# Database Systems

## Distributed OLAP Databases

# ADMINISTRIVIA

**DBMS Potpourri Lecture** on Wednesday Dec 4th

**Project #4** is due Sunday Dec 8th @ 11:59pm

**Homework #6** is due Monday Dec 9th @ 11:59pm

**Final Exam** is on Friday Dec 13th @ 8:30am
→ Early exam will <u>not</u> be offered.
→ Study guide will be released tomorrow.

# UPCOMING DATABASE TALKS

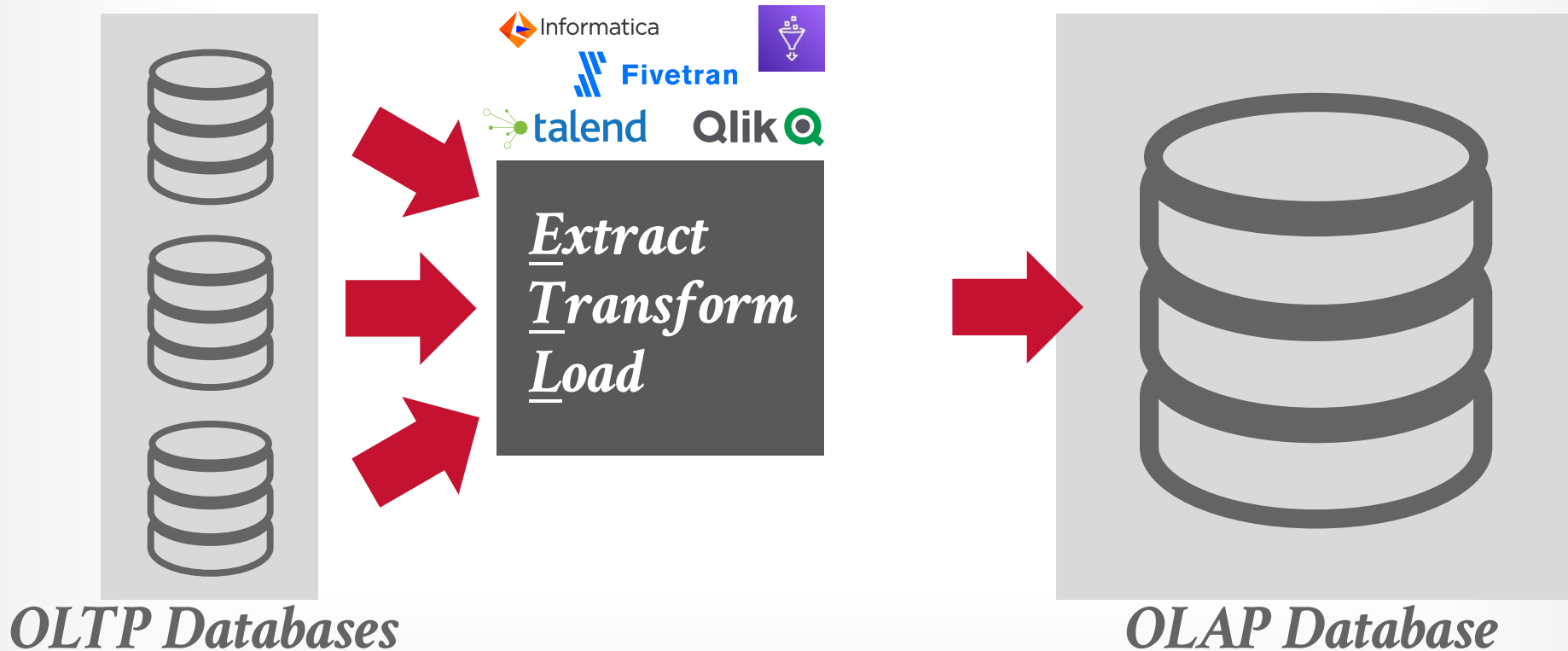**OpenDAL / DataBend** (DB Seminar)
→ Monday Dec 2nd @ 4:30pm
→ Zoom



**GreptimeDB** (DB Seminar)
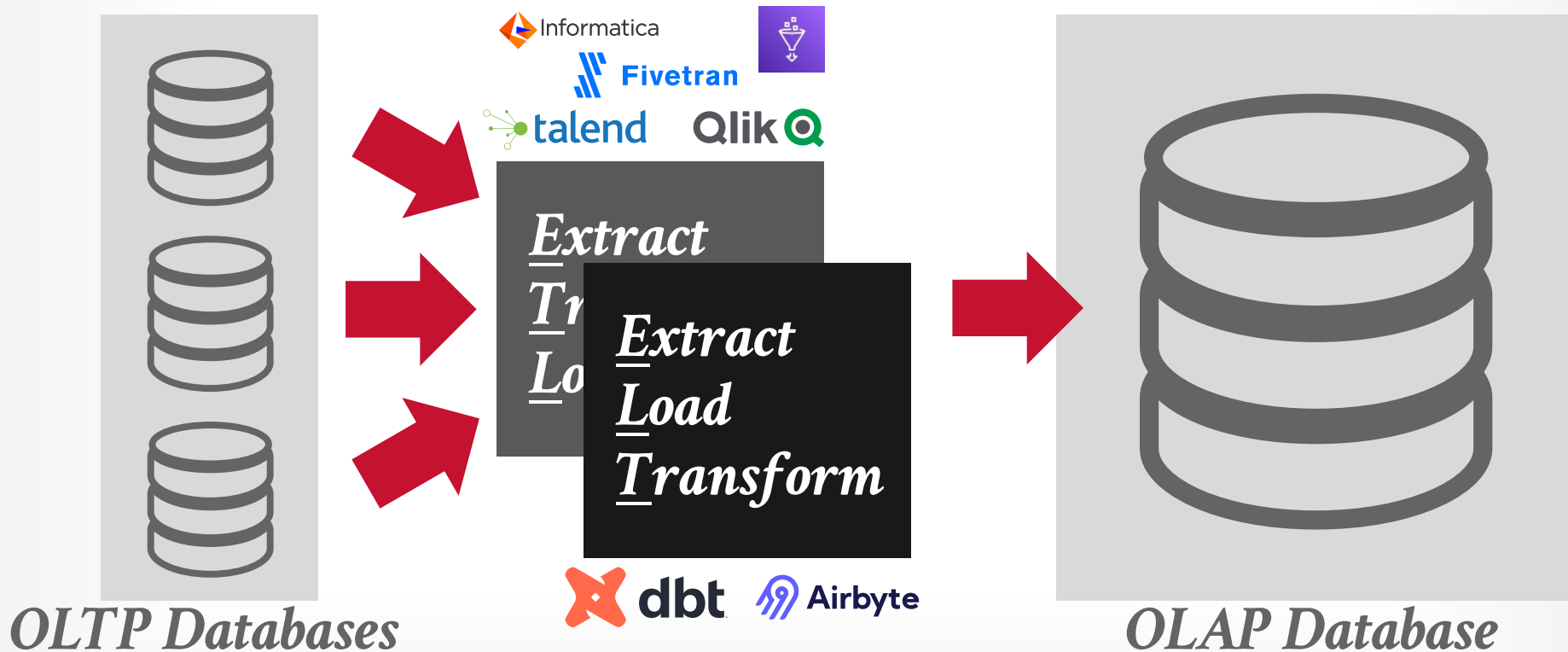→ Monday Dec 9th @ 4:30pm
→ Zoom
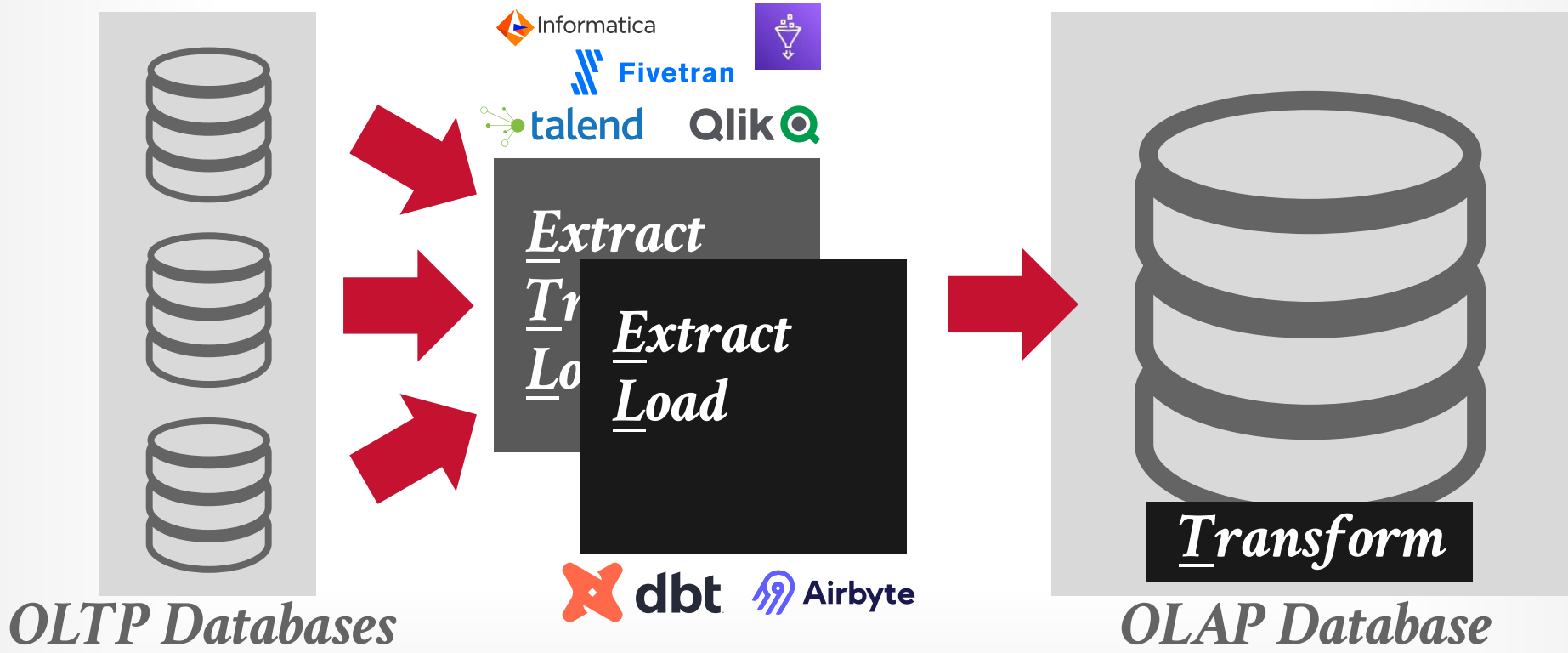
# BIFURCATED ENVIRONMENT



