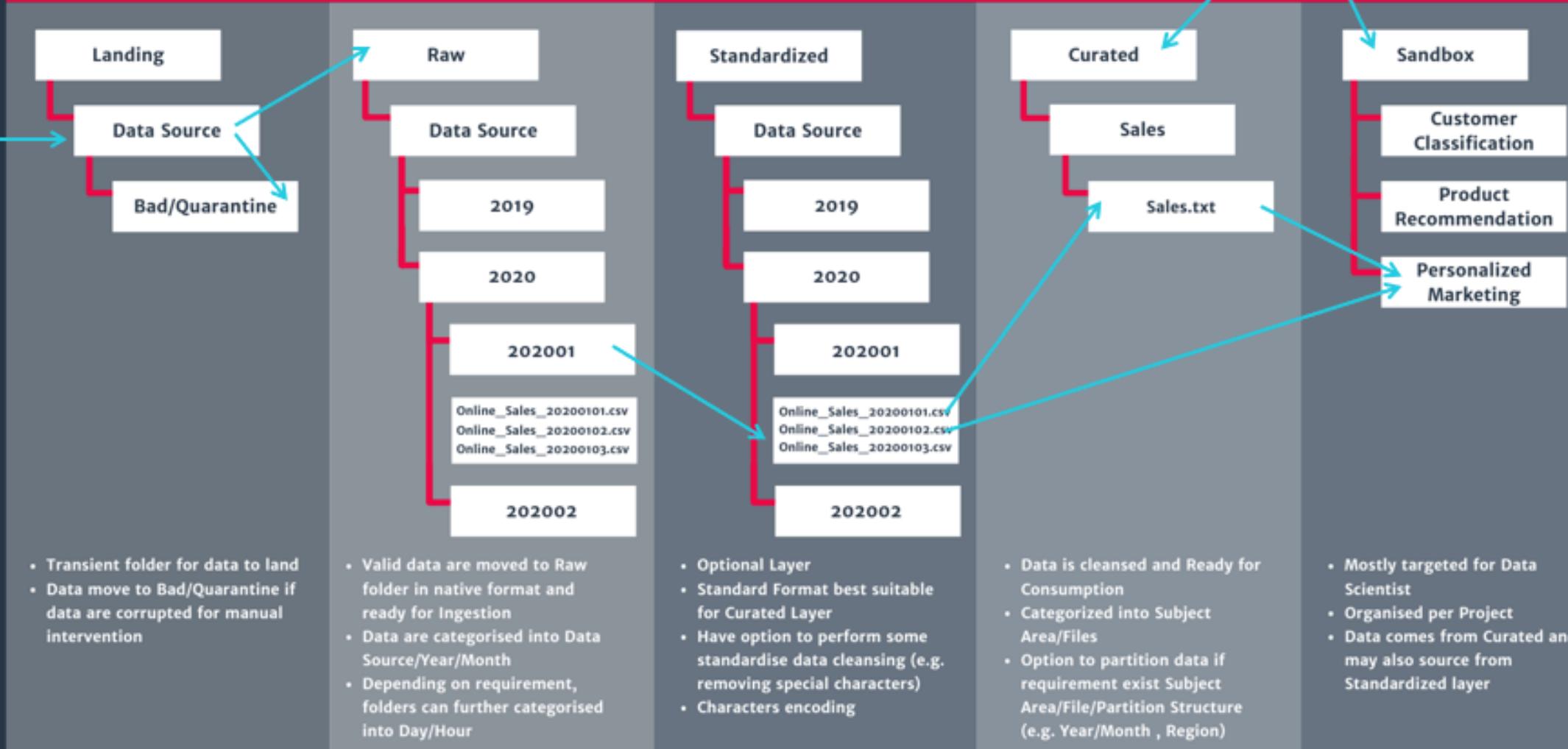


# Data Lake Zones

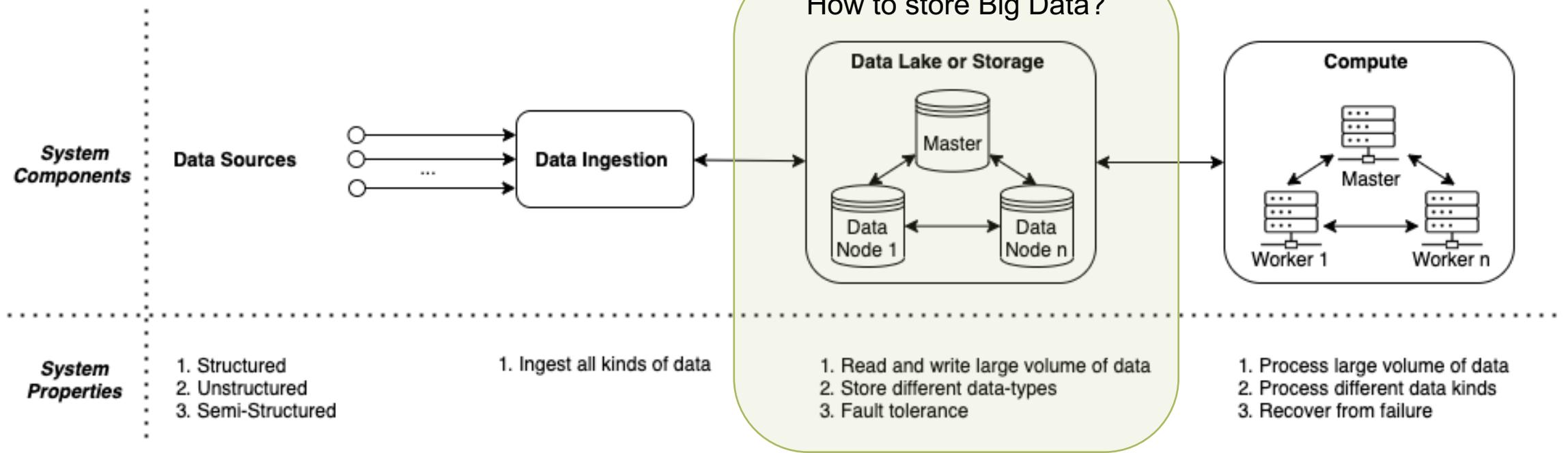




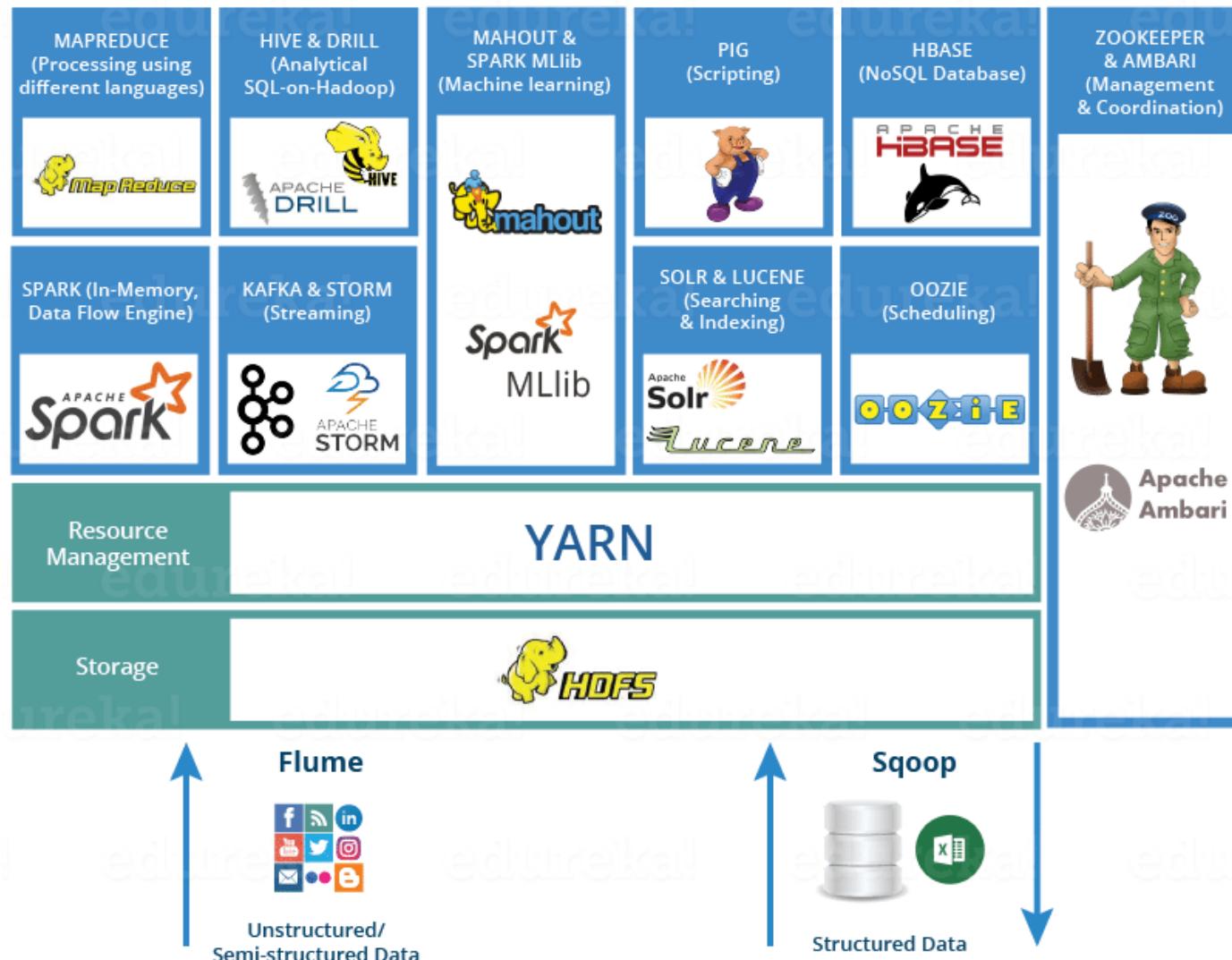
## Data Lake Storage Zones

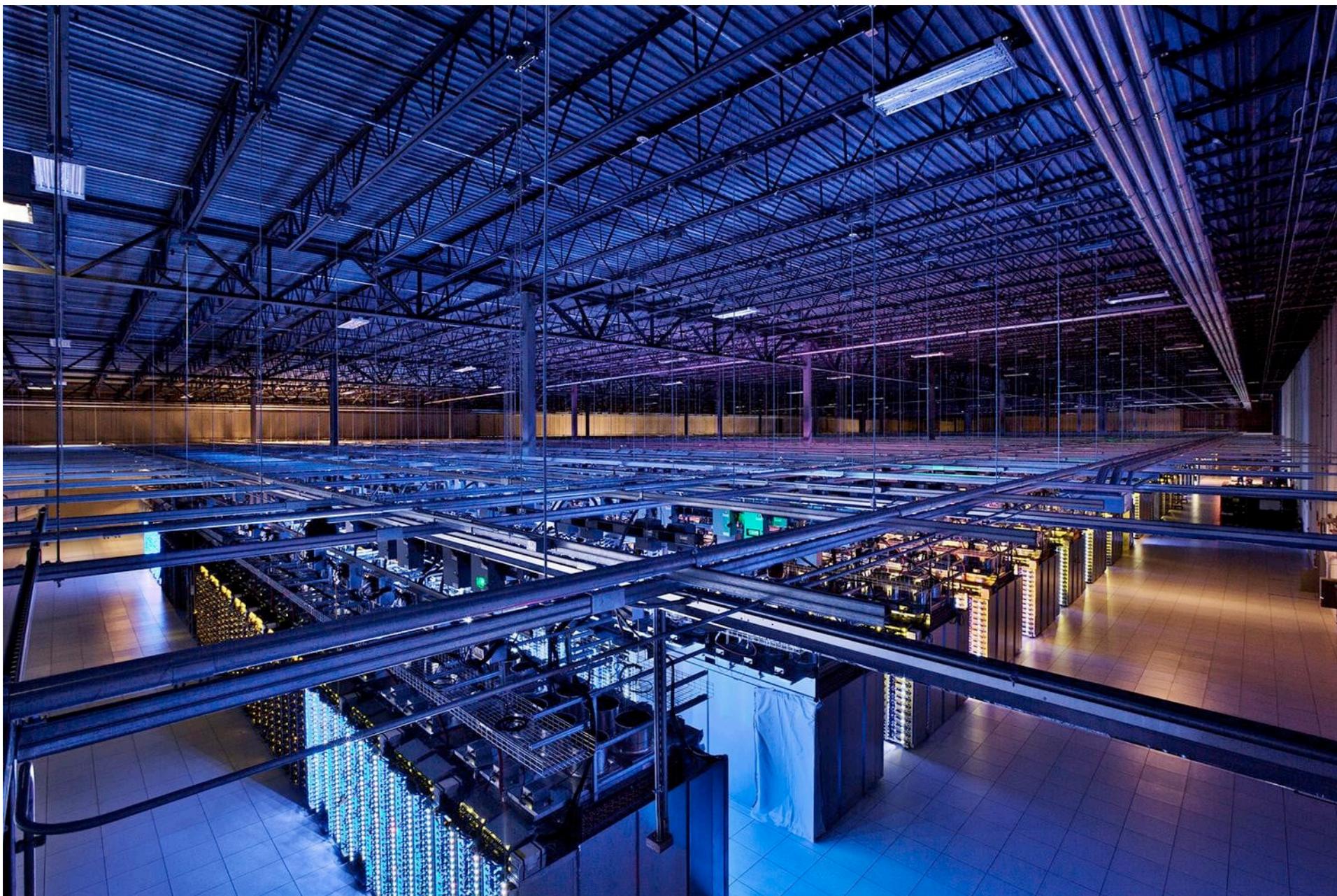


# Hadoop



# Center of the Apache Ecosystem





# The cluster fails

→ Google cluster: 103 commodity Linux servers

- 1-5% of disks will crash per year,
- frequent network failures

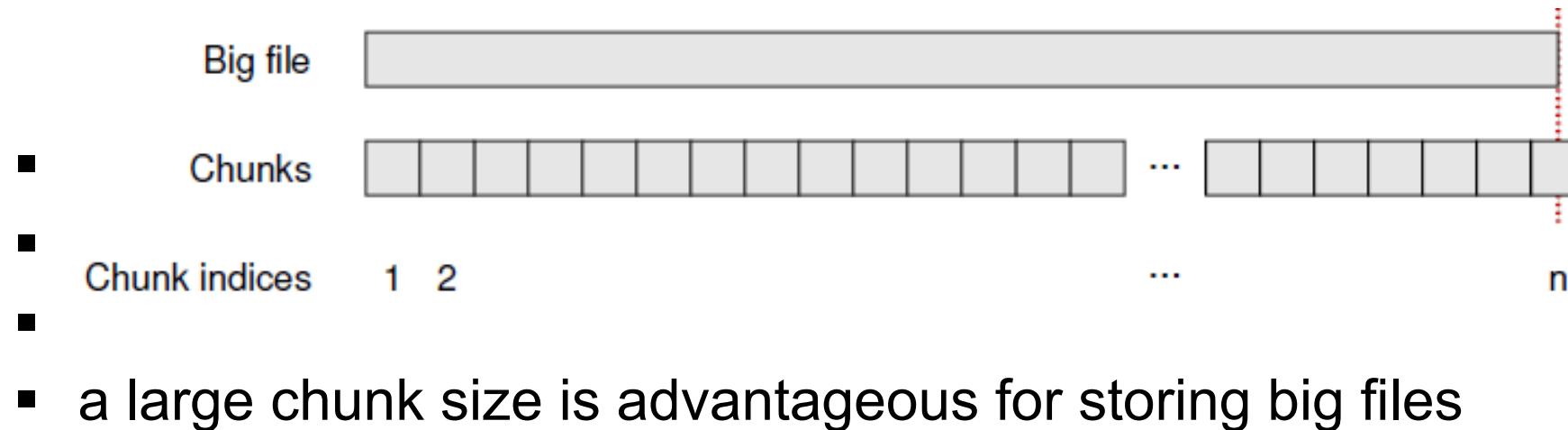
→ Some component of the DFS is always non-functional.

→ A DFS must be highly **fault-tolerant!**

# DFS design

→ Files (> 1 TB) won't fit on a disk, so:

- divide files into **chunks** (64 or 128 MB)



# DFS design

→ Since cluster components fail,

- store all data **redundantly**
- Typically 3 chunk replicas, on different machines, on different racks
- if one copy is unreachable, read another copy.

