# TODAY'S AGENDA

Background

Architecture Overview

Query Execution

Project Discussion

# BACKGROUND

Organizations use **on-line analytical processing** (OLAP) systems to extract new information from existing data sets.

Historically these workloads were run in a monolithic DBMSs that had all an organization's data in centralized managed storage…
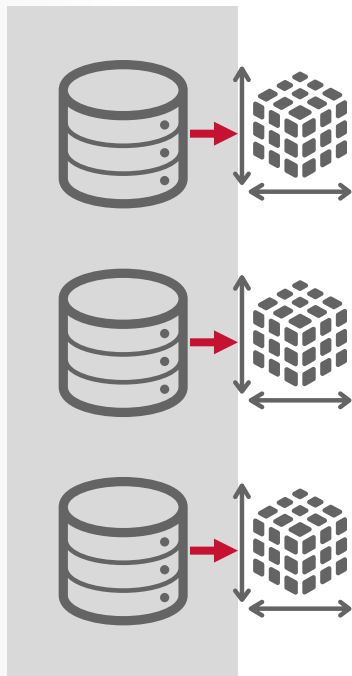
# 1990s – DATA CUBES

DBMSs would maintain multi-dimensional arrays as pre-computed aggregations to speed up queries.
→ Periodically refreshed materialized views.
→ Administrator had to specify cubes ahead of time.

Data cubes were often introduced in existing operational DBMSs originally designed to operate on row-oriented data.

# 1990s – DATA CUBES



Data cubes : simply pre-aggregation

```
SELECT product, region, cdate,
       SUM(amount) AS total_sales
  FROM sales
 GROUP BY CUBE (product, region, cdate);
```

*OLTP Databases*
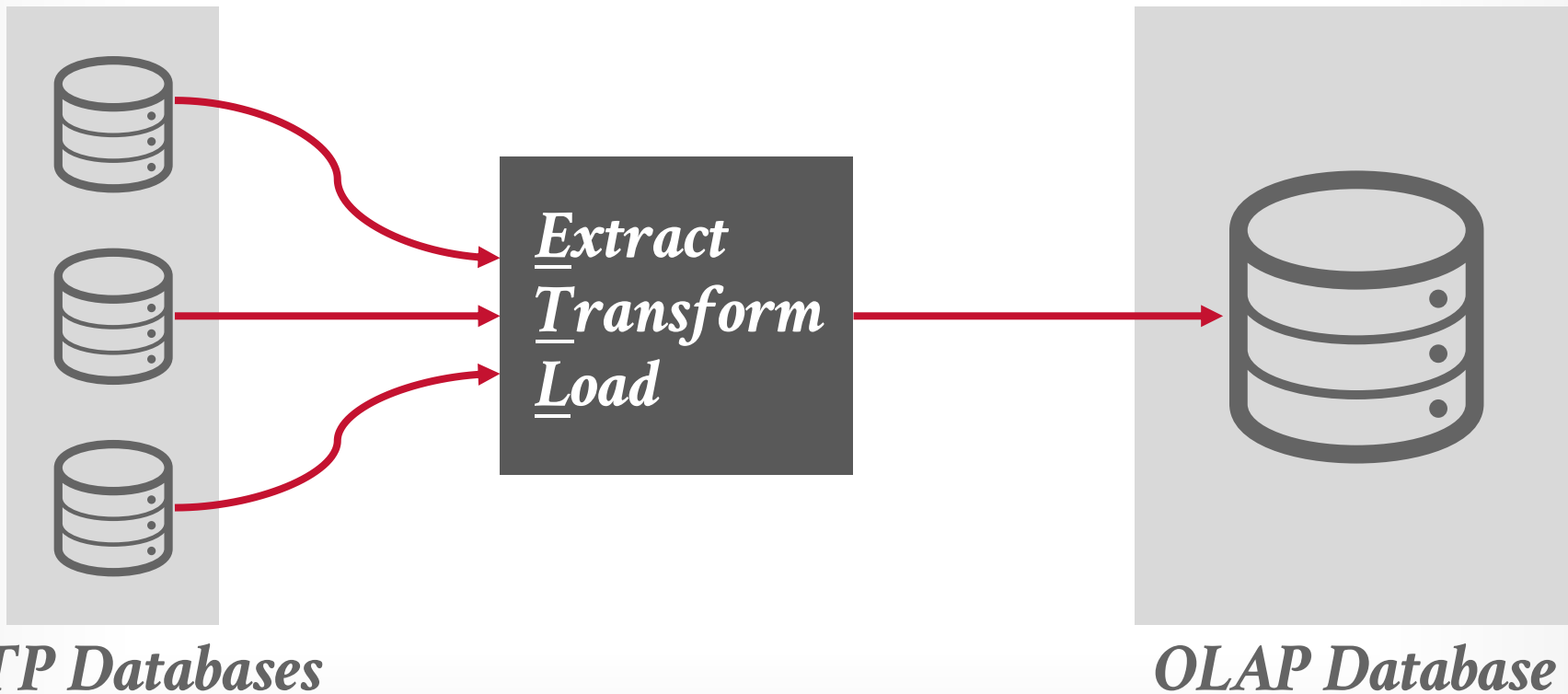
# 2000s - DATA WAREHOUSES

