Carnegie Mellon University Database Systems Final Review & Systems Potpourri

15-445/645 FALL 2024 >> PROF. ANDY PAVLO

ADMINISTRIVIA

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

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

Final Project Submission Deadline: Monday Dec 16th @ 11:59am



SPRING 2025

Jignesh is recruiting impressionable TAs for 15-445/645 in Spring 2025.

- \rightarrow All BusTub projects will remain in C++.
- \rightarrow If you want to work on fixing BusTub over the winter break for money, please let us know.

Sign up here: https://www.ugrad.cs.cmu.edu/ta/S25



COURSE EVALS

Your feedback is strongly needed:

- → <u>https://cmu.smartevals.com</u>
- → <u>https://www.ugrad.cs.cmu.edu/ta/F24/feedback/</u>

Things that we want feedback on:

- \rightarrow Homework Assignments
- \rightarrow Projects
- \rightarrow Reading Materials
- \rightarrow Lectures



OFFICE HOURS

Andy:

- \rightarrow Wednesday Dec 11th @ 3:30-4:30pm (GHC 9019)
- \rightarrow Thursday Dec 12th @ 3:00-4:00pm (GHC 9019)
- \rightarrow Or email me for an appt

Will:

→ Wednesday Dec 11th @ 10:30-11:30am (GHC 5th Floor Commons)

All other TAs will have their office hours up to and including Saturday Dec 7th



FINAL EXAM

Who: You What: Final Exam Where: Baker Hall A51 When: Friday Dec 13th @ 8:30-11:30am Why: https://youtu.be/8tuoIO4CxOw

Email instructors if you need special accommodations.

https://15445.courses.cs.cmu.edu/fall2024/final-guide.html



FINAL EXAM

Everyone should come to BH A51.

You will then be assigned a random location in either A51 or A53.

There will be TAs stationed in each room to give you the exam and to handle questions. Andy will bounce around the rooms during the exam time.



8

FINAL EXAM

What to bring:

- \rightarrow CMU ID
- \rightarrow Pencil + Eraser (!!!)
- \rightarrow Calculator (cellphone is okay)
- \rightarrow One 8.5x11" page of handwritten notes (double-sided)

STUFF BEFORE MID-TERM

SQL Buffer Pool Management Data Structures (Hash Tables, B+Trees) Storage Models **Query Processing Models** Inter-Query Parallelism **Basic Understanding of BusTub Internals**



QUERY OPTIMIZATION

Heuristics

- \rightarrow Predicate Pushdown
- \rightarrow Projection Pushdown
- \rightarrow Nested Sub-Queries: Rewrite and Decompose

Statistics

- \rightarrow Cardinality Estimation
- \rightarrow Histograms

Cost-based search \rightarrow Bottom-up vs. Top-Down



11

TRANSACTIONS

ACID

- Conflict Serializability:
- \rightarrow How to check for correctness?
- \rightarrow How to check for equivalence?
- View Serializability
- \rightarrow Difference with conflict serializability
- Isolation Levels / Anomalies

TRANSACTIONS

- Two-Phase Locking
- \rightarrow Strong Strict 2PL
- \rightarrow Cascading Aborts Problem
- \rightarrow Deadlock Detection & Prevention

Multiple Granularity Locking

- \rightarrow Intention Locks
- \rightarrow Understanding performance trade-offs
- \rightarrow Lock Escalation (i.e., when is it allowed)

TRANSACTIONS

Optimistic Concurrency Control

- \rightarrow Read Phase
- \rightarrow Validation Phase (Backwards vs. Forwards)
- \rightarrow Write Phase

Multi-Version Concurrency Control

- \rightarrow Version Storage / Ordering
- \rightarrow Garbage Collection
- \rightarrow Index Maintenance



CRASH RECOVERY

Buffer Pool Policies: \rightarrow STEAL vs. NO-STEAL \rightarrow FORCE vs. NO-FORCE

Shadow Paging

Write-Ahead Logging

- \rightarrow How it relates to buffer pool management
- \rightarrow Logging Schemes (Physical vs. Logical)



CRASH RECOVERY

Checkpoints \rightarrow Non-Fuzzy vs. Fuzzy

ARIES Recovery

- \rightarrow Dirty Page Table (DPT)
- \rightarrow Active Transaction Table (ATT)
- \rightarrow Analyze, Redo, Undo phases
- \rightarrow Log Sequence Numbers
- \rightarrow CLRs

DISTRIBUTED DATABASES

System Architectures Replication Partitioning Schemes Two-Phase Commit

TOPICS NOT ON EXAM!