*OLTP Databases*

*OLAP Database*

# BIFURCATED ENVIRONMENT



*OLTP Databases*

*Extract*
*Transform*
*Load*

*Extract*
*Load*
*Transform*

*OLAP Database*

# BIFURCATED ENVIRONMENT



*OLTP Databases*

*Extract*
*Transform*
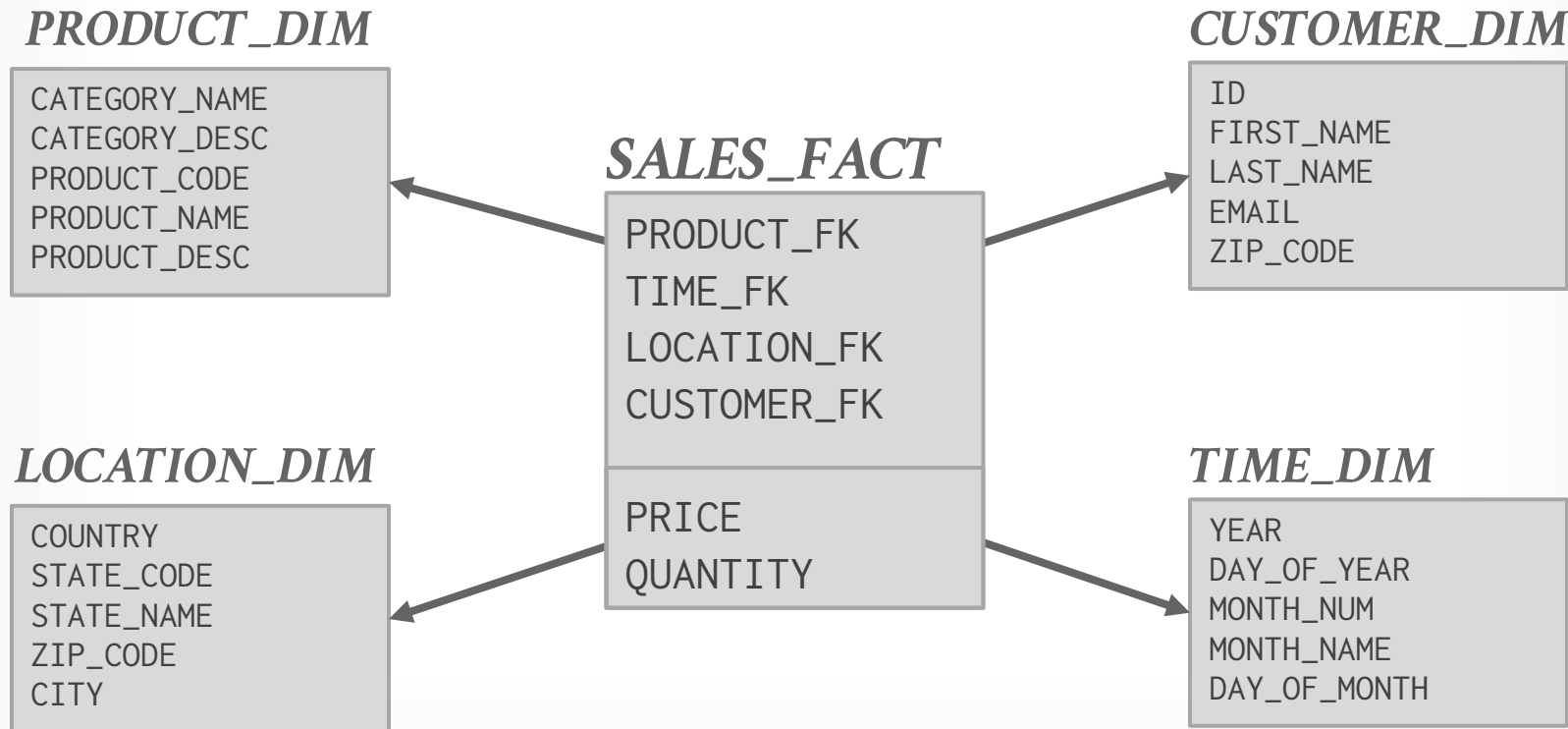*Load*

*Extract*
*Load*

*Transform*

*OLAP Database*

# DECISION SUPPORT SYSTEMS

Applications that serve the management, operations, and planning levels of an organization to help people make decisions about future issues and problems by analyzing historical data.