Monolithic DBMSs designed to efficiently execute OLAP workloads using shared-nothing architectures and column-oriented data.
→ Many systems from this era started as forks of Postgres.

DBMS-managed storage using proprietary data encoding / formats.

# 2000s — DATA WAREHOUSES



**OLTP Databases**

**OLAP Database**

# 2010s – SHARED-DISK ENGINES

Shared-disk DBMS architectures that relied on third-party distributed storage (object stores) instead of using a custom storage manager.

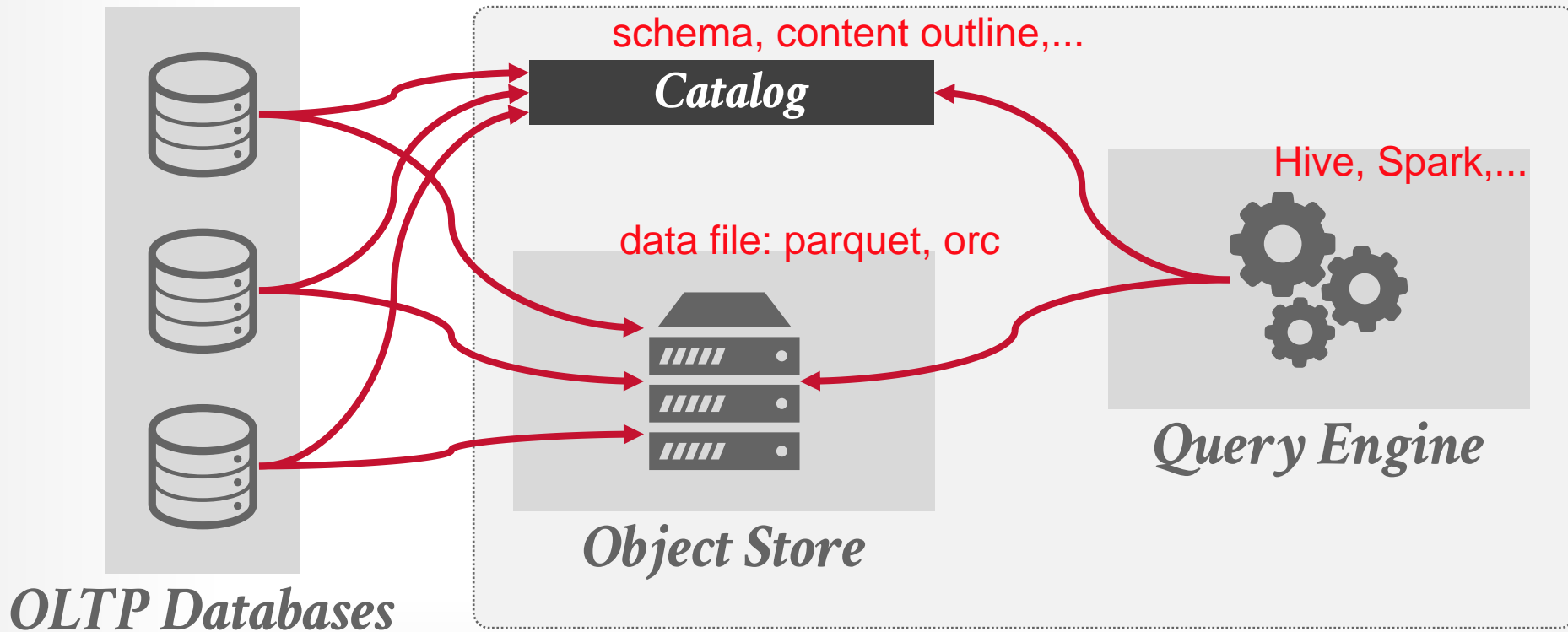First generation of these systems managed data files themselves.

