

## COURSE OUTLINE

#### Storage

- $\rightarrow$  Columnar Storage
- $\rightarrow$  Compression
- $\rightarrow$  Indexes

### **Query Execution:**

- $\rightarrow$  Processing Models
- $\rightarrow$  Scheduling
- $\rightarrow$  Vectorization
- $\rightarrow$  Compilation
- $\rightarrow$  Joins
- $\rightarrow$  Materialized Views

## Query Optimization Network Interfaces

## **Client Interface**

Optimization

**Query Execution** 

Storage

## TODAY'S AGENDA

Query Execution Distributed System Architectures OLAP Commoditization

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

For each operator, the DBMS considers where input is coming from and where to send output.  $\rightarrow$  Table Scans

- $\rightarrow$  Joins
- $\rightarrow$  Aggregations
- $\rightarrow$  Sorting











15-721 (Spring 2023)



15-721 (Spring 2023)

# DATA CATEGORIES

#### **Persistent Data:**

- $\rightarrow$  The "source of record" for the database (e.g., tables).
- $\rightarrow$  Modern systems assume that these data files are immutable but can support updates by rewriting them.

#### Intermediate Data:

- $\rightarrow$  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):  $\rightarrow$  **Push Query to Data** 

 $\rightarrow$  Pull Data to Query





# PUSH VS. PULL

#### Approach #1: Push Query to Data

- $\rightarrow$  Send the query (or a portion of it) to the node that contains the data.
- $\rightarrow$  Perform as much filtering and processing as possible where data resides before transmitting over network.

#### Approach #2: Pull Data to Query

- $\rightarrow$  Bring the data to the node that is executing a query that needs it for processing.
- $\rightarrow$  This is necessary when there is no compute resources available where persistent data files are located.

	Filtering and retrieving da	ata using Amazon S3 Select
Approa Ouery Blob Contents Article • 07/20/2021 • 10 minutes to read • 3 contributors The Query Blob Contents API applies a simple Structur contents and returns only the queried subset of the dar the contents of a version or snapshot.	With Amazon S3 Select wat: Microsoft C Feedback red Query Language (SQL) statement on a blob's ta. You can also call Query Blob Contents to query blob Contents to query	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 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. Dbject Content REST API, the AWS Command Line Interface le limits the amount of data returned to 40 MB. To retrieve
The Query Blob Contents request may be constructed myaccount with the name of your storage account: POST Method Request URI https://myaccount.blob.core.windows.net/mycontainer https://myaccount.blob.core.windows.net/mycontainer https://myaccount.blob.core.windows.net/mycontainer	ed as follows. HTTPS is recommended. Replace /myblob?comp=query HTTP/1.0 r/myblob?comp=query&snapshot= <datetime> HTTP/1.1 r/myblob?comp=query&amp;versionid=<datetime></datetime></datetime>	npute resources located.

# 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.



# 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.



# SYSTEM ARCHITECTURE

#### Choice #1: Shared Nothing:

- $\rightarrow$  Harder to scale capacity (data movement).
- $\rightarrow$  Potentially better performance & efficiency.
- $\rightarrow$  Apply filters where the data resides before transferring.

#### Choice #2: Shared Disk:

- $\rightarrow$  Scale compute layer independently from the storage layer.
- $\rightarrow$  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

Traditionally the storage layer in shared-disk DBMSs were dedicated on-prem NAS.  $\rightarrow$  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

# **OBJECT STORES**

Partition the database's tables (persistent data) into large, immutable files stored in an object store.

- $\rightarrow$  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**).

 $\rightarrow$  Some vendors support predicate pushdown (S3).

## **OBJECT STORES**



**ECMU·DB** 15-721 (Spring 2023)

## ADDITIONAL TOPICS

File Formats Table Partitioning Data Ingestion / Updates / Discovery Scheduling / Adaptivity

#### 15

# OBSERVATION

Snowflake is a monolithic system comprised of components built entirely in-house.

Most of the non-academic DBMSs we will cover this semester will have a similar overall architecture.

But this means that multiple organizations are writing the same DBMS software...

# OLAP COMMODITIZATION

One recent trend of the last decade is the breakout OLAP engine sub-systems into standalone opensource components.

 $\rightarrow$  This is typically done by organizations <u>not</u> in the business of selling DBMS software.

#### **Examples:**

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

## OLAP COMMODITIZATION

One recent trend of the last decad OLAP engine sub-systems into sta source components.

 $\rightarrow$  This is typically done by organization of selling DBMS software.

#### **Examples:**

SECMU-DB 15-721 (Spring 2023)

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

#### **Rethinking Database System Architecture:** Towards a Self-tuning RISC-style Database System

Surajit Chaudhuri Microsoft Research Redmond, WA 98052, USA surajite@microsoft.com

Abstract

Database technology is one of the cornerstones for

the new millennium's IT landscape. However, database systems as a unit of code packaging and

deployment are at a crossroad: commercial

systems have been adding features for a long time

and have now reached complexity that makes them a difficult choice, in terms of their "gain/pain

ratio", as a central platform for value-added

information services such as ERP or e-commerce. It is critical that database systems be easy to

manage, predictable in their performance

characteristics, and ultimately self-tuning. For this

elusive goal, RISC-style simplification of server

functionality and interfaces is absolutely crucial.