→ \_\_\_\_\_

---

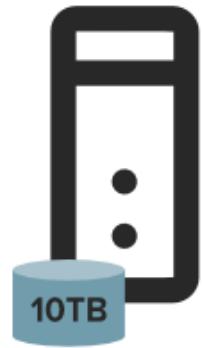
→ To implement redundancy, you need:

→ **metadata** for each file ( that maps the file into chunks)

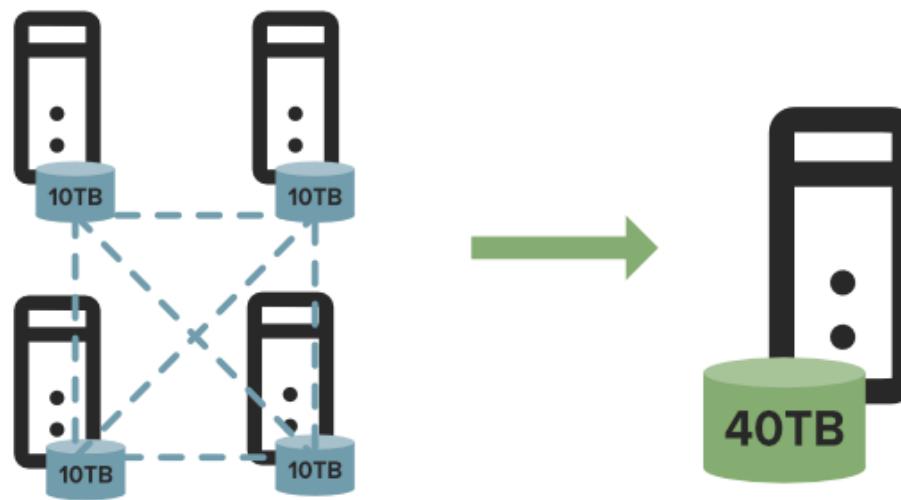
→ a **master** for the DFS (to maintain metadata).