Flash Talks

Seminar Talks

Details of specific database systems (e.g., Postgres) Andy's legal troubles



CMU 15-721 (Spring 2024) SPEED RUN

15721.courses.cs.cmu.edu/spring2024



SEQUENTIAL SCAN: OPTIMIZATIONS

- *Lecture #5* Data Encoding / Compression
- *Lecture #06* Prefetching / Scan Sharing / Buffer Bypass
- Lecture #14 Task Parallelization / Multi-threading
- Lecture #08 Clustering / Sorting
- Lecture #12 Late Materialization
 - Materialized Views / Result Caching
- Lecture #13 Data Skipping
- Lecture #14 Data Parallelization / Vectorization

Code Specialization / Compilation



SELECTION SCANS

SELECT * FROM table WHERE key > \$(low) AND key < \$(high)</pre>

Source: Bogdan Raducanu SCMU-DB 15-445/645 (Fall 2024)

SELECTION SCANS

Scalar (Branching)

```
i = 0
for t in table:
    key = t.key
    if (key>low) && (key<high):
        copy(t, output[i])
        i = i + 1</pre>
```

Scalar (Branchless)

SELECTION SCANS



Source: <u>Bogdan Raducanu</u> **CMU-DB** 15-445/645 (Fall 2024)

Scalar (Branchless)

SELECT * FROM table
WHERE key >= \$low AND key <= \$high</pre>

Vectorized

SELECT * FROM table
WHERE key >= \$low AND key <= \$high</pre>

Vectorized

SELECT * FROM table
WHERE key >= "N" AND key <= "U"</pre>

Vectorized

TID	KEY
100	А
101	Ν
102	D
103	Y
104	Р
105	Ι
106	S
107	
0	

Vectorized

TID	KEY
100	А
101	Ν
102	D
103	Y
104	Р
105	Ι
106	S
107	



SELECT * FROM table
WHERE key >= "N" AND key <= "U"</pre>



Vectorized

TTD	KFY
100	A
101	Ν
102	D
103	Y
104	Р
105	Ι
106	S
107	



SELECT * FROM table
WHERE key >= "N" AND key <= "U"</pre>

ECMU-DB 15-445/645 (Fall 2024)

Vectorized

TID	KEY
100	А
101	Ν
102	D
103	Y
104	Р
105	Ι
106	S
107	



SELECT * FROM table
WHERE key >= "N" AND key <= "U"</pre>



Vectorized





SELECT * FROM table
WHERE key >= "N" AND key <= "U"</pre>

ECMU·DB 15-445/645 (Fall 2024)

Vectorized

TID	KEY
100	А
101	Ν
102	D
103	Y
104	Р
105	Ι
106	S
107	



SELECT * FROM table
WHERE key >= "N" AND key <= "U"</pre>

EFCMU·DB 15-445/645 (Fall 2024)

HIQUE: HOLISTIC CODE GENERATION

For a given query plan, create a C/C++ program that implements that query's execution. \rightarrow Bake in all the predicates and type conversions.

Use an off-shelf compiler to convert the code into a shared object, link it to the DBMS process, and then invoke the exec function.



EFCMU·DB 15-445/645 (Fall 2024)

HIQUE: OPERATOR TEMPLATES

Interpreted Plan

```
for t in range(table.num_tuples):
   tuple = get_tuple(table, t)
   if eval(predicate, tuple, params):
      emit(tuple)
```

- 1. Get schema in catalog for table.
- 2. Calculate offset based on tuple size.
- 3. Return pointer to tuple.



HIQUE: OPERATOR TEMPLATES

Interpreted Plan

```
for t in range(table.num_tuples):
```

```
tuple = get_tuple(table, t)
```

```
if eval(predicate, tuple, params):
    emit(tuple)
```

- 1. Get schema in catalog for table.
- 2. Calculate offset based on tuple size.
- 3. Return pointer to tuple.
- 1. Traverse predicate tree and pull values up.
- 2. If tuple value, calculate the offset of the target attribute.
- 3. Perform casting as needed for comparison operators.
- 4. Return true / false.