**Star Schema vs. Snowflake Schema**

# STAR SCHEMA

**PRODUCT_DIM**

CATEGORY_NAME
CATEGORY_DESC
PRODUCT_CODE
PRODUCT_NAME
PRODUCT_DESC

**SALES_FACT**

PRODUCT_FK
TIME_FK
LOCATION_FK
CUSTOMER_FK

PRICE
QUANTITY

**CUSTOMER_DIM**

ID
FIRST_NAME
LAST_NAME
EMAIL
ZIP_CODE

**LOCATION_DIM**

COUNTRY
STATE_CODE
STATE_NAME
ZIP_CODE
CITY

**TIME_DIM**

YEAR
DAY_OF_YEAR
MONTH_NUM
MONTH_NAME
DAY_OF_MONTH

# SNOWFLAKE SCHEMA

**CAT_LOOKUP**

| |
|---|
| CATEGORY_ID |
| CATEGORY_NAME |
| CATEGORY_DESC |

**PRODUCT_DIM**

| |
|---|
| CATEGORY_FK |
| PRODUCT_CODE |
| PRODUCT_NAME |
| PRODUCT_DESC |

**SALES_FACT**

| |
|---|
| PRODUCT_FK |
| TIME_FK |
| LOCATION_FK |
| CUSTOMER_FK |
| |
| PRICE |
| QUANTITY |

**CUSTOMER_DIM**

| |
|---|
| ID |
| FIRST_NAME |
| LAST_NAME |
| EMAIL |
| ZIP_CODE |

**LOCATION_DIM**

| |
|---|
| COUNTRY |
| STATE_FK |
| ZIP_CODE |
| CITY |

**TIME_DIM**

| |
|---|
| YEAR |
| DAY_OF_YEAR |
| MONTH_FK |
| DAY_OF_MONTH |

**STATE_LOOKUP**

| |
|---|
| STATE_ID |
| STATE_CODE |
| STATE_NAME |

**MONTH_LOOKUP**

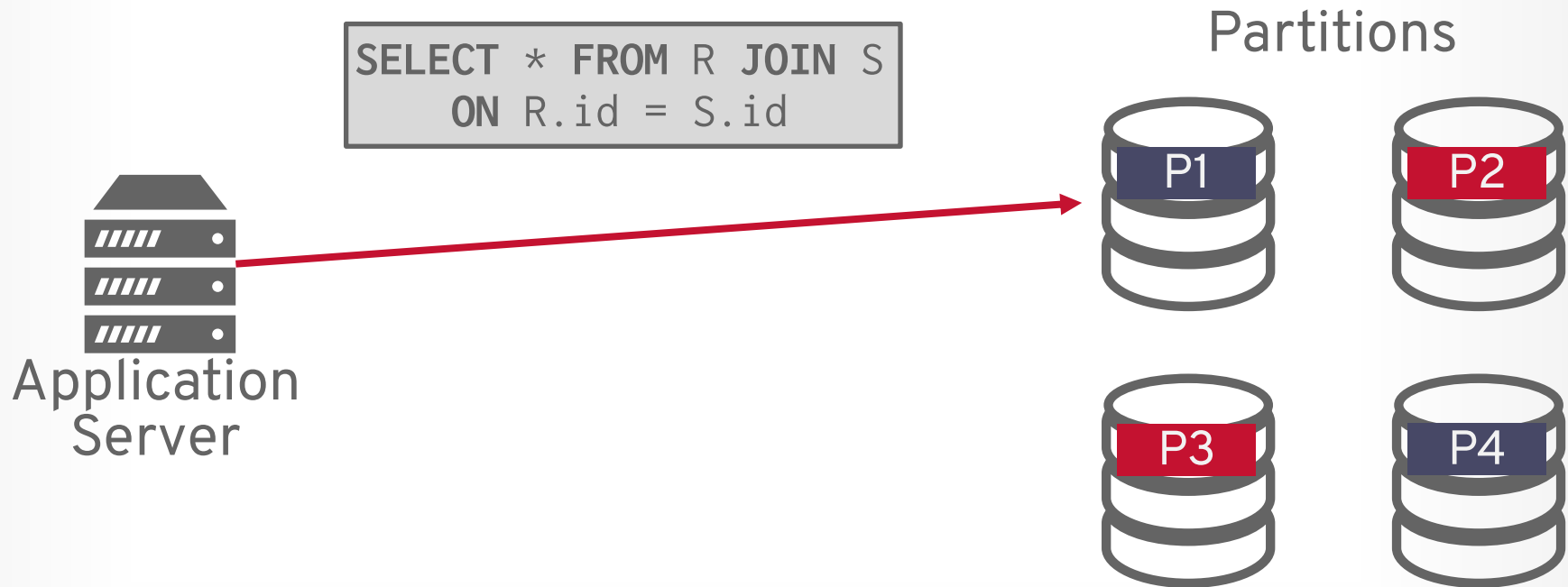| |
|---|
| MONTH_NUM |
| MONTH_NAME |
| MONTH_SEASON |

# STAR VS. SNOWFLAKE SCHEMA

**Issue #1: Normalization**
→ Snowflake schemas take up less storage space.
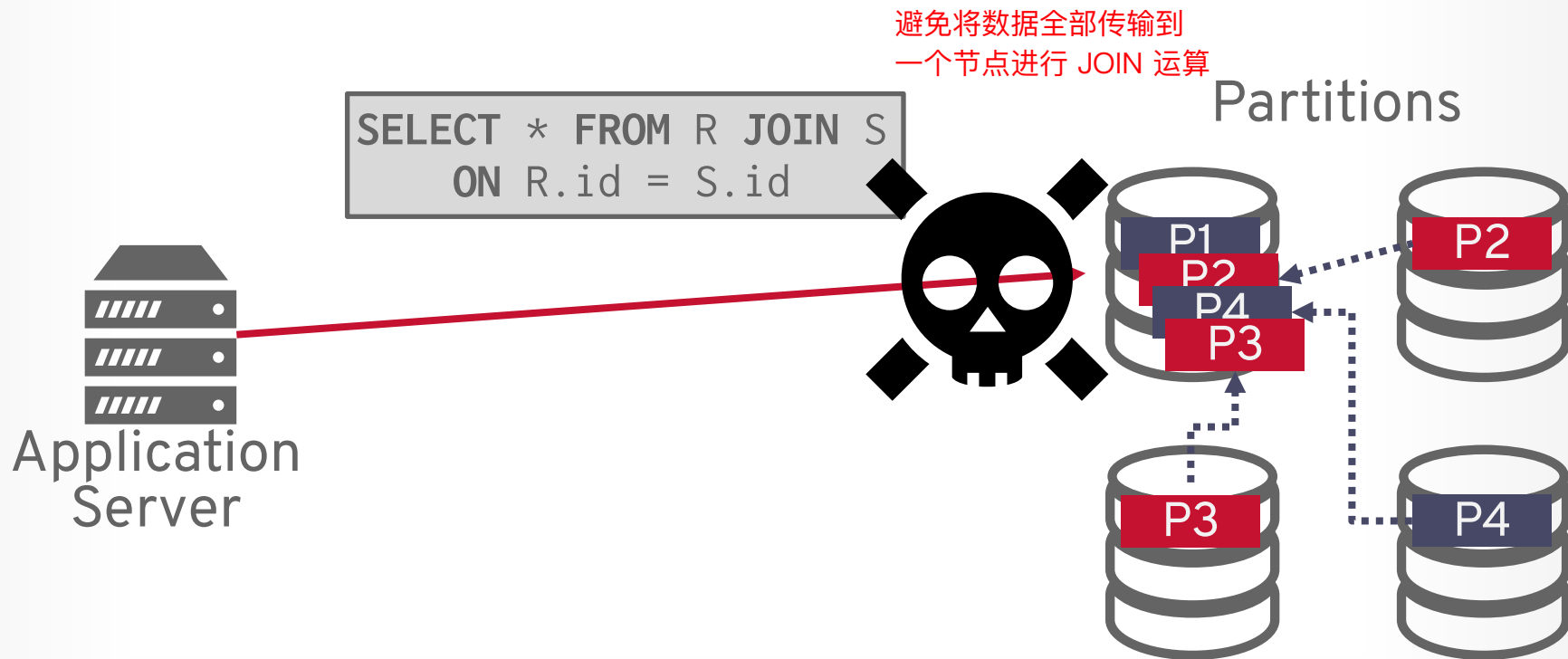→ Denormalized data models may incur integrity and consistency violations.

**Issue #2: Query Complexity**
→ Snowflake schemas require more joins to get the data needed for a query.
→ Queries on star schemas will (usually) be faster.

# PROBLEM SETUP

```
SELECT * FROM R JOIN S
    ON R.id = S.id
```

Partitions



Application
Server

P1  P2

P3  P4

# PROBLEM SETUP

避免将数据全部传输到
一个节点进行 JOIN 运算

Partitions

```
SELECT * FROM R JOIN S
    ON R.id = S.id
```

Application
Server

P1
P2
P4
P3

P2

P3

P4

# TODAY'S AGENDA

Execution Models

Query Planning

Distributed Join Algorithms
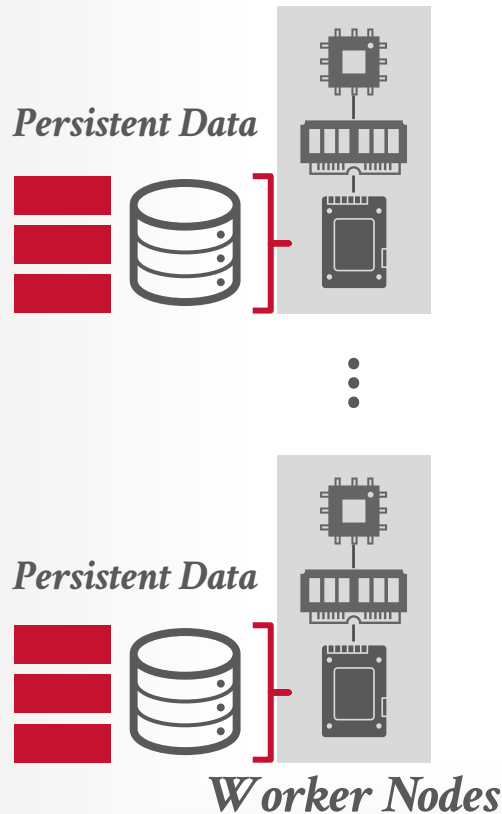
Cloud Systems

# DISTRIBUTED QUERY EXECUTION

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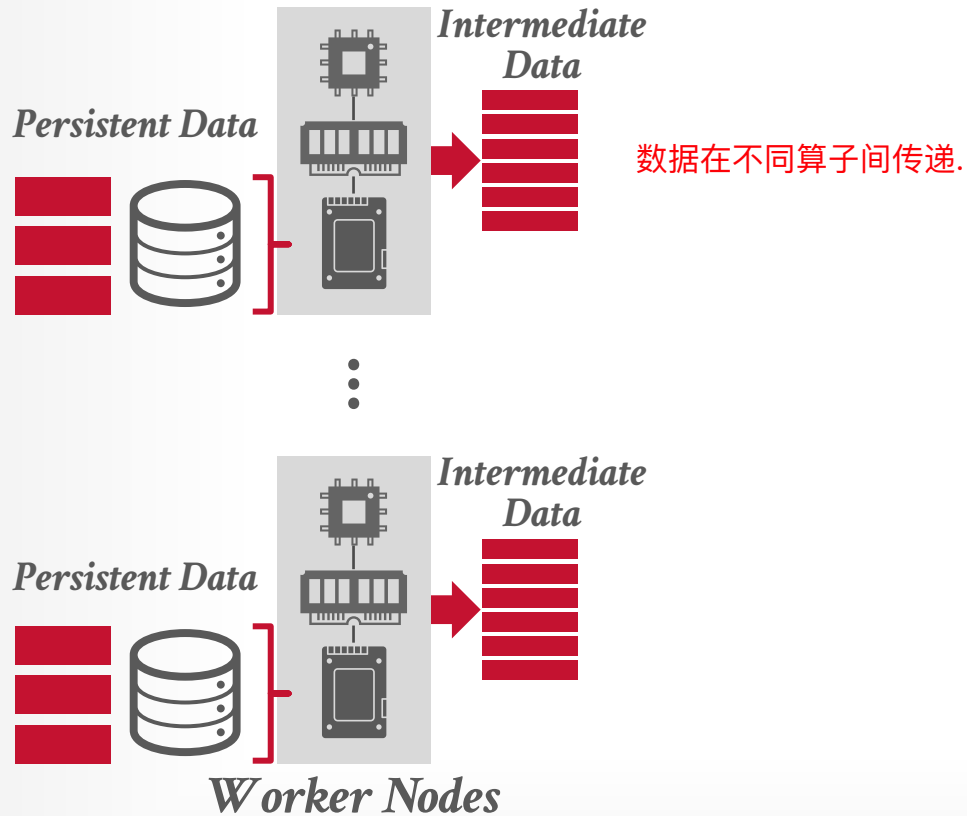→ Query plan is a DAG of physical operators.