# LFS vs DFS

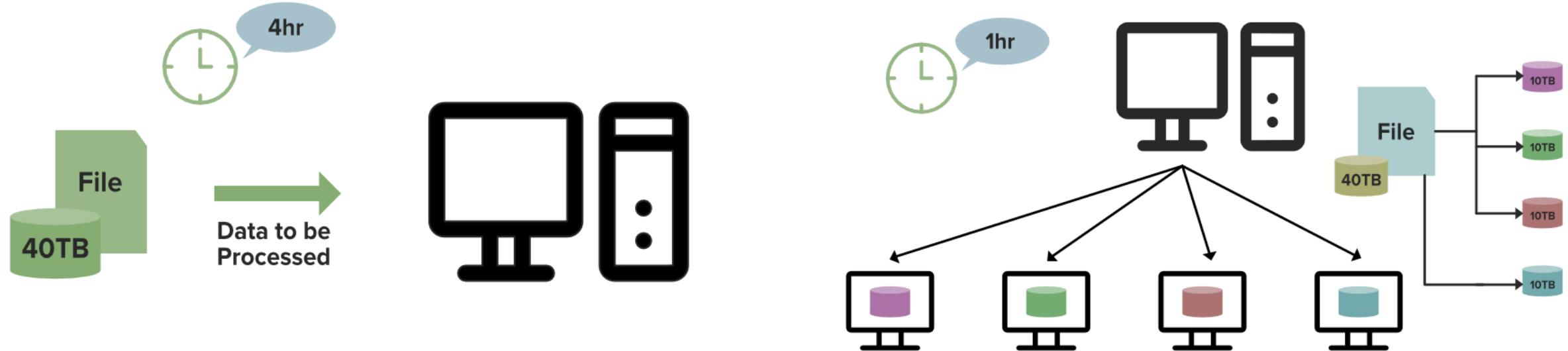
**Local File System**



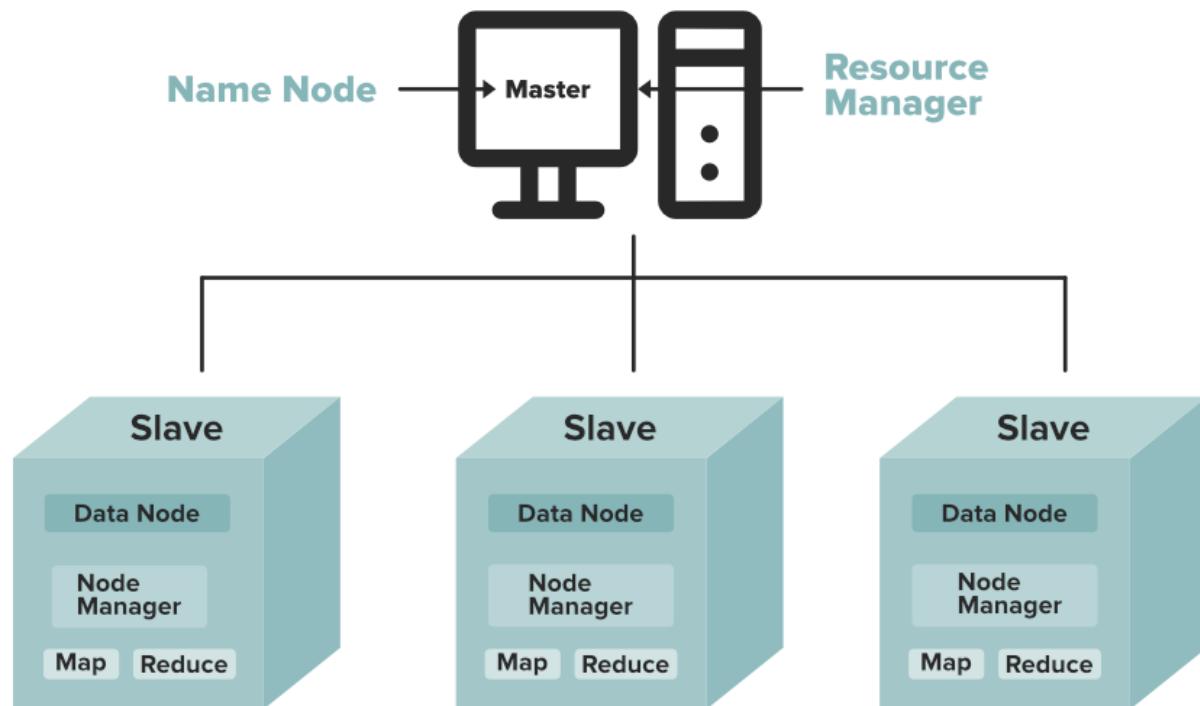
**DFS (Distributed File System)**



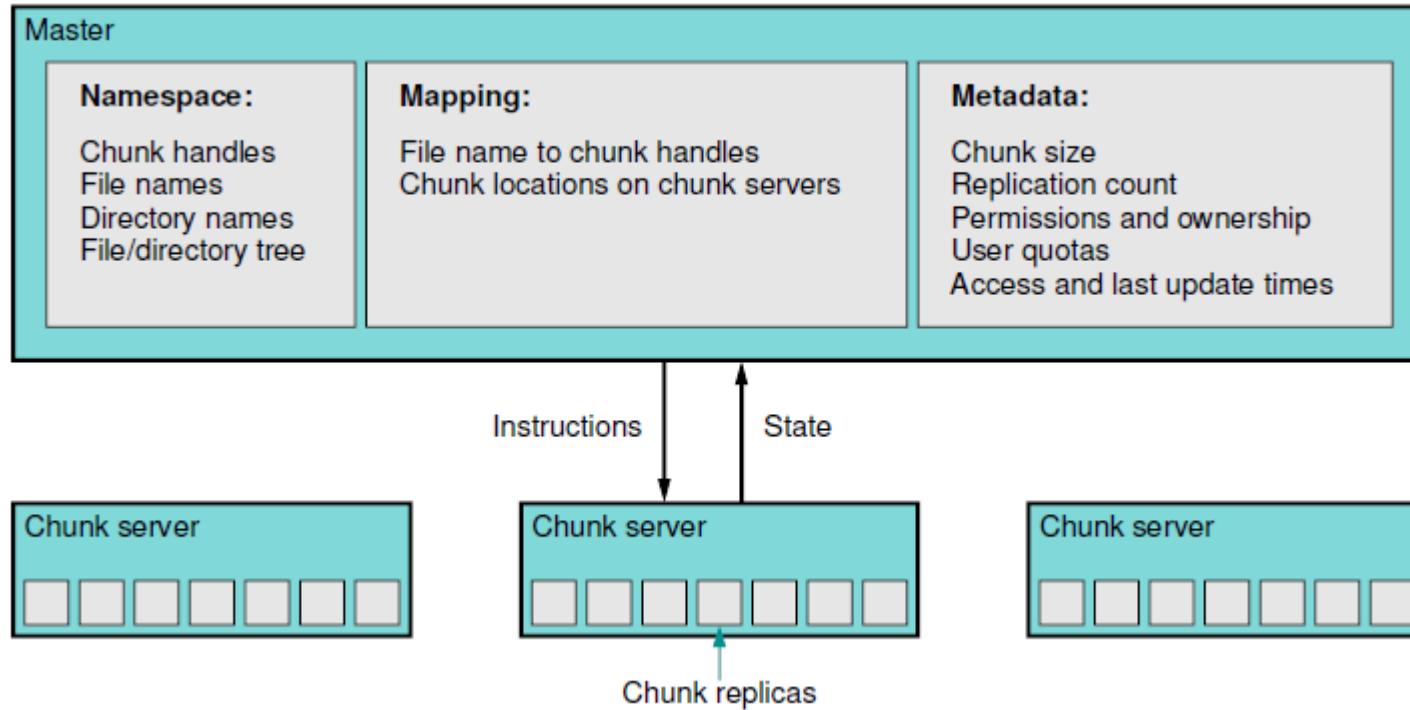
# DFS – Distribution of data (remove bottleneck)



# The Hadoop Setup

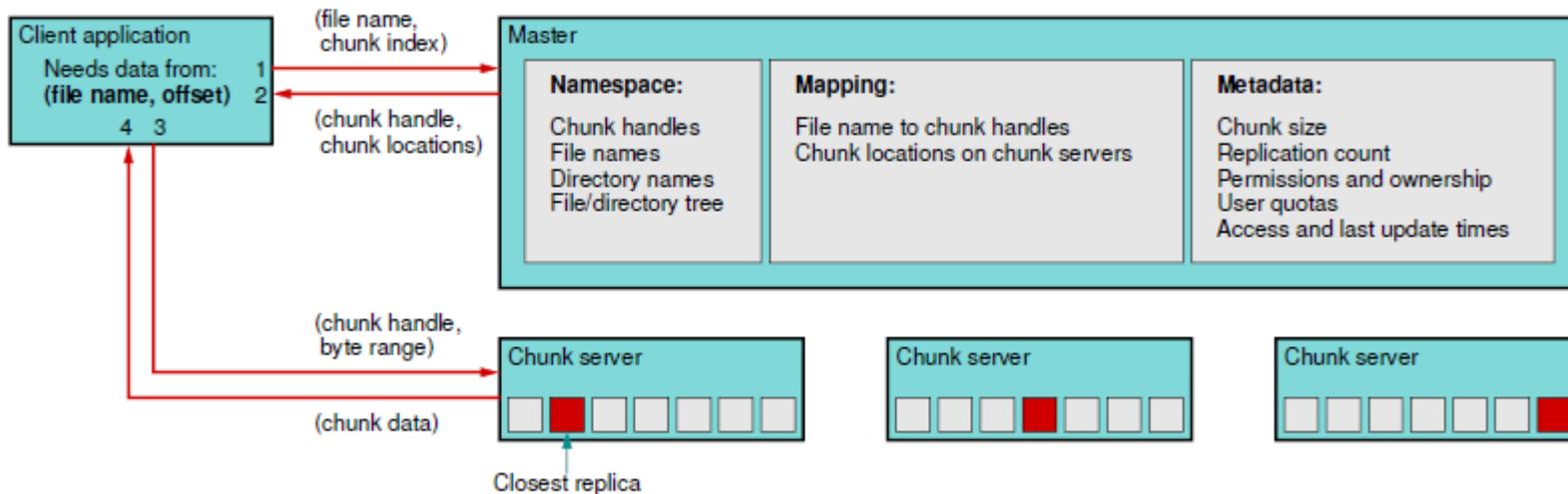


# The Namenode (Master server)



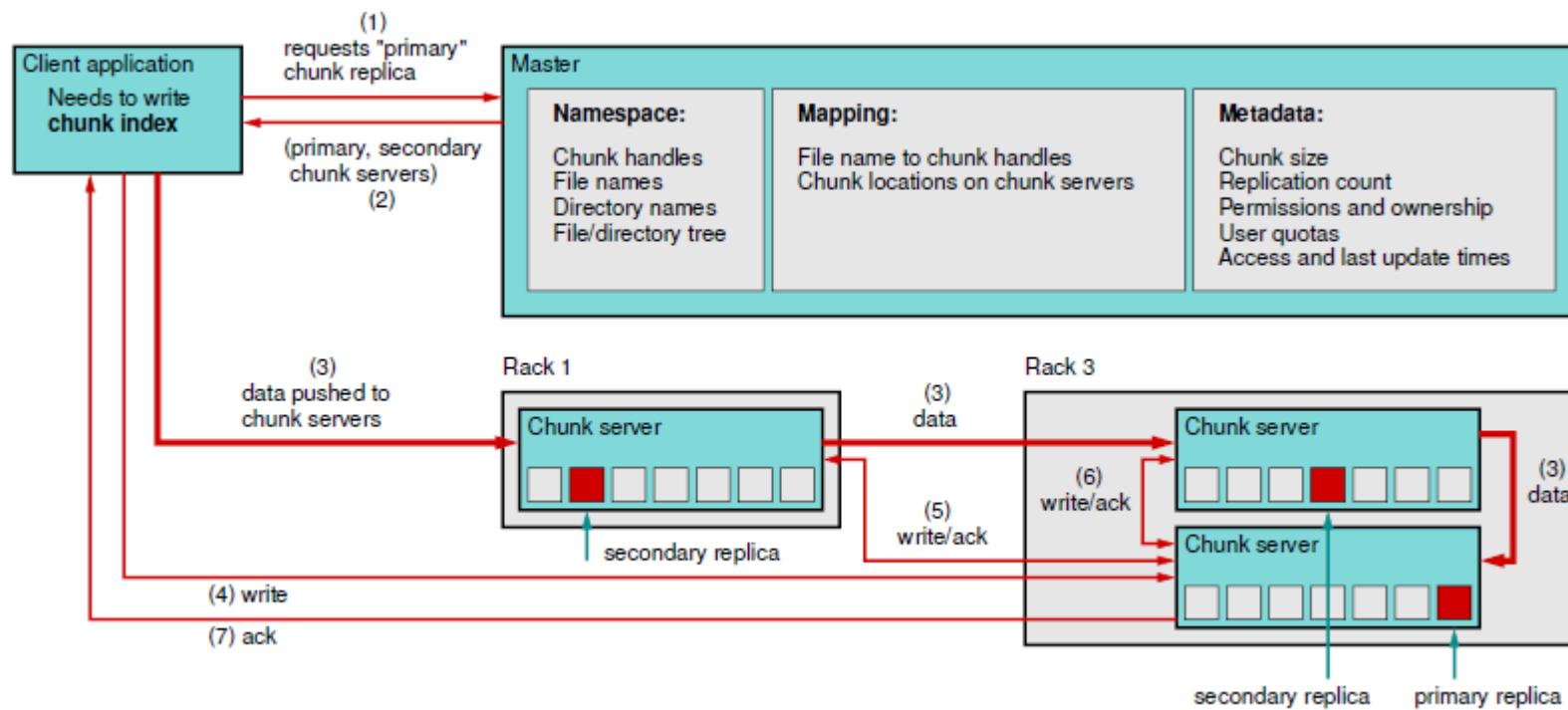
- Global, single master (bottleneck). "Shadow" masters.
- All metadata is kept in memory (fast!). Most is also logged on disk for reliability.
- Exception: the chunk locations. The master simply polls the servers at startup.

