#### CS4221 Cloud Databases III. OLAP optimizations

Yao LU 2024 Semester 2

National University of Singapore School of Computing

## Today's agenda

- Current practice in distributed OLAP
- Cost-intelligent data analytics in the cloud

## The trend: disaggregated OLAP

- One recent trend of the last decade is the breakout of OLAP engine sub-systems into standalone open-source components.
  - This is typically done by organizations <u>not</u> in the business of selling DBMS software.

#### • Examples:

- System Catalogs
- Query Optimizers
- File Format / Access Libraries
- Execution Engines

#### Recall: Snowflake architecture





# 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, it can maintain the catalog incrementally.
  - If an external process adds data files, it also needs to update the catalog so that the DBMS is aware of them.
- Notable implementations:
  - <u>HCatalog</u>
  - Google Data Catalog
  - Amazon Glue Data Catalog

## **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:
  - <u>Velox</u>
  - DataFusion
  - Intel OAP

# Query optimizers

- Extendible search engine framework for heuristic- and cost-based query optimization.
  - Applications provide transformation rules and cost estimates.
  - Framework returns either a logical or physical query plan.
- This is the hardest part to build in any DBMS.
- Notable implementations:
  - Greenplum Orca
  - Apache Calcite

## Query optimization for distributed execution

- 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.
  - Using broadcast join vs. repartition join?
  - Considering the impact of data partitioning.

#### **Distributed query execution**

- Executing an OLAP query in a distributed DBMS is roughly the same as on a single-node DBMS.
  - A query plan is represented as a tree of physical operators.
- For each operator, the DBMS considers where input is coming from and where to send output, like parallel execution in a single node.

### Unique challenges

- Data is partitioned across nodes, so a worker thread does not have access to all data for free: accessing data requires network communication.
- Data is partitioned across nodes, so we need to consider leveraging the distributed partitioning to accelerate query executions.
- A query can run for a long time -> how to continue query processing under node failures.

#### Parallel execution for joins







### Parallel execution for joins





# Distributed join algorithms

- 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.
  - This approach loses 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.

- One table is replicated at every node. Each node joins its local data in parallel and then sends its results to a coordinating node.
  - Which tables to replicate?
  - What is the cost?

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



- Tables are partitioned on the join attribute using the same partitioning function.
- Each node performs the join on local data and then sends it to a coordinator node for coalescing.

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



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

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

• Known as "broadcast join".



- Both tables are partitioned on different keys. The DBMS copies/re-partitions the tables on the fly across nodes.
  - This repartitioned data is generally deleted when the query is done.
- SELECT \* FROM R JOIN S ON R.id = S.id



# 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 the original query into partition-specific queries.
- Allows for local optimization at each node.

# Query plan fragments



## 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 snapshot of the intermediate results for a query during execution to allow it to recover if nodes fail.

### Query fault tolerance



## MapReduce

- MapReduce is a data processing paradigm proposed by Google.
- Motivation
  - There is a large volume of raw (unstructured) data to process
    - crawled documents, web request logs, etc.
  - They want to build a system that can use hundreds to thousands of machines.
  - The system needs to be easy to program so that developers do not need to worry about how to parallelize the data processing.
  - Also, the system needs to tolerate node failures when executing a query.



- MapReduce is inspired by the map and reduce primitives of Lisp and many other functional languages.
- Map function:
  - applying the map function to each input record to compute a set of intermediate key/value pairs.
- Reduce function:
  - apply a reduce operation to all the values that share the same key to combine the derived data for the same key together.

### MapReduce

- Developers only need to define the map and reduce functions:
  - map  $(k1,v1) \rightarrow list(k2,v2)$
  - reduce (k2, list(v2)) )  $\rightarrow list(v3)$
- The system will handle
  - Data distribution
  - Parallel execution
  - Fault tolerance

## Example

map(String key, String value):
 // key: document name
 // value: document contents
 for each word w in value:
 EmitIntermediate(w, "1");

```
reduce(String key, Iterator values):
    // key: a word
    // values: a list of counts
    int result = 0;
    for each v in values:
        result += ParseInt(v);
    Emit(AsString(result));
```



- A master process will manage the execution of mappers and reducers
  - Scheduling
  - Managing metadata: where is the output of the mapper
  - Fault tolerance

- Map Phase:
  - Input data is partitioned into M splits.
    - Each split is 64 MB 256 MB.
  - Each map function will process one split and map functions are spread over multiple machines.
  - The output of a map function is automatically partitioned based on the reduce key (k2) and stored on the local disk.