For each operator, the DBMS considers where input is coming from and where to send output.

根据 pipeline 分而治之

→ Table Scans
→ Joins
→ Aggregations
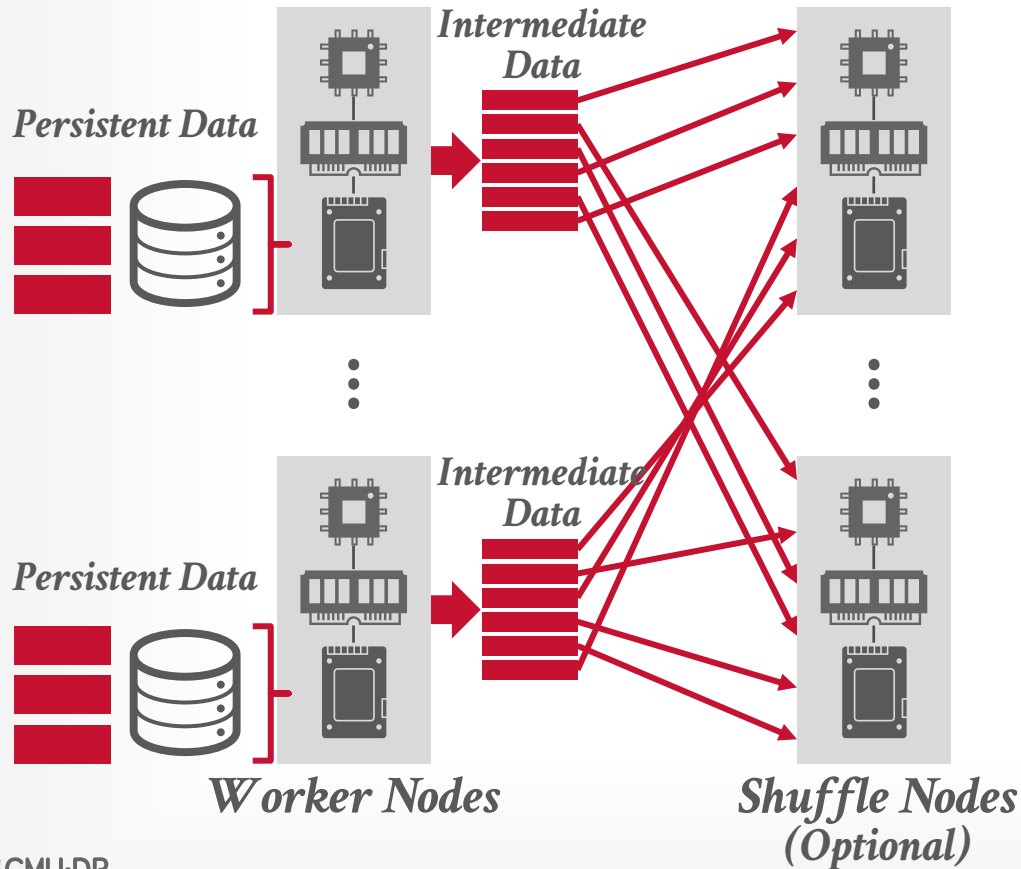→ Sorting

# DISTRIBUTED QUERY EXECUTION



*Persistent Data*

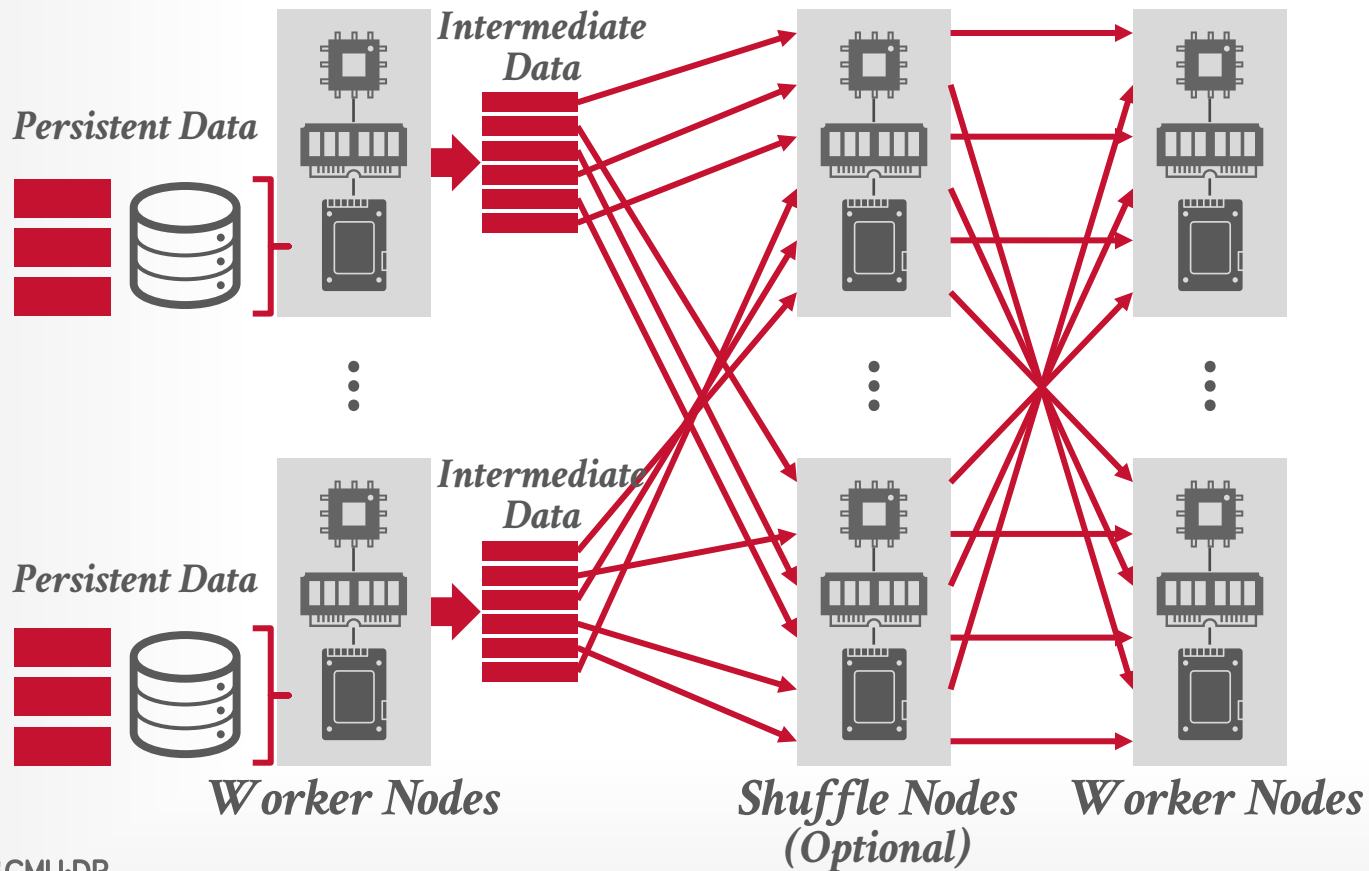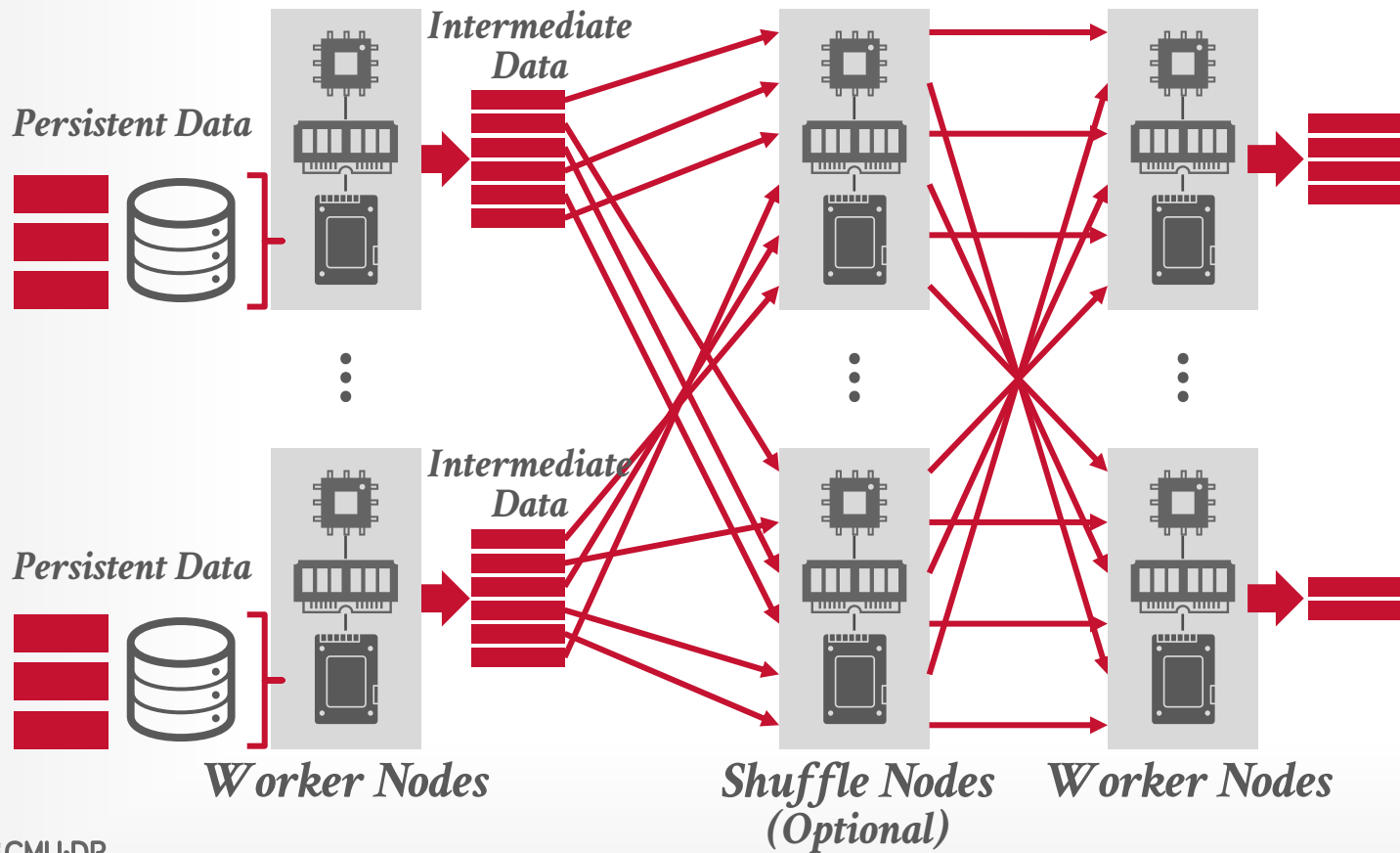*Persistent Data*