# A read operation



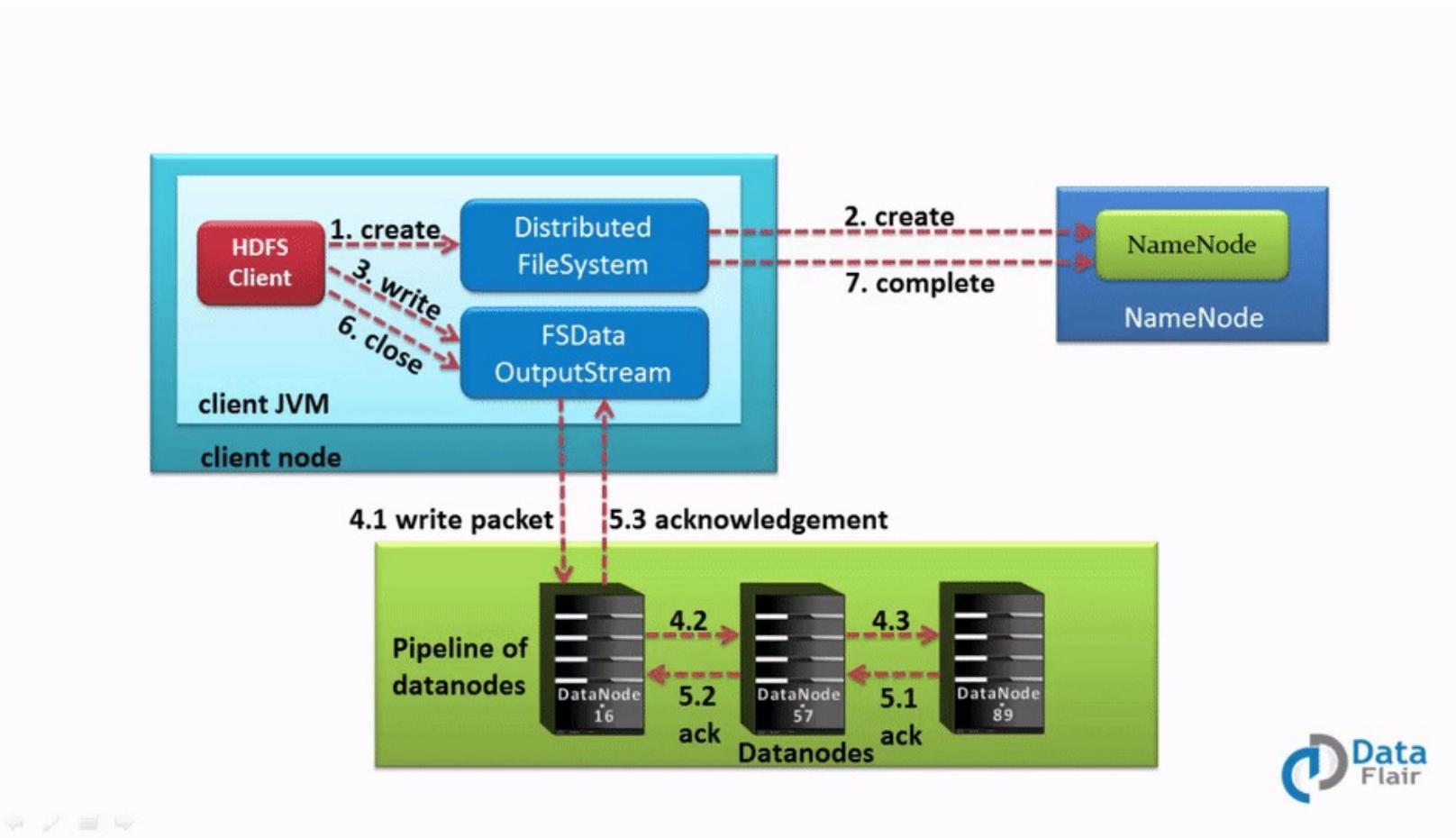
→ Data doesn't pass through the Master.

# A write operation



→ Data doesn't pass through the Master.

# A write operation



→ Src: <https://data-flair.training/blogs/hdfs-data-write-operation/>

# Reading data on a DFS

→ In practice, reads come in two types:

- large contiguous reads (1 MB+)
- small random reads (1 KB).

→ **Random reads** should be avoided:

→ performance-conscious applications **batch and sort their small reads** to advance linearly through the file.

# Writing data on a DFS

- Once written, the files are **mostly read**, and mostly sequentially.
- Most files are mutated by **appending** new data rather than overwriting data.

→ **Random writes** are practically non-existent:

→ they are supported, but are inefficient.

# DFS implementations

**GFS :** *Google File System, the original DFS*

**HDFS :** *Hadoop Distributed File System*  
a GFS clone

open-source, distributed by Apache Software

Foundation

**many others :** see [W4]

→ Facebook HDFS (2012): 100+ petabytes of disk space in one of its largest HDFS clusters (hundreds of thousands of servers).

# Some Sources for further reading:

W1. Facebook is collecting your data - 500 terabytes a day (2012)

<https://goo.gl/x2DZGz> (gigaom.com)

→W2. Where the Internet lives. Take a look inside Google's high-tech data centers

<https://www.google.com/about/datacenters/gallery/>

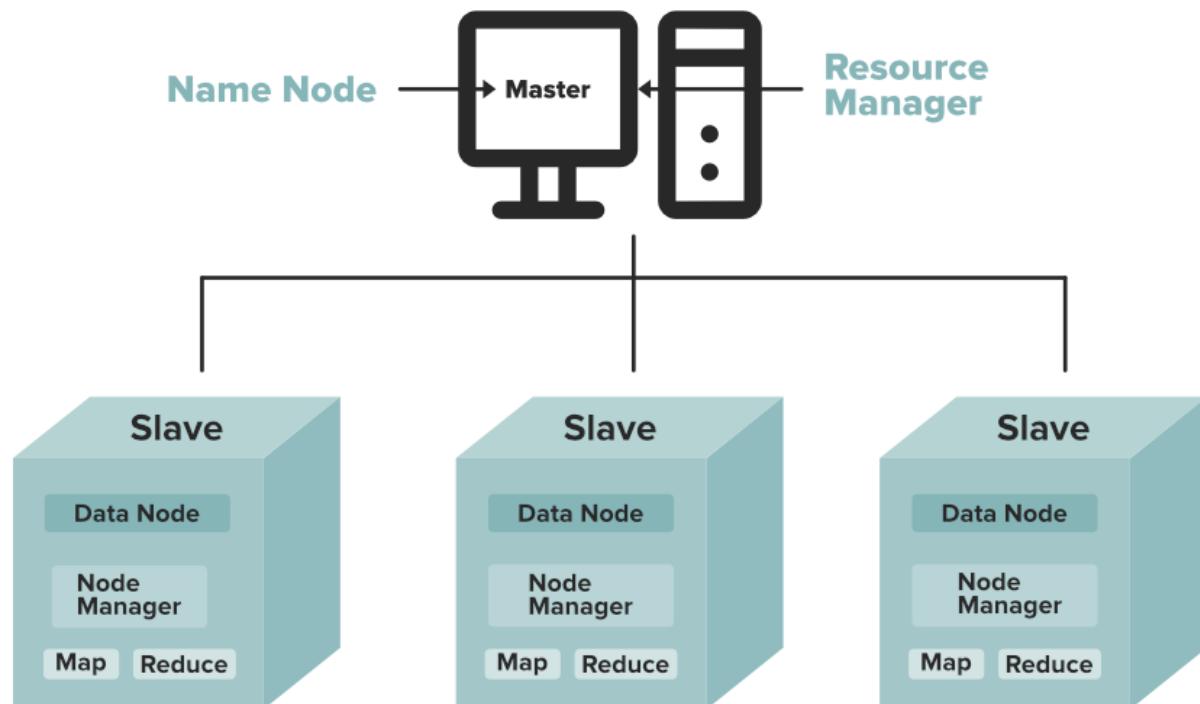
→W3. The Datacenter as a Computer: An Introduction to the Design of Warehouse-Scale Machines (2013)