HIQUE: OPERATOR TEMPLATES

Interpreted Plan

```
for t in range(table.num_tuples):
   tuple = get_tuple(table, t)
   if eval(predicate, tuple, params):
      emit(tuple)
```

- 1. Get schema in catalog for table.
- 2. Calculate offset based on tuple size.
- 3. Return pointer to tuple.
- 1. Traverse predicate tree and pull values up.
- 2. If tuple value, calculate the offset of the target attribute.
- 3. Perform casting as needed for comparison operators.
- 4. Return true / false.

15-445/645 (Fall 2024

Templated Plan

```
tuple_size = ###
predicate_offset = ###
parameter_value = ###
```

for t in range(table.num_tuples):
 tuple = table.data + t * tuple_size
 val = (tuple+predicate_offset)
 if (val == parameter_value + 1):
 emit(tuple)

VECTORWISE: PRECOMPILED PRIMITIVES

Pre-compiles thousands of "primitives" that perform basic operations on typed data. \rightarrow Using simple kernels for each primitive means that they

are easier to vectorize.

The DBMS then executes a query plan that invokes these primitives at runtime.

- \rightarrow Function calls are amortized over multiple tuples.
- \rightarrow The output of a primitive are the offsets of tuples that


VECTORWISE: PRECOMPILED PRIMITIVES





VECTORWISE: PRECOMPILED PRIMITIVES



```
vec<offset> sel_eq_str(vec<string> col, string val) {
vec<offset> positions;
for (offset i = 0; i < col.size(); i++)
  if (col[i] == val) positions.append(i);
return (positions);</pre>
```



VECTORWISE: PRECOMPILED PRIMITIVES





SYSTEMS

Google BigQuery (2011) Snowflake (2013) Amazon Redshift (2014) Yellowbrick (2014) **Databricks Photon** (2022) **DuckDB** (2019) **TabDB** (2019)









GOOGLE BIGQUERY (2011)

Originally developed as "Dremel" in 2006 as a sideproject for analyzing data artifacts generated from other tools.

- \rightarrow The "interactive" goal means that they want to support ad hoc queries on <u>in-situ</u> data files.
- \rightarrow Did <u>not</u> support joins in the first version.

Rewritten in the late 2010s to shared-disk architecture built on top of GFS.

Released as public commercial product (<u>BigQuery</u>) in 2012.



BIGQUERY: OVERVIEW

- Shared-Disk / Disaggregated Storage Vectorized Query Processing Shuffle-based Distributed Query Execution Columnar Storage
- \rightarrow Zone Maps / Filters
- \rightarrow Dictionary + RLE Compression
- \rightarrow Only Allows "Search" Inverted Indexes

Hash Joins Only

Heuristic Optimizer + Adaptive Optimizations





BIGQUERY: OVERVIEW

Shared-Disk / Disaggregated Storage

Vectorized Query Processing

Shuffle-based Distributed Query Execution

Columnar Storage

- \rightarrow Zone Maps / Filters
- \rightarrow Dictionary + RLE Compression
- \rightarrow Only Allows "Search" Inverted Indexes

Hash Joins Only

Heuristic Optimizer + Adaptive Optimizations





BIGQUERY: IN-MEMORY SHUFFLE

The shuffle phases represent checkpoints in a query's lifecycle where that the coordinator makes sure that all tasks are completed.

Fault Tolerance / Straggler Avoidance:

→ If a worker does not produce a task's results within a deadline, the coordinator speculatively executes a redundant task.

Dynamic Resource Allocation:

 \rightarrow Scale up / down the number of workers for the next stage depending size of a stage's output.





SECMU.DB

BIGQUERY: IN-MEMORY SHUFFLE





32



SPCMU·DB 15-445/645 (Fall 2024)

BIGQUERY: IN-MEMORY SHUFFLE













Distributed File System



Stage n+1



SECMU.DB

BIGQUERY: IN-MEMORY SHUFFLE





Stage n+1



BIGQUERY: IN-MEMORY SHUFFLE









BIGQUERY: IN-MEMORY SHUFFLE





Gooale

Bia Querv

BIGQUERY: DYNAMIC REPARTITIONING

BigQuery dynamically load balances and adjusts intermediate result partitioning to adapt to data skew.



Google

Bia Querv

BIGQUERY: DYNAMIC REPARTITIONING

BigQuery dynamically load balances and adjusts intermediate result partitioning to adapt to data skew.