*Worker Nodes*

# DISTRIBUTED QUERY EXECUTION



数据在不同算子间传递.

# DISTRIBUTED QUERY EXECUTION



*Persistent Data*

*Intermediate Data*

*Intermediate Data*

*Persistent Data*

*Worker Nodes*

*Shuffle Nodes (Optional)*

# DISTRIBUTED QUERY EXECUTION



Persistent Data

Intermediate Data

Intermediate Data

Persistent Data

Worker Nodes

Shuffle Nodes (Optional)

Worker Nodes

# DISTRIBUTED QUERY EXECUTION



*Persistent Data*

*Intermediate Data*

*Worker Nodes*

*Shuffle Nodes (Optional)*

*Worker Nodes*

# DISTRIBUTED QUERY EXECUTION



Final
Result

最后发送到一个节点将结
果全部汇总返回给客户端.

# 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 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**

# 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 database files are located.

# PUSH VS. PULL

**Approac**
→ Send th
contain
→ Perfor
data re

**Approa**
→ Bring
needs it for processing.
→ This is necessary when there is no compute resources
available where database files are located.

---

## Filtering and retrieving data using Amazon S3 Select

**amazon**

PDF | RSS

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.

PUSH vs. PULL

Approac...

# Filtering and retrieving data using Amazon S3 Select

PDF | RSS

With Amazon S3 Select, yo...

...uery language (SQL) statements to filter the contents of an
...at you need. By using Amazon S3 Select to filter this data, you can
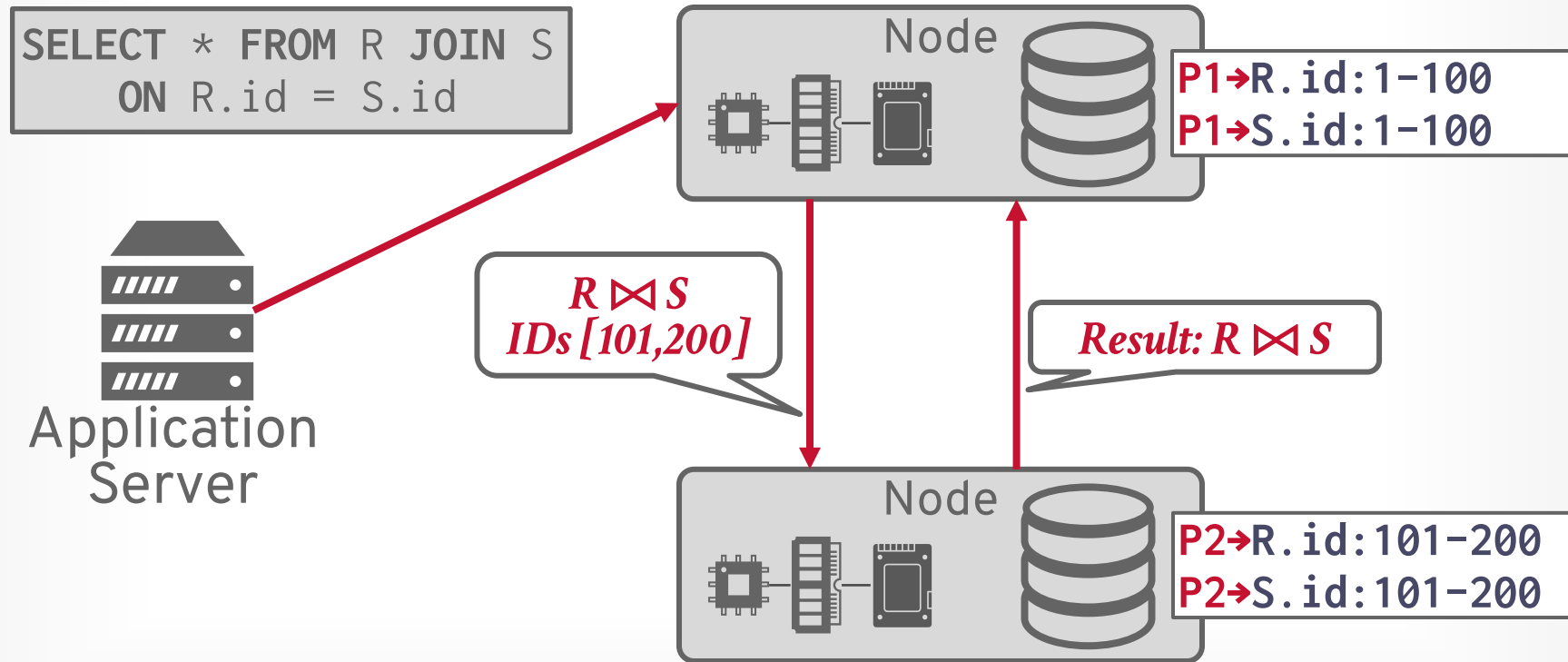...ich reduces the cost and latency to retrieve this data.

...or Apache Parquet format. It also works with objects that are
...s only), and server-side encrypted objects. You can specify the
...etermine how the records in the result are delimited.