<https://goo.gl/zrWA3G>

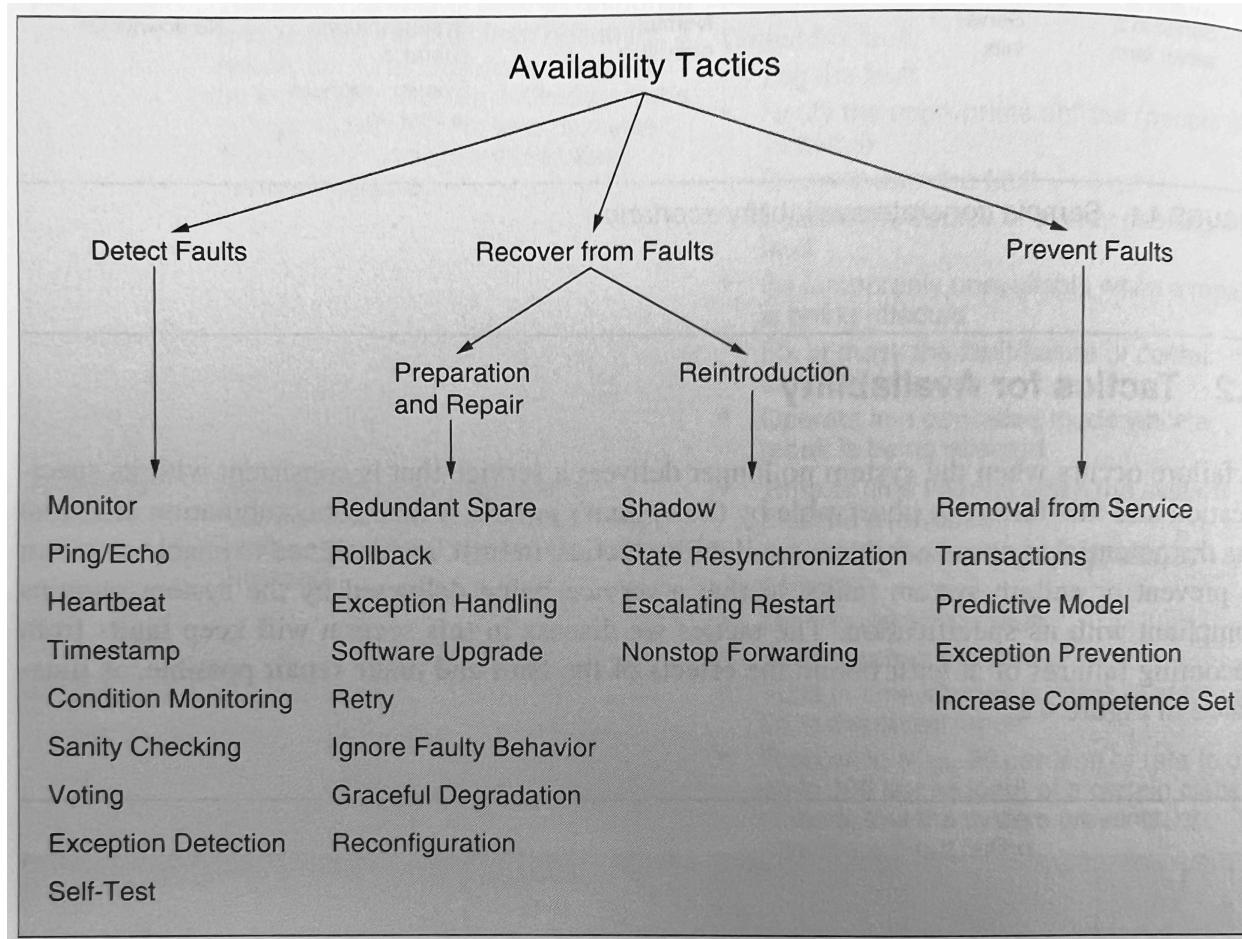
→W4. Distributed parallel fault-tolerant file systems

<https://goo.gl/h8q37d> (Wikipedia)

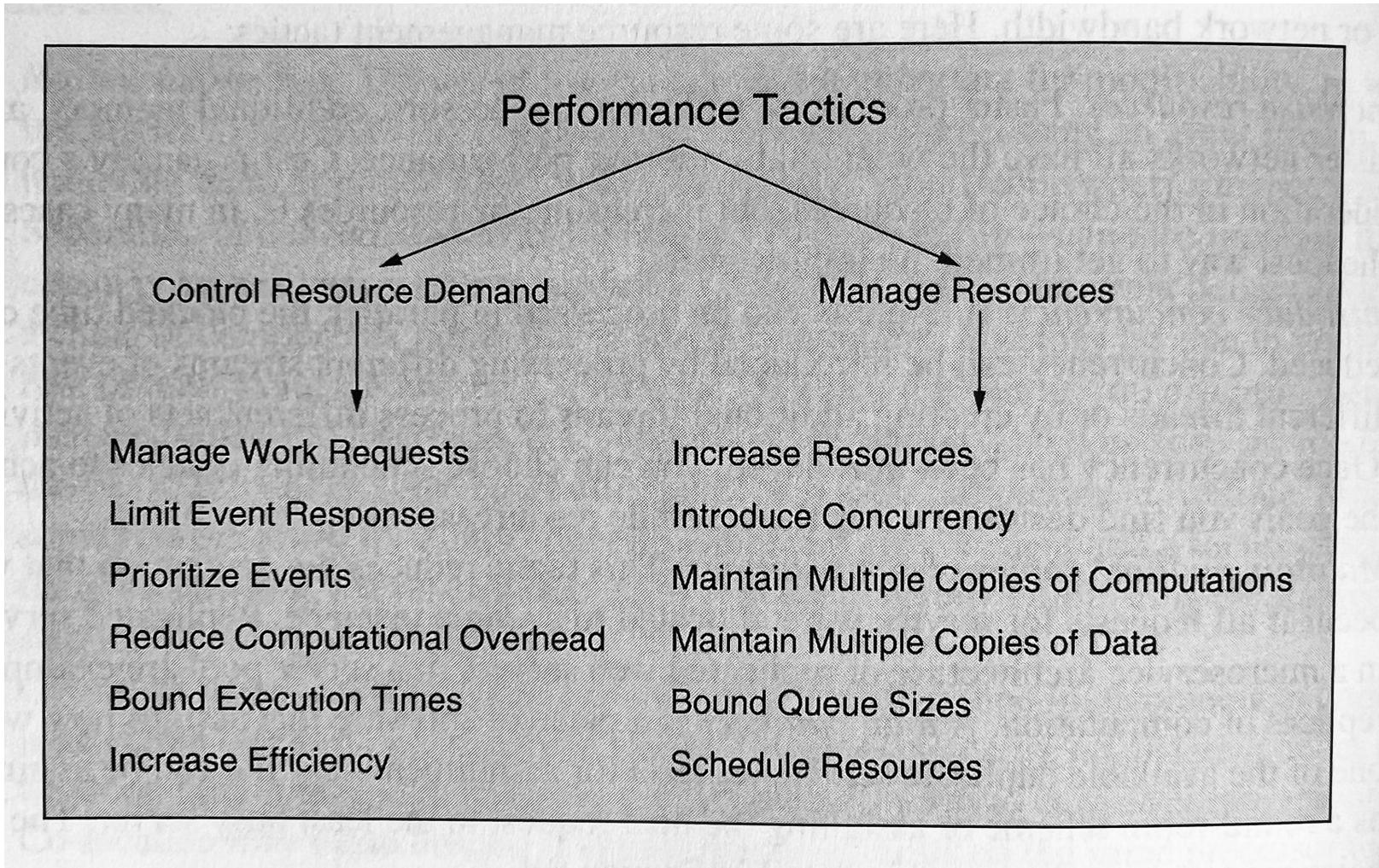
# Architecture: SPOF and Tactics



# Availability Tactics



# Performance



# Fileformats: Avro & Parquet

	 Parquet	 AVRO
Format Type	Column-based	Row-based
Built For/with	Optimized column-wise compression and querying in a splittable file format designed for efficient Map-Reduce processing	Compact binary storage and exchange of records, with schema evolution and support for many different programming languages
Schema Storage	Column metadata stored at the end of the file (allows for fast, one-pass writing)	Stored in human-readable JSON format at the beginning of each message or file
Use Case	Quickly query all the values from a particular column across a very large dataset. For example, compute the average price of purchases over millions of purchase records	Share entire records between applications. For example, event data of in-app purchases for use by many downstream applications such as logging, auditing, and business analytics