Google

Bia Querv

BIGQUERY: DYNAMIC REPARTITIONING

BigQuery dynamically load balances and adjusts intermediate result partitioning to adapt to data skew.



BIGQUERY: DYNAMIC REPARTITIONING

BigQuery dynamically load balances and adjusts intermediate result partitioning to adapt to data skew.





Gooale

Bia Querv

BIGQUERY: DYNAMIC REPARTITIONING

BigQuery dynamically load balances and adjusts intermediate result partitioning to adapt to data skew.

DBMS detects whether shuffle partition gets too full and then instructs workers to adjust their partitioning scheme.





Gooale

Bia Querv

BIGQUERY: DYNAMIC REPARTITIONING

BigQuery dynamically load balances and adjusts intermediate result partitioning to adapt to data skew.

DBMS detects whether shuffle partition gets too full and then instructs workers to adjust their partitioning scheme.



Google

Bia Querv

Repartition

BIGQUERY: DYNAMIC REPARTITIONING

BigQuery dynamically load balances and adjusts intermediate result partitioning to adapt to data skew.

DBMS detects whether shuffle partition gets too full and then instructs workers to adjust their partitioning scheme.



Worker

Worker

Gooale

Bia Querv

BIGQUERY: DYNAMIC REPARTITIONING

BigQuery dynamically load balances and adjusts intermediate result partitioning to adapt to data skew.

DBMS detects whether shuffle partition gets too full and then instructs workers to adjust their partitioning scheme.



Google

Bia Querv

BIGQUERY: DYNAMIC REPARTITIONING

BigQuery dynamically load balances and adjusts intermediate result partitioning to adapt to data skew.







SNOWFLAKE (2013)

Managed OLAP DBMS written in C++.

- \rightarrow Shared-disk architecture with aggressive compute-side local caching.
- → Written from scratch. Did not borrow components from existing systems.
- \rightarrow Custom SQL dialect and client-server network protocols.

The OG cloud-native data warehouse.





SNOWFLAKE: OVERVIEW

Cloud-native OLAP DBMS written in C++. Shared-Disk / Disaggregated Storage Push-based Vectorized Query Processing **Precompiled Operator Primitives** Separate Table Data from Meta-Data No Buffer Pool PAX Columnar Storage





SNOWFLAKE: QUERY PROCESSING

Snowflake is a push-based vectorized engine that uses precompiled primitives for operator kernels.

- \rightarrow Pre-compile variants using C++ templates for different vector data types.
- → Only uses codegen (via LLVM) for tuple serialization/deserialization between workers.

Does not support partial query retries \rightarrow If a worker fails, then the entire query has to restart.



SNOWFLAKE: ADAPTIVE OPTIMIZATION

After determining join ordering, Snowflake's optimizer identifies aggregation operators to push down into the plan below joins.

The optimizer adds the downstream aggregations but then the DBMS only enables them at runtime according to statistics observed during execution.



Source: Bowei Chen CMU-DB 15-445/645 (Fall 2024)

SNOWFLAKE: ADAPTIVE OPTIMIZATION

After determining join ordering, Snowflake's optimizer identifies aggregation operators to push down into the plan below joins.

The optimizer adds the downstream aggregations but then the DBMS only enables them at runtime according to statistics observed during execution.



38

Source: Bowei Chen SCMU-DB 15-445/645 (Fall 2024)

SNOWFLAKE: ADAPTIVE OPTIMIZATION

After determining join ordering, Snowflake's optimizer identifies aggregation operators to push down into the plan below joins.

The optimizer adds the downstream aggregations but then the DBMS only enables them at runtime according to statistics observed during execution.



38

Source: Bowei Chen SCMU-DB 15-445/645 (Fall 2024)

SNOWFLAKE: AD

After determining join orderi Snowflake's optimizer identif aggregation operators to pusl into the plan below joins.

The optimizer adds the down aggregations but then the DI enables them at runtime acco statistics observed during ex

Source: Bowei Chen SCMU-DB 15-445/645 (Fall 2024)

Aggregation Placement — An Adaptive Query Optimization for Snowflake



• Medium

Bowei Chen · Follow Published in Snowflake · 8 min read · Aug 10, 2023

Q Search

Snowflake's Data Cloud is backed by a data platform designed from the ground up to leverage cloud computing technology. The platform is delivered as a fully managed service, providing a user-friendly experience to run complex analytical workloads easily and efficiently without the burden of managing on-premise infrastructure. Snowflake's architecture separates the compute layer from the storage layer. Compute workloads on the same dataset can scale independently and run in isolation without interfering with each other, and compute resources could be allocated and scaled on demand within seconds. The cloud-native architecture makes Snowflake a powerful platform for data warehousing, data engineering, data science, and many other types of applications. More about Snowflake architecture can be found in <u>Key Concepts & Architecture documentation</u> and the <u>Snowflake Elastic</u> Data Warehouse research paper.