E.x: Snowflake data format + data store in S3

Newer systems allow external entities to add new data files to storage without enforcing schema (**lakehouse**).

Idea: Seperate compute and storage layer --> use shared-disk engines

# 2010s – SHARED-DISK ENGINES



schema, content outline,...

**Catalog**

Hive, Spark,...

data file: parquet, orc

*Query Engine*

*Object Store*

*OLTP Databases*

# 2020S – LAKEHOUSE SYSTEMS

Middleware for data lakes that adds support for ==better schema control / versioning== with transactional CRUD operations.

<span style="color:red">Schema evolution (re-create table)</span>

→ Store changes in row-oriented log-structured files with indexes.

→ Periodically compact recently added data into read-only columnar files.

We will <u>not</u> be covering this aspect of these systems in this course.

databricks

Apache hudi

ICEBERG

snowflake

# 2020S — LAKEHOUSE SYSTEMS

**Observation #1: People want to execute more than just SQL on data.**

**Observation #2: Decoupling data storage from DBMS reduces ingest/egress barriers.**

**Observation #3: Most data is unstructured / semi-structured.**
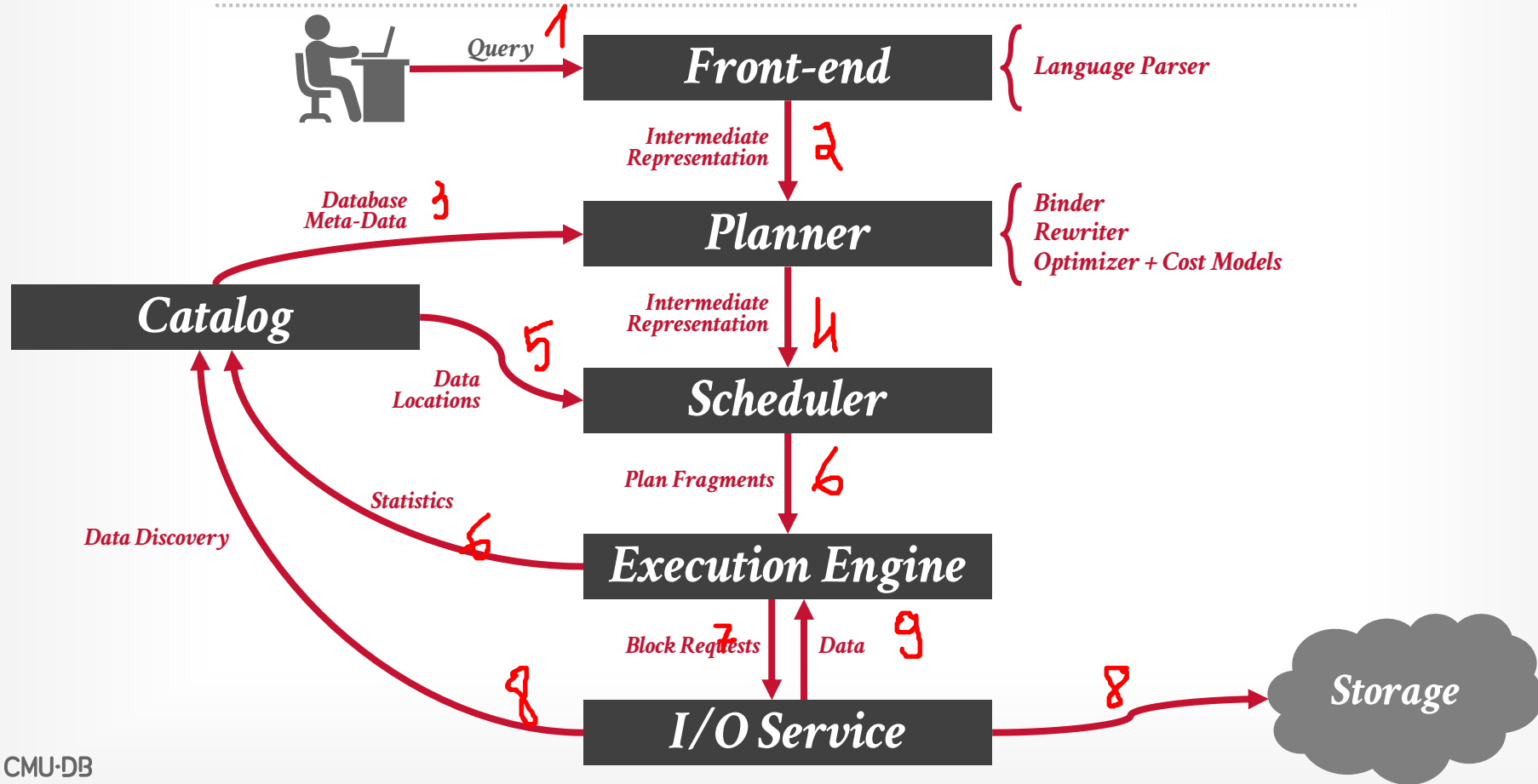
# OLAP DBMS COMPONENTS