- Reduce Phase:
  - A reducer will copy a partition of data from each mapper.
  - After it gets the partitions from all mappers, it sorts the data by the reduce key k2 to generate a list of (k2, list(v2)).
  - The reduce function is called for each (k2, list(v2)).

### Parallel execution of aggregation

- MapReduce is very similar to the parallel execution of aggregation
  - The difference is that both map/reduce are UDFs, providing the flexibility of processing raw (unstructured) data.



#### Fault tolerance

- The master process pings every worker periodically to detect worker failures.
- Any map tasks executed by a failed worker are re-executed
  - This is because the local disk for the failed worker is not inaccessible.
- Any reduce tasks in progress are re-executed
- Completed reduce tasks do not need to be executed
  - The output of a reducer is on a distributed file system.

# Data locality

- The input data is stored in a distributed file system
  - Each split is replicated (typically 3 copies)
- It is preferable to schedule a mapper to a node that has a copy it needs to process.
- If not, schedule the mapper closer to that copy
  - Within a rack or sharing a switch
- This is to reduce network bandwidth usage.



- With a large number of mappers and reducers, there will be some of them executing very slowly, called stragglers
  - Scheduling problems
  - Slow/bad local disks
- Solution: backup tasks
  - The master process will execute the same task on multiple workers.
  - If one of them finishes, the task is complete.

#### Data file formats

- Most DBMSs use a proprietary on-disk binary file format for their databases.
- 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.
## **Recall:** Parquet

- Apache Parquet is an open-source, column-oriented data file format.
- It provides high-performance compression and encoding schemes to handle complex data in bulk.
- A parquet file consists of one or more row groups.
  - A row group is a partition of rows and includes a column chunk for each column in the dataset.
  - A column chunk is guaranteed to be contiguous in the file and divided up into pages.
    - A page is conceptually an indivisible unit (in terms of compression and encoding).
    - There can be multiple page types that are interleaved in a column chunk.

### Parquet file format

- There are N columns in this table, split into M row groups.
- The file metadata contains the locations of all the column chunk start locations.
- File metadata is written after the data to allow for single-pass writing.
- Readers first read the file metadata to find all the column chunks they are interested in.
- The column chunks should then be read sequentially.

```
4-byte magic number "PAR1"
<Column 1 Chunk 1>
<Column 2 Chunk 1>
. . .
<Column N Chunk 1>
<Column 1 Chunk 2>
<Column 2 Chunk 2>
. . .
<Column N Chunk 2>
. . .
<Column 1 Chunk M>
<Column 2 Chunk M>
. . .
<Column N Chunk M>
File Metadata
4-byte length in bytes of file metadata (little endian)
4-byte magic number "PAR1"
```

#### File metadata



## Today's agenda

- Current practice in distributed OLAP
- Cost-intelligent data analytics in the cloud

**Recall cloud databases** 



Database optimization in the cloud era





**Resource Pool** 



**Resource Pool** 





**Resource Pool** 



#### Cost control is still difficult

Creating as 🛋 ACCOUNTADMIN	
Name	Size ⑦
technicallyWarehouse	X-Large 16 credits/hour
Comment (optional)	X-Small 1 credit/hour
	Small 2 credits/hour
	Medium 4 credits/hour
Advanced Warehouse Options ∧	Large 8 credits/hour
	✓ X-Large 16 credits/hour
Auto Resume	2X-Large 32 credits/hour
Auto Suspend	3X-Large 64 credits/hour
Suspend After (min)	4X-Large 128 credits/hour

→ Users tend to over-provision

#### → Fixed cluster size over the entire workload

### Cost control is still difficult

New W	/arehouse
Creating as 🕻	ACCOUNTADMIN
Name	Size ⑦
technicallyWarehouse	X-Large 16 credits/hour
Comment (optional)	X-Small 1 credit/hour
	Small 2 credits/hour
	Medium 4 credits/hour
Advanced Warehouse Options ∧ Auto Resume	Large 8 credits/hour
	✓ X-Large 16 credits/hour
	2X-Large 32 credits/hour
Auto Suspend	3X-Large 64 credits/hour
Suspand After (min)	4X-Large 128 credits/hour
Suspend Arter (mm)	