Write



SNOWFLAKE: FLEXIBLE COMPUTE

If a query plan fragment will process a large amount of data, then the DBMS can temporarily deploy additional worker nodes to accelerate its performance.

Flexible compute worker nodes write results to storage as if it was a table.



Source: Libo Wang CMU-DB 15-445/645 (Fall 2024)



SNOWFLAKE: FLEXIBLE COMPUTE

If a query plan fragment will process a large amount of data, then the DBMS can temporarily deploy additional worker nodes to accelerate its performance.

Flexible compute worker nodes write results to storage as if it was a table.



Source: Libo Wang SCMU-DB 15-445/645 (Fall 2024)



amazon REDSHIFT





AMAZON REDSHIFT (2014)

Amazon's flagship OLAP DBaaS.

- \rightarrow Based on ParAccel's original shared-nothing architecture.
- \rightarrow Switched to support disaggregated storage (S3) in 2017.
- \rightarrow Added <u>serverless</u> deployments in 2022.

Redshift is a more traditional data warehouse compared to BigQuery/Spark where it wants to control all the data.

Overarching design goal is to remove as much administration + configuration choices from users.

AMAZON REDSHIFT RE-INVENTED





REDSHIFT: OVERVIEW

Shared-Disk / Disaggregated Storage

Push-based Vectorized Query Processing

Transpilation Query Codegen (C++)

Precompiled Primitives

Compute-side Caching

PAX Columnar Storage

Sort-Merge + Hash Joins

Hardware Acceleration (AQUA)

Stratified Query Optimizer




Separate nodes to compile query plans using GCC and aggressive caching.

- → DBMS checks whether a compiled version of each templated fragment already exists in customer's local cache.
- \rightarrow If fragment does not exist in the local cache, then it checks a global cache for the **entire** fleet of Redshift customers.

Background workers proactively recompile plans when new version of DBMS is released.





mazon DSHIFT **REDSHIFT: HARDWARE ACCELERATION**

AWS introduced the **AQUA** (Advanced Query Accelerator) for Redshift (Spectrum?) in 2021.

Separate compute/cache nodes that use FPGAs to evaluate predicates.

AQUA was phased out and replaced with Nitro cards on compute nodes



Yellowbrick





YELLOWBRICK (2014)

46

OLAP DBMS written on C++ and derived from a hardfork of PostgreSQL v9.5.

- → Uses PostgreSQL's front-end (networking, parser, catalog) to handle incoming SQL requests.
- \rightarrow They <u>hate</u> the OS as much as I do.

Originally started as an on-prem appliance with FPGA acceleration. Switched to DBaaS in 2021.

Cloud-version uses Kubernetes for all components.



15-445/645 (Fall 2024)



YELLOWBRICK

Shared-Disk / Disaggregated Storage

Push-based Vectorized Query Processing

Transpilation Query Codegen (C++)

Compute-side Caching

Separate Row + PAX Columnar Storage

Sort-Merge + Hash Joins

PostgreSQL Query Optimizer++

Insane Systems Engineering





YELLOWBRICK: ARCHITECTURE



Source: Mark Cusack



YELLOWBRICK: QUERY EXECUTION

Pushed-based vectorized query processing that supports both row- and columnar-oriented data with early materialization.

→ Introduces transpose operators to convert data back and forth between row and columnar formats.

Holistic query compilation via source-to-source transpilation.

Yellowbrick's architecture goal is for workers to always process data residing in the CPU's L3 cache and not memory.

YELLOWBRICK: MEMORY ALLOCATOR

Custom NUMA-aware, latch-free allocator that gets all the memory needed upfront at start-up

- \rightarrow Using **mmap** with **mlock** with <u>huge pages</u>.
- \rightarrow Allocations are grouped by query to avoid fragmentation.
- \rightarrow Claims their allocator is 100x faster than libc **malloc**.

Each worker also has a <u>buffer pool manager</u> that uses MySQL-style approximate LRU-K to store cached data files.