One recent trend of the last decade is the breakout of OLAP DBMS components into standalone services and libraries:
→ System Catalogs
→ Intermediate Representation
→ Query Optimizers
→ File Format / Access Libraries
→ Execution Engines / Fabrics

Lots of engineering challenges to make these components interoperable + performant.
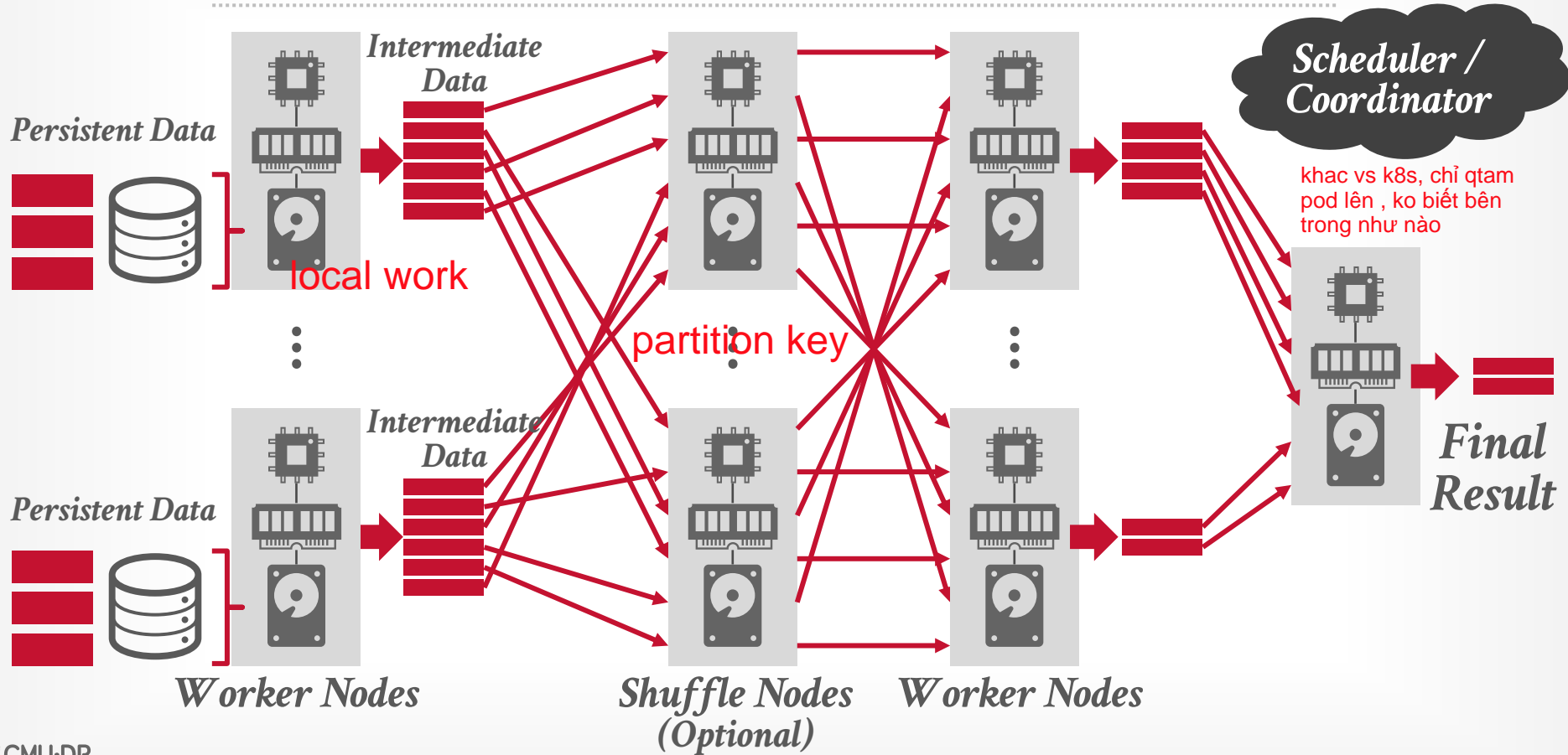
# ARCHITECTURE OVERVIEW