# How do they compare?



- Data format for Data Serialization
- Tuple based
  - Great when you read entire records
- Schema Based
- Schama Evolution
  - Can change current props and add



## Parquet

- Data format for Data Serialization
- Column based
  - Great for wide tables, where only few are used for each query
- Schema Based
- Schema Evolution
  - Can add props, but not change them
- Good for nested datasets

# Avro Primitives

Type	Description	Schema
null	The absence of a value	"null"
boolean	A binary value	"boolean"
int	32-bit signed integer	"int"
long	64-bit signed integer	"long"
float	Single-precision (32-bit) IEEE 754 floating-point number	"float"
double	Double-precision (64-bit) IEEE 754 floating-point number	"double"
bytes	Sequence of 8-bit unsigned bytes	"bytes"
string	Sequence of Unicode characters	"string"

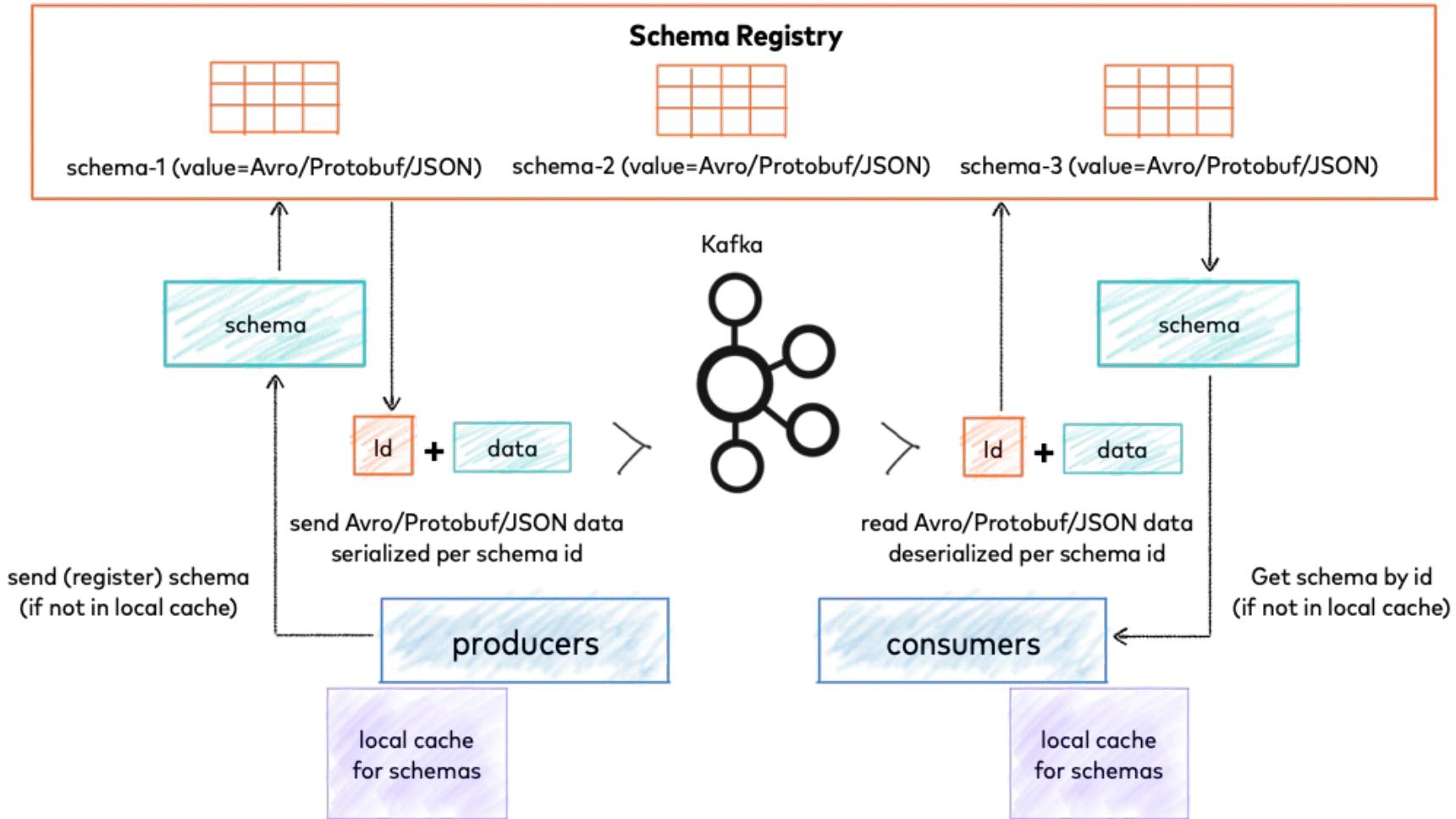
Type	Description	Schema example
array	An ordered collection of objects. All objects in a particular array must have the same schema.	{ " <b>type</b> ": "array", " <b>items</b> ": "long" }
map	An unordered collection of key-value pairs. Keys must be strings and values may be any type, although within a particular map, all values must have the same schema.	{ " <b>type</b> ": "map", " <b>values</b> ": "string" }
record	A collection of named fields of any type.	{ " <b>type</b> ": "record", " <b>name</b> ": "WeatherRecord", " <b>doc</b> ": "A weather reading.", " <b>fields</b> ": [ { <b>"name"</b> : "year", <b>"type"</b> : "int"}, { <b>"name"</b> : "temperature", <b>"type"</b> : "int"}, { <b>"name"</b> : "stationId", <b>"type"</b> : "string"} ] }
enum	A set of named values.	{ " <b>type</b> ": "enum", " <b>name</b> ": "Cutlery", " <b>doc</b> ": "An eating utensil.", " <b>symbols</b> ": ["KNIFE", "FORK", "SPOON"] }
fixed	A fixed number of 8-bit unsigned bytes.	{ " <b>type</b> ": "fixed", " <b>name</b> ": "Md5Hash", " <b>size</b> ": 16 }