Custom reliable UDP network protocol with kernel-bypass (DPDK)

for internal communication.

- → Each CPU has its own receive/transmit queues that it polls asynchronously.
- → Only sends data to a "partner" CPU at other workers.

YELLOWBRICK: DEVICE DRIVERS

Custom NVMe / NIC drivers that run in user-space to avoid memory copy overheads.

 \rightarrow Falls back to Linux drivers if necessary.









DATABRICKS PHOTON (2022)

Single-threaded C++ execution engine embedded into **Databricks Runtime** (DBR) via **JNI**.

- \rightarrow Overrides existing engine when appropriate.
- → Support both Spark's earlier SQL engine and Spark's DataFrame API.
- \rightarrow Seamlessly handle impedance mismatch between roworiented DBR and column-oriented Photon.

Accelerate execution of query plans over "raw / uncurated" files in a data lake.





DATABRICKS PHOTON (2022)

Photon: A Fast Query Engine for Lakehouse Systems

Alexander Behm, Shoumik Palkar, Utkarsh Agarwal, Timothy Armstrong, David Cashman, Ankur Dave, Todd Greenstein, Shant Hovsepian, Ryan Johnson, Arvind Sai Krishnan, Paul Leventis, Ala Luszczak, Prashanth Menon, Mostafa Mokhtar, Gene Pang, Sameer Paranjpye, Greg Rahn, Bart Samwel, Tom van Bussel, Herman van Hovell, Maryann Xue, Reynold Xin, Matei Zaharia photon-paper-authors@databricks.com Databricks Inc.

ABSTRACT

Many organizations are shifting to a data management paradigm called the "Lakehouse," which implements the functionality of structured data warehouses on top of unstructured data lakes. This from SQL to machine learning. Traditionally, for the most demanding SQL workloads, enterprises have also moved a curated subset of their data into data warehouses to get high performance, governance and concurrency. However, this two-tier architecture is



ECMU·DB 15-445/645 (Fall 2024)



PHOTON: OVERVIEW

Shared-Disk / Disaggregated Storage

Pull-based Vectorized Query Processing

Precompiled Primitives + Expression Fusion Shuffle-based Distributed Query Execution Sort-Merge + Hash Joins Unified Query Optimizer + Adaptive Optimizations





PHOTON: VECTORIZED PROCESSING

Photon is a pull-based vectorized engine that uses precompiled **operator kernels** (primitives).

 \rightarrow Converts physical plan into a list of pointers to functions that perform low-level operations on column batches.

Databricks: It is easier to build/maintain a vectorized engine than a JIT engine.

- → Engineers spend more time creating specialized codepaths to get closer to JIT performance.
- \rightarrow With codegen, engineers write tooling and observability hooks instead of writing the engine.



PHOTON: EXPRESSION FUSION

SELECT * FROM foo
WHERE cdate BETWEEN '2024-01-01' AND '2024-04-01';





PHOTON: EXPRESSION FUSION



```
vec<offset> sel_geq_date(vec<date> batch, date val) {
  vec<offset> positions;
  for (offset i = 0; i < batch.size(); i++)
    if (batch[i] >= val) positions.append(i);
  return (positions);
```

```
vec<offset> sel_leq_date(vec<date> batch, date val) {
  vec<offset> positions;
  for (offset i = 0; i < batch.size(); i++)
     if (batch[i] <= val) positions.append(i);
  return (positions);</pre>
```



PHOTON: EXPRESSION FUSION







Spark (over-)allocates a large number
of shuffle partitions for each stage.
→ Number needs to be large enough to avoid
one partitioning from filling up too much.

After the shuffle completes, the DBMS then combines underutilized partitions using heuristics.





Source: Maryann Xue



Spark (over-)allocates a large number
of shuffle partitions for each stage.
→ Number needs to be large enough to avoid
one partitioning from filling up too much.

After the shuffle completes, the DBMS then combines underutilized partitions using heuristics.





Spark (over-)allocates a large number
of shuffle partitions for each stage.
→ Number needs to be large enough to avoid
one partitioning from filling up too much.

After the shuffle completes, the DBMS then combines underutilized partitions using heuristics.







Spark (over-)allocates a large number of shuffle partitions for each stage.
→ Number needs to be large enough to avoid one partitioning from filling up too much.

After the shuffle completes, the DBMS then combines underutilized partitions using heuristics.