**Query** 1

## Front-end

{ *Language Parser*

*Intermediate Representation* 2

## Planner

{ *Binder*
*Rewriter*
*Optimizer + Cost Models*

**Database Meta-Data** 3

## Catalog

*Intermediate Representation* 4

**Data Locations** 5

## Scheduler

**Plan Fragments** 6

**Statistics** 6

## Execution Engine

**Data Discovery**

**Block Requests** 7 **Data** 9

## I/O Service

8

## Storage

# DISTRIBUTED QUERY EXECUTION

Executing an OLAP query in a distributed DBMS is roughly the same as on a single-node DBMS.
→ Query plan is a DAG of physical operators.

For each operator, the DBMS considers where input is coming from and where to send output.
→ Table Scans
→ Joins
→ Aggregations
→ Sorting

# DISTRIBUTED QUERY EXECUTION



Intermediate Data

Persistent Data

local work

partition key

Scheduler / Coordinator

khac vs k8s, chỉ qtam pod lên , ko biết bên trong như nào

Final Result

Intermediate Data

Persistent Data

Worker Nodes

Shuffle Nodes (Optional)

Worker Nodes

# DATA CATEGORIES

**Persistent Data:**
→ The "source of record" for the database (e.g., tables).
→ Modern systems assume that these data files are immutable but can support updates by rewriting them.

**Intermediate Data:**
→ Short-lived artifacts produced by query operators during execution and then consumed by other operators.
→ The amount of intermediate data that a query generates has little to no correlation to amount of persistent data that it reads or the execution time.

# DISTRIBUTED SYSTEM ARCHITECTURE

A distributed DBMS's system architecture specifies the location of the database's persistent data files. This affects how nodes coordinate with each other and where they retrieve/store objects in the database.

Two approaches (not mutually exclusive):
→ **Push Query to Data**
→ **Pull Data to Query**

THE CASE FOR SHARED NOTHING
HPTS 1985

# PUSH VS. PULL

**Approach #1: Push Query to Data**
→ Send the query (or a portion of it) to the node that contains the data.
→ Perform as much filtering and processing as possible where data resides before transmitting over network.

**Approach #2: Pull Data to Query**
→ Bring the data to the node that is executing a query that needs it for processing.
→ This is necessary when there is no compute resources available where persistent data files are located.

# PUSH VS. PULL

## Approac[h]
→ Send th[e]
   contain[...]
→ Perfor[m]
   data re[...]

## Approa[ch]
→ Bring [...]
   needs it for processing.
→ T
   his is necessary when there is no compute resources
   available where persistent data files are located.

---

Filtering and retrieving data using Amazon S3 Select

PDF | RSS

amazon

With Amazon S3 Select, you can use simple structured query language (SQL) statements to filter the contents of an Amazon S3 object and retrieve just the subset of data that you need. By using Amazon S3 Select to filter this data, you can reduce the amount of data that Amazon S3 transfers, which reduces the cost and latency to retrieve this data.

Amazon S3 Select works on objects stored in CSV, JSON, or Apache Parquet format. It also works with objects that are compressed with GZIP or BZIP2 (for CSV and JSON objects only), and server-side encrypted objects. You can specify the format of the results as either CSV or JSON, and you can determine how the records in the result are delimited.

You pass SQL expressions to Amazon S3 in the request. Amazon S3 Select supports a subset of SQL. For more information about the SQL elements that are supported by Amazon S3 Select, see SQL reference for Amazon S3 Select.