# Avro Complex Data Types

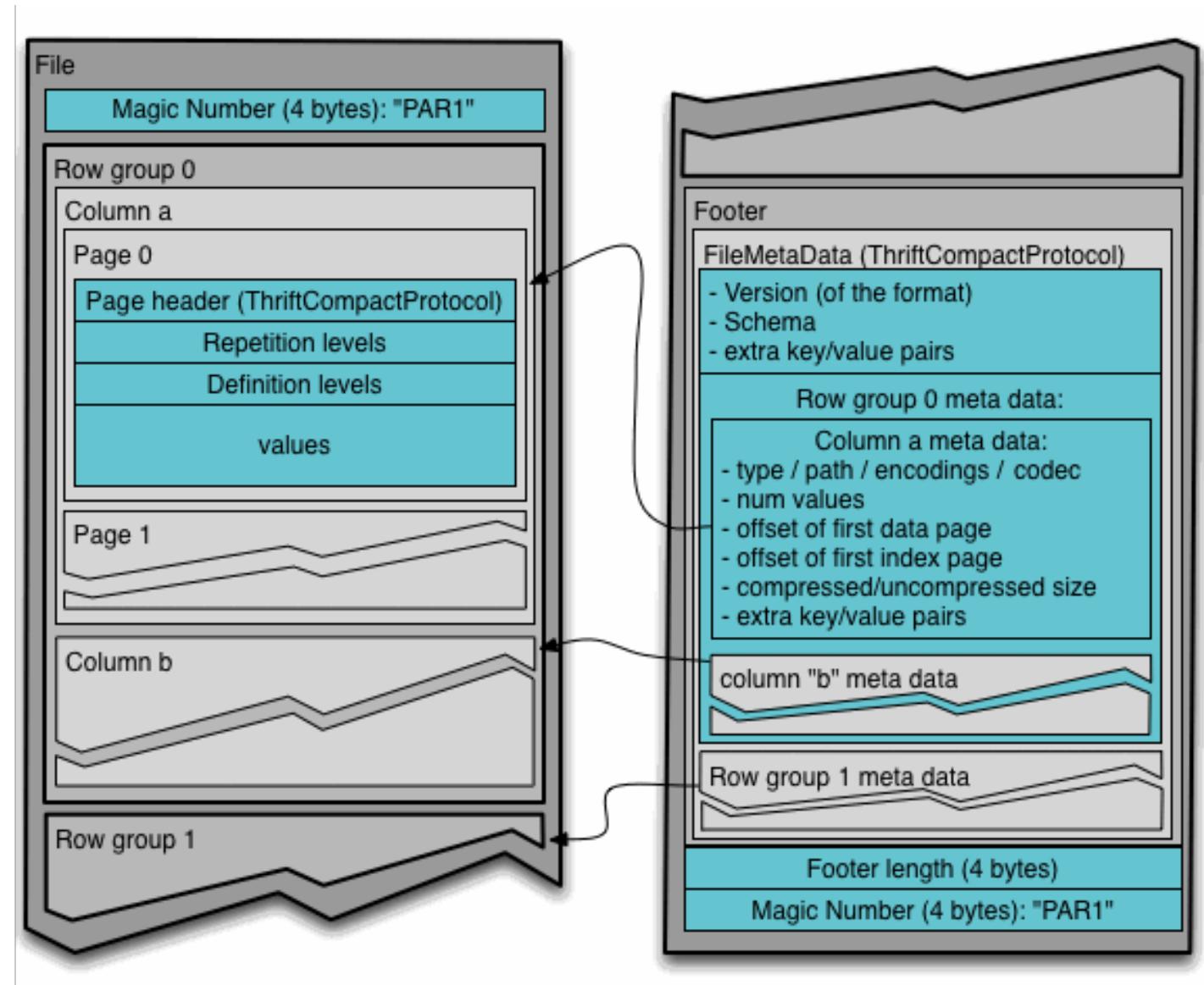
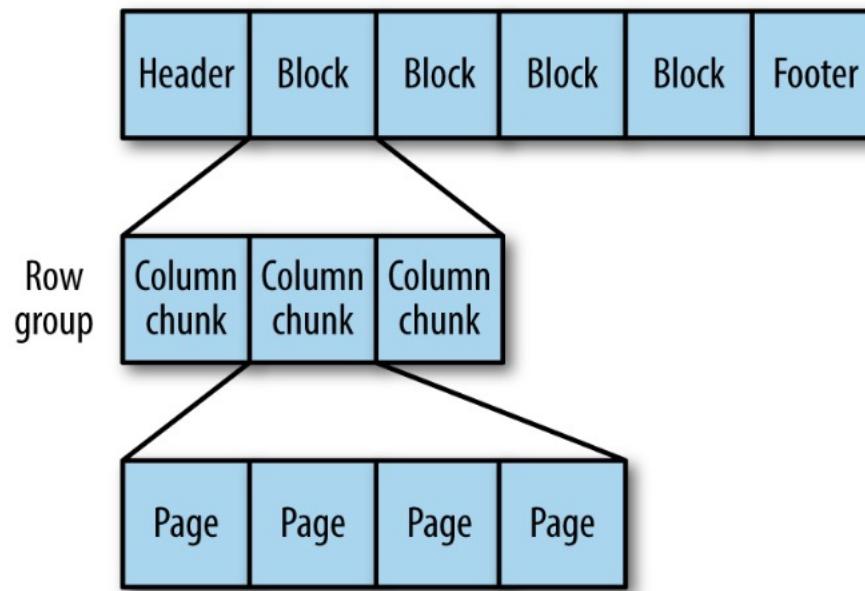
# Avro Schema

```
{"namespace": "example.avro",
"type": "record",
"name": "User",
"fields": [
    {"name": "name", "type": "string"},
    {"name": "favorite_number", "type": ["int", "null"]},
    {"name": "favorite_color", "type": ["string", "null"]}
]
}
```

# Usage example



# Parquet Schema and Structure



# Parquet Primitives

Type	Description
boolean	Binary value
int32	32-bit signed integer
int64	64-bit signed integer
int96	96-bit signed integer
float	Single-precision (32-bit) IEEE 754 floating-point number
double	Double-precision (64-bit) IEEE 754 floating-point number
binary	Sequence of 8-bit unsigned bytes
fixed_len_byte_array	Fixed number of 8-bit unsigned bytes

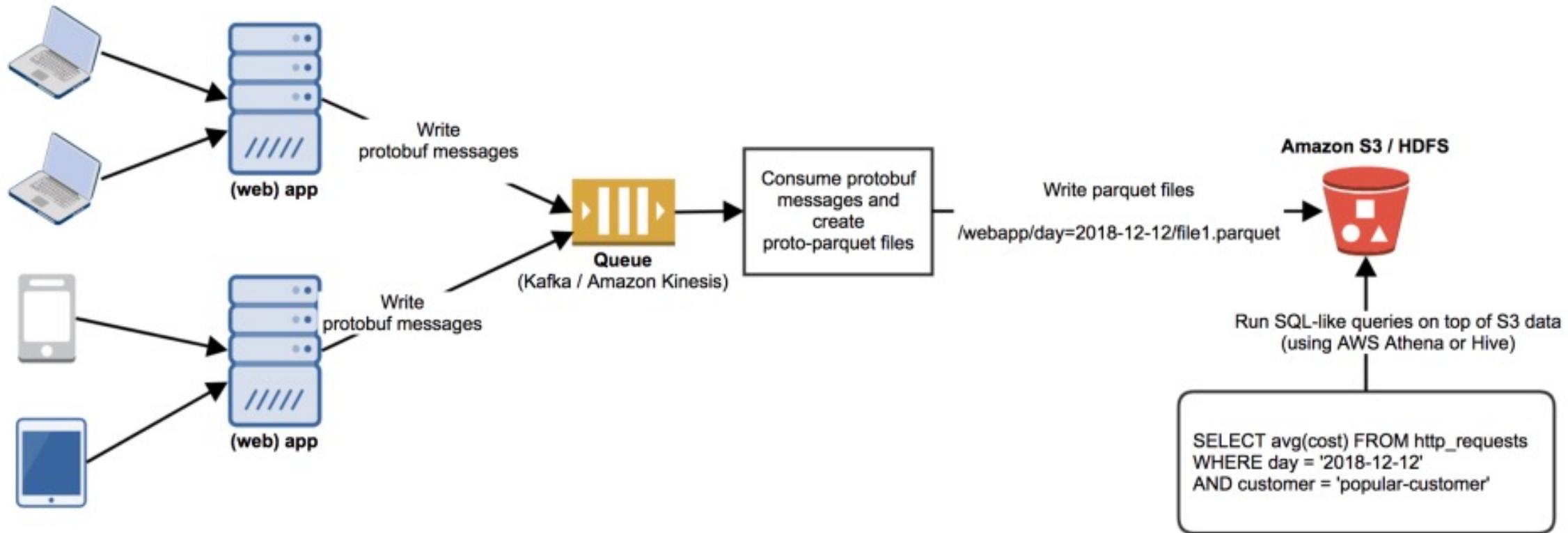
# Parquet Logical Types

Logical type annotation	Description	Schema example
UTF8	A UTF-8 character string. Annotates binary.	<pre>message m {     required binary a (UTF8); }</pre>
ENUM	A set of named values. Annotates binary.	<pre>message m {     required binary a (ENUM); }</pre>
DECIMAL( <i>precision, scale</i> )	An arbitrary-precision signed decimal number. Annotates int32, int64, binary, or fixed_len_byte_array.	<pre>message m {     required int32 a (DECIMAL(5,2)); }</pre>
DATE	A date with no time value. Annotates int32. Represented by the number of days since the Unix epoch (January 1, 1970).	<pre>message m {     required int32 a (DATE); }</pre>
LIST	An ordered collection of values. Annotates group.	<pre>message m {     required group a (LIST) {         repeated group list {             required int32 element;         }     } }</pre>
MAP	An unordered collection of key-value pairs. Annotates group.	<pre>message m {     required group a (MAP) {         repeated group key_value {             required binary key (UTF8);             optional int32 value;         }     } }</pre>

# Parquet Schema Example

```
message WeatherRecord {  
    required int32 year;  
    required int32 temperature;  
    required binary stationId (UTF8);  
}
```

# Parquet Usage Example



# Like to know more about parquet?

→Read the Dremel Paper:

→<https://storage.googleapis.com/pub-tools-public-publication-data/pdf/36632.pdf>

# **Exercises!**