Spark (over-)allocates a large number of shuffle partitions for each stage. → Number needs to be large enough to avoid one partitioning from filling up too much.

After the shuffle completes, the DBMS then combines underutilized partitions using heuristics.









DUCKDB (2019)

Multi-threaded embedded (in-process, serverless) DBMS that executes SQL over disparate data files. \rightarrow PostgreSQL-like dialect with quality-of-life enhancements. \rightarrow "SQLite for Analytics"

Provides zero-copy access to query results via Arrow to client code running in same process.

The core DBMS is nearly all custom C++ code with little to no third-party dependencies. \rightarrow Relies on extensions ecosystem to expand capabilities.





DUCKDB (2019)

@frasergeorgew

Multi-threaded emb DBMS that executes \rightarrow PostgreSQL-like dia \rightarrow "SQLite for Analytic

Provides zero-copy Arrow to client coo

The core DBMS is little to no third-pa \rightarrow Relies on extension My second big finding is the vast majority of queries are tiny, and virtually all queries could fit on a large single node. We maybe don't need MPP systems anymore?



CMU·DB 15-445/645 (Fall 2024)



DUCKDB: OVERVIEW

Shared-Everything

Push-based Vectorized Query Processing

Precompiled Primitives

Multi-Version Concurrency Control

Morsel Parallelism + Scheduling

PAX Columnar Storage

Sort-Merge + Hash Joins

Stratified Query Optimizer





DUCKDB: PUSH-BASED PROCESSING

System originally used pull-based vectorized query processing but found it unwieldly to expand to support more complex parallelism. \rightarrow Cannot invoke multiple pipelines simultaneously.

Switched to a push-based query processing model in 2021. Each operator determines whether it will execute in parallel on its own instead of a centralized executor.





5-445/645 (Fall 2024)

DUCKDB: PUSH-BASED PROCESSING

System originally used pulprocessing but found it un support more complex patholic multiple pip

Switched to a push-based 2021. Each operator dete execute in parallel on its centralized executor.

Switch to Push-Based Execution Model #2393 <> Code 👻 ⊱ Merged) duckdb:master ← Mytherin:pushbasedmodel [] on Oct10,2021 🛇 v0.3.1 ③ • □ Conversation 3 -O- Commits 124 E Checks 0 Files changed 212 +6,097 -3,002 Mytherin commented on Oct 9, 2021 • edited 👻 This PR implements <u>#1583</u> and switches to a push-based execution model. A summary of the • All PhysicalOperators are reworked to use a push-based API. GetChunkInternal is replaced by two separate interfaces, a Source interface and an Operator interface. The Sink interface is mostly kept as-is. See below for more detail. • Pipelines are no longer scheduled as-is. Instead, pipelines are split up into "events" and events • By default DuckDB will default to using all available cores (i.e. PRAGMA threads=X is no longer necessary unless you want to reduce the number of threads DuckDB uses). Several bugs related to parallelism are fixed (primarily relating to recursive CTEs and some edge UNION nodes now support parallelism FULL/RIGHT OUTER join probes now support parallelism Duplicate eliminated joins now support parallelism Whether or not an operator supports parallelism is now determined in the operator itself, rather Several fixes for the query profiler so that the correct number of tuples/timing is now output Pipelines can now be pretty-printed as well (TODO: this should probably be added to the • Simplification for the Arrow scan - since parallel init is always called in the main thread the extra locking/thread-checks are no longer required.



DUCKDB: VECTORS

Custom internal vector layout for intermediate results that is compatible with Velox. Supports multiple vector types:



Source: Mark Raasveldt



DUCKDB: VECTORS

DuckDB uses a unified format to process all vector types without needing to decompress them first. \rightarrow Reduce # of specialized primitives per vector type



Source: Mark Raasveldt SCMU-DB 15-445/645 (Fall 2024)





TabDB is a relational DBMS that stores data in your browser's tab title fields.

It uses <u>Emscripten</u> to convert SQLite's C code into JavaScript.

It then splits the SQLite database file into strings and stores them in your browser tabs.

https://tabdb.io/



CONCLUDING REMARKS

Databases are awesome.

- \rightarrow They cover all facets of computer science.
- \rightarrow We have barely scratched the surface...

Going forth, you should now have a good understanding how these systems work.

This will allow you to make informed decisions throughout your entire career. \rightarrow Avoid premature optimizations.