You can perform SQL queries using AWS SDKs, the SELECT Object Content REST API, the AWS Command Line Interface (AWS CLI), or the Amazon S3 console. The Amazon S3 console limits the amount of data returned to 40 MB. To retrieve more data, use the AWS CLI or the API.
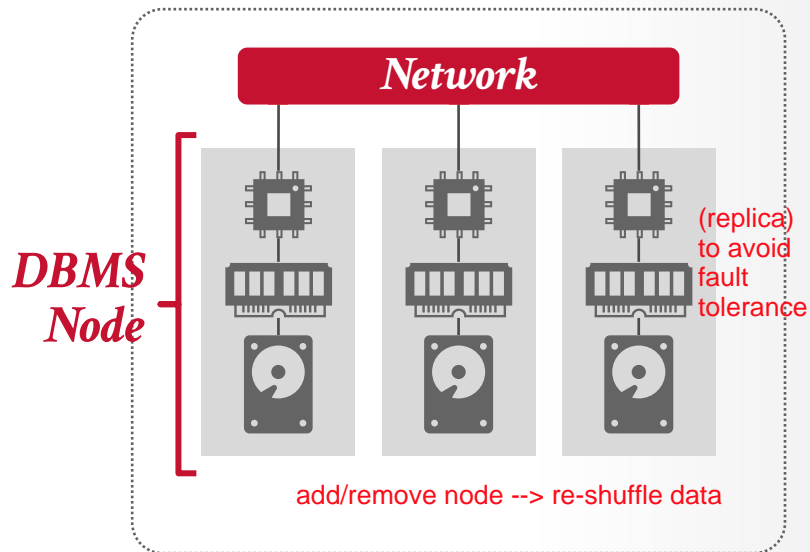
# SHARED-NOTHING

Each DBMS instance has its own
CPU, memory, locally-attached disk.
→ Nodes only communicate with each other
via network.

Database is partitioned into disjoint
subsets across nodes.
→ Adding a new node requires physically
moving data between nodes.

Since data is local, the DBMS can
access it via POSIX API.

f.write | f.read
do not know the data is actually be stored

**Network**

**DBMS Node**

(replica)
to avoid
fault
tolerance

add/remove node --> re-shuffle data

object store and NFS (network file system) provide
same interface except that object store provide more

# SHARED-DISK

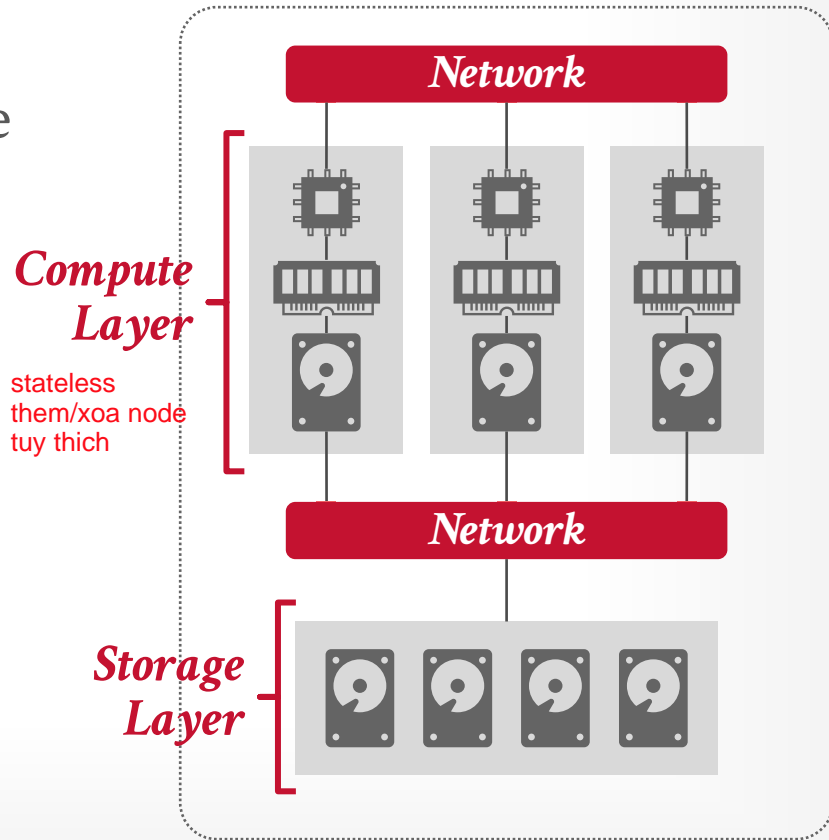Each node accesses a single logical disk via an interconnect, but also have their own private memory and ephemeral storage.
→ Must send messages between nodes to learn about their current state.