- → Users tend to over-provision
- → Fixed cluster size over the entire workload

#### **Resource Waste!**

# Cost control is still difficult

#### **New Warehouse**

Creating as 🛓 ACCOUNTADMIN

Name	Size ⑦
technicallyWarehouse	X-Large 16 credits/hour
Comment (optional)	X-Small 1 credit/hour
	Small 2 credits/hour
	Medium 4 credits/hour
Advanced Warehouse Options ∧ Auto Resume	Large 8 credits/hour
	✓ X-Large 16 credits/hour
	2X-Large 32 credits/hour
Auto Suspend	3X-Large 64 credits/hour
Suspend After (min)	4X-Large 128 credits/hour



Cancel

## Cost intelligence



The system's ability to **self-adapt** to stay at the **Pareto Frontier** in the performance- cost trade-off under different workloads and user constraints.



**Cost Efficiency** 

#### An Idea UI



**Build Indexes** ( / Build Materialized Views Re-partition Data Re-train a Learned Module DBA **\$\$\$** 

An Idea UI

Workload i solution Time: 10s -i 10min Cost: \$2 -i \$0.1



#### Base system architecture



Automatic resource deployment





Automatic resource deployment









## **Cost-oriented database auto-tuning**



**Build Indexes** 



**Build Materialized Views** 

Re-partition Data



Re-train a Learned Module



## Database tuning under fixed resources



- → Speeds up a subset of queries
- → MV update slows down writes











# Towards cost intelligence





## Al-enhanced system components

• QO hint recommendation on workload

Hint set 1

Hint set 2

Hint set 3

Bao

Optimizer

uery

0

- Build a model for query hints
- SIGMOD best paper award

SQL

Parser



## Al-enhanced system components



# Query optimization & cardinality estimation

SELECT \* FROM A,B,C ON A.t = B.t = C.t WHERE A.a=X and B.b=Y and C.c=Z



Cheaper if the first join has fewer <u>output rows</u> (cardinality).

- Fast & accurate cardinality estimation is essential:
  - A wrong plan may result in1000x performance regression.
- Cardinality estimation: a tradeoff between costs & accuracy:
  - Actual execution / multi-dim histograms / samples expensive.
  - Production systems: column-wise stats less accurate.



# **ML-based cardinality estimation**

- Recent ML-enhanced cardinality estimation solutions are instance-optimized.
  - Train a model per workload: MSCN [CIDR 18], Learned Models [VLDB 19, 20].
  - Costly to label & to train models. Dataset/workload changes  $\rightarrow$  retrain.



- Useful cardinality estimation:
  - Small building latency.
  - Adapting to changes quickly.

- $\mu$ s latency per cardinality estimate.
- High accuracy.

## Pretrained models structured tables



Adapting to changes quickly.

• High accuracy.

Yao Lu, Srikanth Kandula, Arnd Christian Konig, Surajit Chaudhuri. Pre-training Summarization Models of Structured Datasets for Cardinality Estimation. VLDB 2022.

#### Frequency distribution & correlation patterns in tables

cardinality = 20

- Images and text have common patterns.
- Accurate cardinality estimation depends on:
  - Multi-dim frequency distribution & correlation patterns (not string semantics).
- Example: [Weather Col<sub>4</sub>, Weather Col<sub>5</sub>]

- Pretrained on 10Ks of relations from UCI ML repo.
  - Similar/better accuracy @ small storage budget.
  - 1-2 orders of magnitude faster to build. Per estimate in  $\mu$ s.
  - Better cardinality estimation = end-to-end performance gains.



Yao Lu, Srikanth Kandula, Arnd Christian Konig, Surajit Chaudhuri. Pre-training Summarization Models of Structured Datasets for Cardinality Estimation. VLDB 2022.
## Today's agenda

- Current practice in distributed OLAP
- Cost-intelligent data analytics in the cloud

## Credits

- Huanchen Zhang, Tsinghua
- Andy Pavlo, CMU