We suggest a radical architectural departure in which database technology is packaged into much

smaller RISC-style data managers with lean,

specialized APIs, and with built-in self-assessment

Database technology has an extremely successful track

record as a backbone of information technology (IT) throughout the last three decades. High-level declarative query languages like SQL and atomic transactions are key

assets in the cost-effective development and maintenance

of information systems. Furthermore, database technology

continues to play a major role in the trends of our modern

cyberspace society with applications ranging from web-

based applications/services, and digital libraries to

information mining on business as well as scientific data.

Thus, database technology has impressively proven its

benefits and seems to remain crucially relevant in the new

Permission to copy without fee all or part of this material is

granted provided that the copies are not made or distributed for

direct commercial advantage, the VLDB copyright notice and

the title of the publication and its date appear, and notice is

given that copying is by permission of the Very Large Data Base Endowment. To copy otherwise, or to republish, requires a fee

and/or special permission from the Endowment. Proceedings of the 26th International Conference on Very

Large Databases, Cairo, Egypt, 2000

and auto-tuning capabilities

millennium as well.

1. The Need for a New Departure

#### Gerhard Weikum University of the Saarland D-66123 Saarbruecken, Germany weikum@cs.uni-sb.de

Success is a lousy teacher (to paraphrase Bill Gates), and therefore we should not conclude that the database system, as the unit of engineering, deploying, and operating packaged database technology, is in good shape. A closer look at some important application areas and major trends in the software industry strongly indicates that database systems have an overly low "gain/pain ratio". First, with the dramatic drop of hardware and software prices, the expenses due to human administration and tuning staff dominate the cost of ownership for a database system. The complexity and cost of these feed-and-care tasks is likely to prohibit database systems from further playing their traditionally prominent role in the future IT infrastructure. Next, database technology is more likely to be adopted in unbundled and dispersed form within higher-level application services.

Both of the above problems stem from packaging all database technology into a single unit of development, maintenance, deployment, and operation. We argue that this architecture is no longer appropriate for the new age of cyberspace applications. The alternative approach that we envision and advocate in this paper is to provide RISC-style, functionally restricted, specialized data managers that have a narrow interface as well as a smaller footprint and are more amenable to automatic tuning.

The rest of the paper is organized as follows. Section 2 puts together some important observations indicating that database systems in their traditional form are in crisis. Section 3 briefly reviews earlier attempts for a new architectural departure along the lines of the current paper, and discusses why they did not catch on. Section 4 outlines the envisioned architecture with emphasis on RISC-style simplification of data-management components and consequences for the viability of autotuning. Section 5 outlines a possible research agenda towards our vision.

#### 2. Crisis Indicators

To begin our analysis, let us put together a few important observations on how database systems are perceived by customers, vendors, and the research community,

Observation 1: Featurism drives products beyond manageability. Database systems offer more and more features, leading to extremely broad and thus complex

# SYSTEM CATALOGS

A DBMS tracks a database's schema (table, columns) and data files in its catalog.

- $\rightarrow$  If the DBMS is on the data ingestion path, then it can maintain the catalog incrementally.
- $\rightarrow$  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:

- $\rightarrow$  <u>HCatalog</u>
- $\rightarrow$  <u>Google Data Catalog</u>
- $\rightarrow$  Amazon Glue Data Catalog

# QUERY OPTIMIZERS

Extendible search engine framework for heuristicand cost-based query optimization.

- $\rightarrow$  DBMS provides transformation rules and cost estimates.
- $\rightarrow$  Framework returns either a logical or physical query plan.

This is the hardest part to build in any DBMS.

Notable implementations:

 $\rightarrow$  <u>Greenplum Orca</u>

 $\rightarrow$  <u>Apache Calcite</u>



15-721 (Spring 2023)

# 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  $\rightarrow$  Examples: CSV, JSON, XML

There are open-source binary file formats that make it easier to access data across systems and libraries for extracting data from files.

→ Libraries provide an iterator interface to retrieve (batched) columns from files.

# UNIVERSAL FORMATS

#### Apache Parquet (2013)

→ Compressed columnar storage from Cloudera/Twitter

### Apache ORC (2013)

→ Compressed columnar storage from Apache Hive.

## Apache CarbonData (2013)

→ Compressed columnar storage with indexes from Huawei.

## Apache Iceberg (2017)

→ Flexible data format that supports schema evolution from Netflix.

#### **HDF5** (1998)

→ Multi-dimensional arrays for scientific workloads.

## Apache Arrow (2016)

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

## **EXECUTION ENGINES**

Standalone libraries for executing vectorized query operators on columnar data.

- $\rightarrow$  Input is a DAG of physical operators.
- $\rightarrow$  Require external scheduling and orchestration.

Notable implementations:

- $\rightarrow$  <u>Velox</u>
- $\rightarrow$  DataFusion
- $\rightarrow$  <u>Intel OAP</u>

VELOX: META'S UNIFIED EXECUTION ENGINE



# 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 Partitioning Catalogs