Instead of a POSIX API, the DBMS accesses disk using a userspace API.



*Network*

*Compute Layer*

stateless
them/xoa node
tuy thich

Network

*Storage Layer*

tmp

# SYSTEM ARCHITECTURE

**Choice #1: Shared-Nothing:**
→ <mark>Harder to scale capacity due to data movement.</mark>
→ Potentially better performance & efficiency.
→ Apply filters where the data resides before transferring.

**Choice #2: Shared-Disk:**
→ Scale compute layer independently from the storage layer.
→ Easy to shutdown idle compute layer resources.
→ May need to pull uncached persistent data from storage
   layer to compute layer before applying filters.

# SHARED-DISK IMPLEMENTATIONS

Traditionally the storage layer in shared-disk DBMSs were dedicated on-prem NAS.
→ Example: Oracle Exadata

Cloud **object stores** are now the prevailing storage target for modern OLAP DBMSs because they are "infinitely" scalable.
→ Examples: Amazon S3, Azure Blob, Google Cloud Storage

# SHARED-DISK IMPLEMENTATIONS

Traditionally the storag
DBMSs were dedicated
→ Example: Oracle Exadata

Cloud **object stores** ar
target for modern OLA
"infinitely" scalable.
→ Examples: Amazon S3,

# OBJECT STORES

Partition the database's tables (persistent data) into large, immutable files stored in an object store.
→ All attributes for a tuple are stored in the same file in a columnar layout (PAX).
→ Header (or footer) contains meta-data about columnar offsets, compression schemes, indexes, and zone maps.

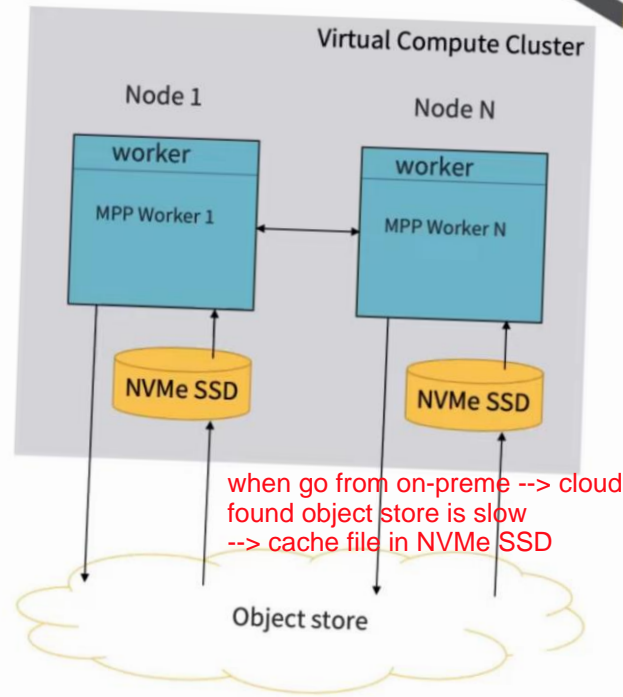The DBMS retrieves a block's header to determine what byte ranges it needs to retrieve (if any).

Each cloud vendor provides their own proprietary API to access data (**PUT**, **GET**, **DELETE**).

# OBJECT STORES

Partit
large,
→ All
    col
→ He
    off

The
what

Each
API



## Workers

- Separated compute / storage
- One Worker pod per compute node
  - Executes portions of the query plan
- Custom network protocol over UDP
  - Data distribution between workers
  - Uses Intel DPDK
  - 50% higher throughput on AWS over TCP/IP
- Shard files cached in local NVMe SSD
- Shards persisted in object store
  - Custom AWS S3 access library
  - 3X better throughput than stock S3 lib

Yellowbrick

**Virtual Compute Cluster**

Node 1 — worker — MPP Worker 1 — NVMe SSD

Node N — worker — MPP Worker N — NVMe SSD

when go from on-preme --> cloud
found object store is slow
--> cache file in NVMe SSD

Object store

# CONCLUSION

Today was about understanding the high-level
context of what modern OLAP DBMSs look like.
→ Fundamentally these new DBMSs are not different than
previous distributed/parallel DBMSs except for the
prevalence of a cloud-based object store for shared disk.

Our focus for the rest of the semester will be about
state-of-the-art implementations of these systems'
components.

# NEXT CLASS

Storage Models

Data Representation

Encoding

Compression