...azon S3 Select supports a subset of SQL. For more information
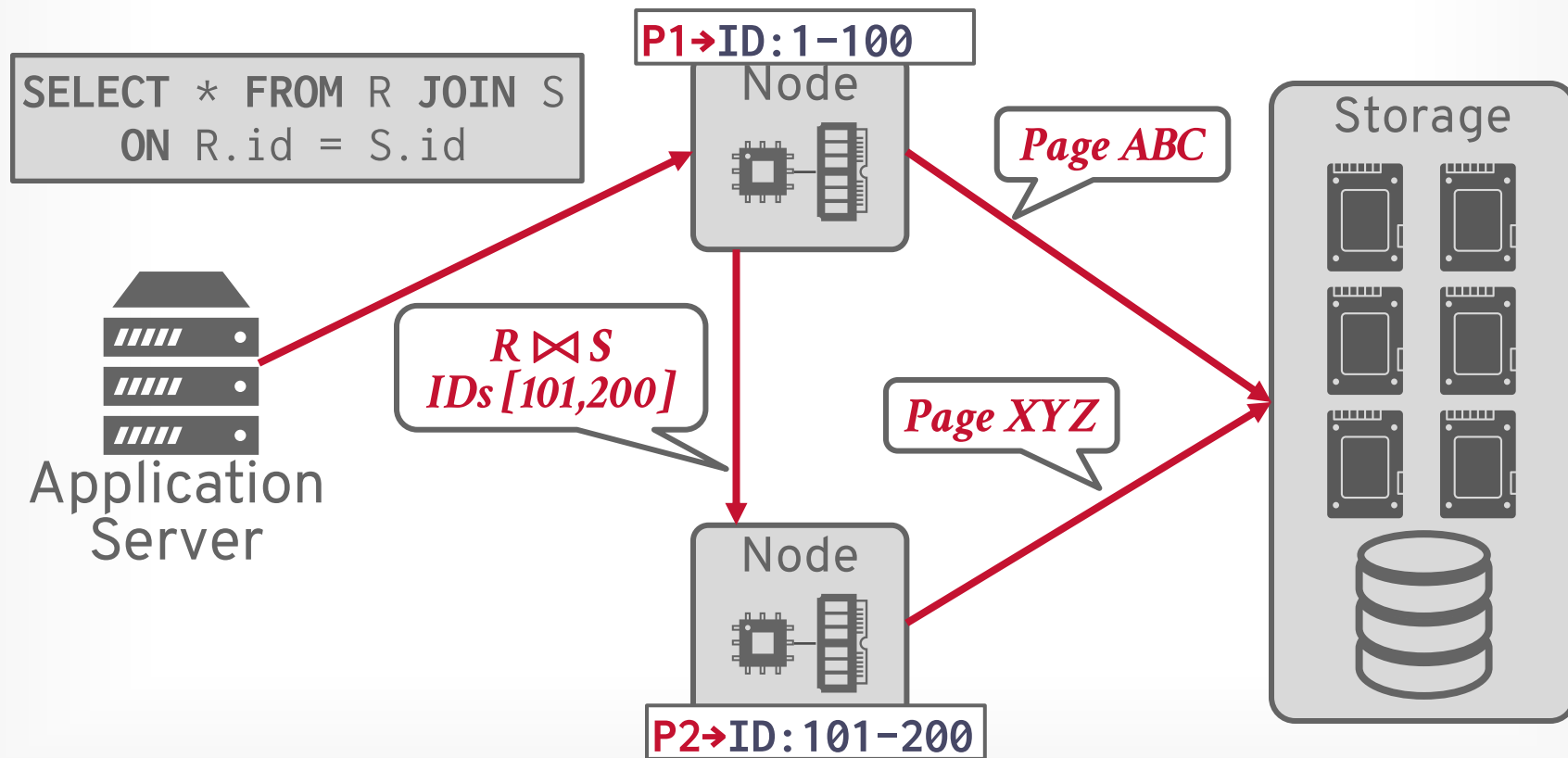...Select, see SQL reference for Amazon S3 Select.

...Object Content REST API, the AWS Command Line Interface
...e limits the amount of data returned to 40 MB. To retrieve

## Query Blob Contents

Microsoft

👍 Feedback

Article • 07/20/2021 • 10 minutes to read • 3 contributors

The `Query Blob Contents` API applies a simple Structured Query Language (SQL) statement on a blob's
contents and returns only the queried subset of the data. You can also call `Query Blob Contents` to query
the contents of a version or snapshot.

## Request

The `Query Blob Contents` request may be constructed as follows. HTTPS is recommended. Replace
*myaccount* with the name of your storage account:

| POST Method Request URI | HTTP Version |
|---|---|
| https://myaccount.blob.core.windows.net/mycontainer/myblob?comp=query | HTTP/1.0 |
| https://myaccount.blob.core.windows.net/mycontainer/myblob?comp=query&snapshot=<DateTime> | HTTP/1.1 |
| https://myaccount.blob.core.windows.net/mycontainer/myblob?comp=query&versionid=<DateTime> | |

...mpute resources
...ed.

# PUSH QUERY TO DATA

```
SELECT * FROM R JOIN S
     ON R.id = S.id
```

Node

**P1➔R.id:1-100**
**P1➔S.id:1-100**

Application
Server

$R \bowtie S$
$IDs\,[101,200]$

$Result: R \bowtie S$

Node

**P2➔R.id:101-200**
**P2➔S.id:101-200**

# PULL DATA TO QUERY

```
SELECT * FROM R JOIN S
    ON R.id = S.id
```

P1➔ID:1-100

Node

Storage

*Page ABC*

*R ⋈ S
IDs [101,200]*

*Page XYZ*

Application
Server

Node

P2➔ID:101-200

# PULL DATA TO QUERY

# PULL DATA TO QUERY



`P1→ID:1-100`

Node

```
SELECT * FROM R JOIN S
       ON R.id = S.id
```

$R \bowtie S$
*IDs [101,200]*

*Result: $R \bowtie S$*

Application
Server

Node

`P2→ID:101-200`

Storage

# OBSERVATION

The data that a node receives from remote sources are cached in the buffer pool.
→ This allows the DBMS to support intermediate results that are large than the amount of memory available.
→ Ephemeral pages are <u>not</u> persisted after a restart.

What happens to a long-running OLAP query if a node crashes during execution?

# QUERY FAULT TOLERANCE

Most shared-nothing distributed OLAP DBMSs are designed to assume that nodes do not fail during query execution.
→ If one node fails during query execution, then the whole query fails.

The DBMS could take a <u>snapshot</u> of the 为查询的中间结果提供故障容错机制. intermediate results for a query during execution to allow it to recover if nodes fail.

# QUERY FAULT TOLERANCE



```
SELECT * FROM R JOIN S
    ON R.id = S.id
```

Node

Application
Server

$R \bowtie S$

$Result: R \bowtie S$

Storage

Node

将中间结果刷盘.

# QUERY FAULT TOLERANCE



```
SELECT * FROM R JOIN S
    ON R.id = S.id
```

Node

*Result: $R \bowtie S$*

Storage

Application Server

当该节点宕机后，协调器将检索存储介质是否已保存计算结果，如果 yes，则获取该结果并继续执行后续的查询计划.

Node

# QUERY PLANNING

All the optimizations that we talked about before are still applicable in a distributed environment.
→ Predicate Pushdown
→ Projection Pushdown
→ Optimal Join Orderings

Distributed query optimization is even harder because it must consider the physical location of data and network transfer costs. 成本估算模型.

# QUERY PLAN FRAGMENTS

## Approach #1: Physical Operators
→ Generate a single query plan and then break it up into partition-specific fragments.
→ Most systems implement this approach.

## Approach #2: SQL
→ Rewrite original query into partition-specific queries.
→ Allows for local optimization at each node.
→ SingleStore + Vitess are the only systems we know that use this approach.

# QUERY PLAN FRAGMENTS

```
SELECT * FROM R JOIN S
     ON R.id = S.id
```

```
SELECT * FROM R JOIN S
     ON R.id = S.id
 WHERE R.id BETWEEN 1 AND 100
```

```
SELECT * FROM R JOIN S
     ON R.id = S.id
 WHERE R.id BETWEEN 101 AND 200
```

```
SELECT * FROM R JOIN S
     ON R.id = S.id
 WHERE R.id BETWEEN 201 AND 300
```

id:1-100

id:101-200

id:201-300

# QUERY PLAN FRAGMENTS

> *Union the output of each join to produce final result.*

```
FROM R JOIN S
  ON R.id = S.id
```

```
SELECT * FROM R JOIN S
    ON R.id = S.id
 WHERE R.id BETWEEN 1 AND 100
```

```
SELECT * FROM R JOIN S
    ON R.id = S.id
 WHERE R.id BETWEEN 101 AND 200
```

```
SELECT * FROM R JOIN S
    ON R.id = S.id
 WHERE R.id BETWEEN 201 AND 300
```

id:1-100

id:101-200

id:201-300

# OBSERVATION

The efficiency of a distributed join depends on the target tables' partitioning schemes.

One approach is to put entire tables on a single node and then perform the join.
→ You lose the parallelism of a distributed DBMS.
→ Costly data transfer over the network.

# DISTRIBUTED JOIN ALGORITHMS

To join tables **R** and **S**, the DBMS needs to get the proper tuples on the same node.

Once the data is at the node, the DBMS then executes the same join algorithms that we discussed earlier in the semester.
→ Need to produce the correct answer as if all the data is located in a single node system.

# SCENARIO #1

The entire copy of one data set is
replicated at every node.
→ Think of it as a small dimension table.

```
SELECT * FROM R JOIN S
    ON R.id = S.id
```

Each node joins its local data in
parallel and then sends their results to
a coordinating node.

*Partition Key*

表 R 根据 JOIN 键将数据分区到不同节点.

id:1-100

R{Id}

*Replicated*

S

R{Id}

id:101-200

S

*Replicated*

# SCENARIO #1

The entire copy of one data set is replicated at every node.
→ Think of it as a small dimension table.

```
SELECT * FROM R JOIN S
    ON R.id = S.id
```

Each node joins its local data in parallel and then sends their results to a coordinating node. 在每个节点上表 R 分别和 S 进行 JOIN 计算输出结果.



**P1:R⋈S**   **P2:R⋈S**

*Partition Key*

id:1-100   R{Id}   R{Id}   id:101-200

*Replicated*   S   S   *Replicated*

小表 S 复制到每个节点.

# SCENARIO #1

The entire copy of one data set is replicated at every node.
→ Think of it as a small dimension table.

Each node joins its local data in parallel and then sends their results to a coordinating node.

```
SELECT * FROM R JOIN S
    ON R.id = S.id
```

# SCENARIO #2

Both data sets are partitioned on the join attribute. Each node performs the join on local data and then sends to a coordinator node for coalescing.

```
SELECT * FROM R JOIN S
    ON R.id = S.id
```

# SCENARIO #2

Both data sets are partitioned on the join attribute. Each node performs the join on local data and then sends to a coordinator node for coalescing.

```
SELECT * FROM R JOIN S
    ON R.id = S.id
```

# SCENARIO #3 – BROADCAST JOIN

Both data sets are partitioned on different keys. If one of the data sets is small, then the DBMS "broadcasts" that data to all nodes. 两个表按不同的 key 进行数据分区.

```
SELECT * FROM R JOIN S
    ON R.id = S.id
```

例如，小表 S 按非连接键 val 列进行数据分区. 将一个节点上 S 的部分数据广播到
另一个节点，使得这个节点拥有表 S 的全部数据，然后再将 S 的全部数据复制到
其他节点进行 JOIN 运算.

# SCENARIO #3 – BROADCAST JOIN

Both data sets are partitioned on different keys. If one of the data sets is small, then the DBMS "broadcasts" that data to all nodes.

```
SELECT * FROM R JOIN S
    ON R.id = S.id
```

# SCENARIO #3 – BROADCAST JOIN

Both data sets are partitioned on different keys. If one of the data sets is small, then the DBMS "broadcasts" that data to all nodes.

```
SELECT * FROM R JOIN S
    ON R.id = S.id
```



P1:R⋈S

S

id:1-100

R{id}

val:1-50

S{val}

P2:R⋈S

S

R{id}

id:101-200

S{val}

val:51-100

# SCENARIO #3 – BROADCAST JOIN

Both data sets are partitioned on different keys. If one of the data sets is small, then the DBMS "broadcasts" that data to all nodes.

```
SELECT * FROM R JOIN S
    ON R.id = S.id
```
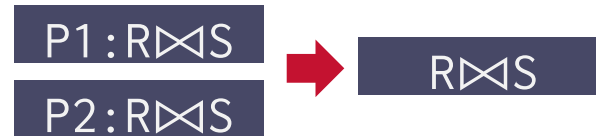
# SCENARIO #4 – <u>SHUFFLE JOIN</u>

Both data sets are <u>not</u> partitioned on the join key. The DBMS copies/re-partitions the data on-the-fly across nodes.
→ The repartitioned data copy is generally deleted when the query is done.

```
SELECT * FROM R JOIN S
    ON R.id = S.id
```
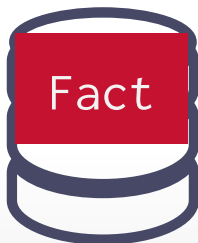
| name:A-M | R{name} |
| val:1-50 | S{val} |

| R{name} | name:N-Z |
| S{val} | val:51-100 |

# SCENARIO #4 – SHUFFLE JOIN

Both data sets are <u>not</u> partitioned on the join key. The DBMS copies/re-partitions the data on-the-fly across nodes.
→ The repartitioned data copy is generally deleted when the query is done.

```
SELECT * FROM R JOIN S
    ON R.id = S.id
```

SHUFFLE：表 R 首先根据所有数据的 id 列进行 hash 操作，然后进行范围分区，从而将表 R 数据重新分区.

R{id}          id:101-200

name:A-M    R{name}        R{name}        name:N-Z

val:1-50    S{val}         S{val}         val:51-100

# SCENARIO #4 – SHUFFLE JOIN

Both data sets are <u>not</u> partitioned on the join key. The DBMS copies/re-partitions the data on-the-fly across nodes.

→ The repartitioned data copy is generally deleted when the query is done.

```
SELECT * FROM R JOIN S
    ON R.id = S.id
```

# SCENARIO #4 – SHUFFLE JOIN

Both data sets are <u>not</u> partitioned on the join key. The DBMS copies/re-partitions the data on-the-fly across nodes.

→ The repartitioned data copy is generally deleted when the query is done.

```
SELECT * FROM R JOIN S
    ON R.id = S.id
```



| id:1-100 | R{id} | | R{id} | id:101-200 |

同理表 S 也进行同样操作.

| name:A-M | R{name} | | R{name} | name:N-Z |
| val:1-50 | S{val} | | S{val} | val:51-100 |

S{id} — id:101-200

# SCENARIO #4 – SHUFFLE JOIN

Both data sets are <u>not</u> partitioned on the join key. The DBMS copies/re-partitions the data on-the-fly across nodes.

→ The repartitioned data copy is generally deleted when the query is done.

```
SELECT * FROM R JOIN S
    ON R.id = S.id
```

| id:1-100 | R{id} | | R{id} | id:101-200 |
| id:1-100 | S{id} | | S{id} | id:101-200 |
| name:A-M | R{name} | | R{name} | name:N-Z |
| val:1-50 | S{val} | | S{val} | val:51-100 |

# SCENARIO #4 – SHUFFLE JOIN

Both data sets are <u>not</u> partitioned on the join key. The DBMS copies/re-partitions the data on-the-fly across nodes.

→ The repartitioned data copy is generally deleted when the query is done.

```
SELECT * FROM R JOIN S
    ON R.id = S.id
```

P1:R⋈S

| id:1-100 | R{id} |
| id:1-100 | S{id} |
| name:A-M | R{name} |
| val:1-50 | S{val} |

P2:R⋈S

| R{id} | id:101-200 |
| S{id} | id:101-200 |
| R{name} | name:N-Z |
| S{val} | val:51-100 |

# SCENARIO #4 – SHUFFLE JOIN

Both data sets are <u>not</u> partitioned on the join key. The DBMS copies/re-partitions the data on-the-fly across nodes.

→ The repartitioned data copy is generally deleted when the query is done.

```
SELECT * FROM R JOIN S
    ON R.id = S.id
```

P1:R⋈S  ➡  R⋈S
P2:R⋈S

| id:1-100 | R{id} | | R{id} | id:101-200 |
|---|---|---|---|---|
| id:1-100 | S{id} | | S{id} | id:101-200 |
| name:A-M | R{name} | | R{name} | name:N-Z |
| val:1-50 | S{val} | | S{val} | val:51-100 |

# SEMI-JOIN OPTIMIZATION

Before pulling data from another node, send a **semi-join filter** to reduce data movement.

→ Perform a join on the bare minimum data needed to avoid unnecessary transfers.
→ Could use an approximate filter (Bloom Join).

```
SELECT Fact.price, Dim.*
  FROM Fact JOIN Dim
    ON Fact.id = Dim.id
 WHERE Dim.zip = 15213
```

# SEMI-JOIN OPTIMIZATION

Before pulling data from another node, send a **semi-join filter** to reduce data movement. 减少数据量的移动.

→ Perform a join on the bare minimum data needed to avoid unnecessary transfers.
→ Could use an approximate filter (Bloom Join).

```
SELECT Fact.price, Dim.*
  FROM Fact JOIN Dim
    ON Fact.id = Dim.id
 WHERE Dim.zip = 15213
```

首先对表进行过滤和投影（仅连接键），获得它的最小形式的数据.

$$\mathbf{Dim_{semi}} = \mathbf{\Pi_{id}} (\mathbf{\sigma_{zip = 15213}} \mathbf{Dim})$$

# SEMI-JOIN OPTIMIZATION

Before pulling data from another node, send a **semi-join filter** to reduce data movement.

→ Perform a join on the bare minimum data needed to avoid unnecessary transfers.

→ Could use an approximate filter (Bloom Join).

```
SELECT Fact.price, Dim.*
  FROM Fact JOIN Dim
    ON Fact.id = Dim.id
 WHERE Dim.zip = 15213
```

将过滤后的 Dim 表数据发送到表 Fact 位于的节点进行 JOIN 计算.

$$Dim_{semi} = \Pi_{id} \left( \sigma_{zip = 15213} \ Dim \right)$$

$$F\text{-small} = Fact \bowtie Dim_{semi}$$

Fact

$Dim_{semi}$

Dim

# SEMI-JOIN OPTIMIZATION

Before pulling data from another node, send a **semi-join filter** to reduce data movement.
→ Perform a join on the bare minimum data needed to avoid unnecessary transfers.
→ Could use an approximate filter (Bloom Join).

```
SELECT Fact.price, Dim.*
  FROM Fact JOIN Dim
    ON Fact.id = Dim.id
 WHERE Dim.zip = 15213
```

$$\textbf{Dim}_{semi} = \Pi_{id} (\sigma_{zip = 15213} \textbf{Dim})$$

$$\textbf{F-small} = \textbf{Fact} \bowtie \textbf{Dim}_{semi}$$

Fact$_{small}$

Fact

Dim$_{semi}$

Dim

# SEMI-JOIN OPTIMIZATION

Before pulling data from another node, send a **semi-join filter** to reduce data movement.

→ Perform a join on the bare minimum data needed to avoid unnecessary transfers.

→ Could use an approximate filter (Bloom Join).

```
SELECT Fact.price, Dim.*
  FROM Fact JOIN Dim
    ON Fact.id = Dim.id
 WHERE Dim.zip = 15213
```

$$\textbf{Result} = \mathbf{\Pi_{price}(Dim \bowtie Fact_{small})}$$

根据过滤后的表 Fact 数据再次进行 JOIN 计算.

Fact

Dim

Fact<sub>small</sub>

# OBSERVATION

Direct communication between compute nodes means the DBMS knows which nodes will participate in query execution ahead of time.

But data skew can cause imbalances…

A better approach is to dynamically adjust compute resources on the fly as a query executes.

# SHUFFLE PHASE

Redistribute of intermediate data across nodes between query plan pipelines. pipeline breaker.
→ Can repartition / rebalance data based on observed data characteristics.

使用专门的硬件进行 SHUFFLE.

Some DBMSs support standalone fault-tolerant shuffle services.
→ Example: You can replace Spark's built-in in-memory shuffle implementation or replace it with a separate service.

# SHUFFLE PHASE



$hash_1(key) \% n$

Shuffle Nodes

Worker — Consumer / Producer

Stage n

Shared-Disk

# SHUFFLE PHASE



Shuffle Nodes

$hash_1(key) \% n$

Stage n

Shared-Disk

Stage n+1

# SHUFFLE PHASE



**Shuffle Nodes**

$hash_1(key) \% n$

*Stage n*

*Stage n+1*

**Shared-Disk**

# SHUFFLE PHASE



Shuffle Nodes

$hash_1(key) \% n$

Worker
Consumer
Producer

Stage n

Shared-Disk

Stage n+1

# SHUFFLE PHASE

**Consumer** *Worker*

**Consumer** *Worker*

**Consumer** *Worker*

**Consumer** *Worker*

*Stage*

*Shared-Disk*

## EXCHANGE OPERATOR

### Exchange Type #1 – Gather
→ Combine the results from multiple workers into a single output stream.

### Exchange Type #2 – Distribute
→ Split a single input stream into multiple output streams.

### Exchange Type #3 – Repartition
→ Shuffle multiple input streams across multiple output streams.
→ Some DBMSs always perform this step after every pipeline (e.g., Dremel/BigQuery).

Source: Craig Freedman

| | |
| Gather |
| Operator | Operator | Operator |

| Operator | Operator | Operator |
| Distribute |

| Operator | | Operator |
| Repartition |
| Operator | Operator | Operator |

# CLOUD SYSTEMS

Vendors provide **database-as-a-service** (DBaaS) offerings that are managed DBMS environments.

Newer systems are starting to blur the lines between shared-nothing and shared-disk.
→ Example: You can do simple filtering on Amazon S3 before copying data to compute nodes.

# CLOUD SYSTEMS

## Approach #1: Managed DBMSs
→ No significant modification to the DBMS to be "aware" that it is running in a cloud environment.
→ Examples: Most vendors

## Approach #2: Cloud-Native DBMS
→ System designed explicitly to run in a cloud environment.
→ Usually based on a shared-disk architecture.
→ Examples: Snowflake, Google BigQuery

# SERVERLESS DATABASES

Rather than always maintaining compute resources for each customer, a "serverless" DBMS evicts tenants when they become idle.



Application Server

Node

# SERVERLESS DATABASES

Rather than always maintaining compute resources for each customer, a "serverless" DBMS evicts tenants when they become idle.

Application Server

Node

# SERVERLESS DATABASES

Rather than always maintaining compute resources for each customer, a "serverless" DBMS evicts tenants when they become idle.

# SERVERLESS DATABASES

Rather than always maintaining compute resources for each customer, a "serverless" DBMS evicts tenants when they become idle.



Application Server

Node

Storage

# SERVERLESS DATABASES

Rather than always maintaining compute resources for each customer, a "serverless" DBMS evicts tenants when they become idle.



*Buffer Pool Page Table*

Storage

Node

Application Server

# SERVERLESS DATABASES

Rather than always maintaining compute resources for each customer, a "serverless" DBMS evicts tenants when they become idle.

Application
Server

Storage

# SERVERLESS DATABASES

Rather than always maintaining compute resources for each customer, a "serverless" DBMS evicts tenants when they become idle.

planetscale

CockroachDB

NEON

amazon

fauna

Microsoft® SQL Azure™

Application Server

Node

Storage

*Buffer Pool Page Table*

# DATA LAKES

Repository for storing large amounts of structured, semi-structured, and unstructured data without having to define a schema or ingest the data into proprietary internal formats.

# DATA LAKES

Repository for storing large amounts of structured, semi-structured, and unstructured data without having to define a schema or ingest the data into proprietary internal formats.



```
CREATE TABLE foo (...);
```

```
INSERT INTO foo VALUES (...);
```

Node

Storage

# DATA LAKES

Repository for storing large amounts of structured, semi-structured, and unstructured data without having to define a schema or ingest the data into proprietary internal formats.

# DATA LAKES

Repository for storing large amounts of structured, semi-structured, and unstructured data without having to define a schema or ingest the data into proprietary internal formats.

```
SELECT * FROM foo
```

Node

Storage

# DATA LAKES

Repository for storing large amounts of structured, semi-structured, and unstructured data without having to define a schema or ingest the data into proprietary internal formats.

```
SELECT * FROM foo
```

Node

Data Lake

JSON

# DATA LAKES

Repository for storing large amounts of structured, semi-structured, and unstructured data without having to define a schema or ingest the data into proprietary internal formats.

```
SELECT * FROM foo
```

Node

Data Lake

JSON

# 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.

# SYSTEM CATALOGS

A DBMS tracks a database's schema (table, columns) and data files in its catalog.
→ If the DBMS is on the data ingestion path, then it can maintain the catalog incrementally.
→ If an external process adds data files, then it also needs to update the catalog so that the DBMS is aware of them.

Notable implementations:
→ HCatalog
→ Google Data Catalog
→ Amazon Glue Data Catalog
→ Databricks Unity
→ Apache Iceberg

# SYSTEM CATALOG

A DBMS tracks a database's
and data files in its catalog.
→ If the DBMS is on the data i
maintain the catalog increme
→ If an external process adds d
update the catalog so that th

Notable implementations:
→ HCatalog
→ Google Data Catalog
→ Amazon Glue Data Catalog
→ Databricks Unity
→ Apache Iceberg



DEFINITE

June 26, 2024    10 minute read

Why Databricks paid $1B for a 40 person startup (Tabular)

Steven Wang

In case you missed it, earlier this month Databricks acquired Tabular, the company behind the open source project Iceberg, for over $1 billion. The acquisition, which was announced during Snowflake's 2024 Summit conference and amid rumors of Snowflake's interest in purchasing Tabular, caught many by surprise especially since Databricks already offers a competing product, Delta Lake. So, what is Iceberg, how does it compare to Delta Lake, and what does the project's future look like post-acquisition?

# QUERY OPTIMIZERS

<mark>Extendible search engine framework for heuristic- and cost-based query optimization.</mark>
→ DBMS provides transformation rules and cost estimates.
→ Framework returns either a logical or physical query plan.

Notable implementations:
→ Greenplum Orca
→ Apache Calcite

**This is what 15-799 will cover next semester!**

# DATA FILE FORMATS

Most DBMSs use a proprietary on-disk binary file format for their databases.
→ Think of the <u>BusTub</u> page types…

The only way to share data between systems is to convert data into a common text-based format
→ Examples: CSV, JSON, XML

There are new open-source binary file formats that make it easier to access data across systems.

# DATA FILE FORMATS

## Apache Parquet
→ Compressed columnar storage from Cloudera/Twitter

## Apache ORC
→ Compressed columnar storage from Apache Hive.

## Apache CarbonData
→ Compressed columnar storage with indexes from Huawei.

## Apache Iceberg
→ Flexible data format that supports schema evolution from Netflix.

## HDF5
→ Multi-dimensional arrays for scientific workloads.

## Apache Arrow
→ In-memory compressed columnar storage from Pandas/Dremio.

# DATA FILE F

**Apache Parquet**
→ Compressed columnar storage from Cloudera/Twitter

**Apache ORC**
→ Compressed columnar storage from Apache Hive.

**Apache CarbonData**
→ Compressed columnar storage with indexes from Huawei.



An Empirical Evaluation of Columnar Storage Formats

Xinyu Zeng
Tsinghua University
zeng-xy21@mails.tsinghua.edu.cn

Yulong Hui
Tsinghua University
huiyl22@mails.tsinghua.edu.cn

Jiahong Shen
Tsinghua University
shen-jh20@mails.tsinghua.edu.cn

Andrew Pavlo
Carnegie Mellon University
pavlo@cs.cmu.edu

Wes McKinney
Voltron Data
wes@voltrondata.com

Huanchen Zhang*
Tsinghua University
huanchen@tsinghua.edu.cn

# EXECUTION ENGINES

Standalone libraries for executing vectorized query operators on columnar data.
→ Input is a DAG of physical operators.
→ Require external scheduling and orchestration.

Notable implementations:
→ Velox
→ DataFusion
→ Intel OAP

# CONCLUSION

The cloud has made the distributed OLAP DBMS market flourish. Lots of vendors. Lots of money.

But more money, more data, more problems…

# NEXT CLASS

Final Review

15-721 in a single lecture!