# Benchmarking Data and AI Platforms

Shannon Barrow, Lead Solutions Architect, Databricks
Joe Harris, Sr Software Engineer, Databricks

# Benchmarking Data and AI Platforms

Joe Harris, Sr Software Engineer

Shannon Barrow, Lead Solutions Architect

# ~~Benchmarking Data and AI Platforms~~

## How Much You Bench, Bro?

**Shannon Barrow, Lead Solutions Architect**

- Joined Databricks in March 2019

- Previously: Principal - Innovation and Thought Leadership, Accenture Applied Analytics

- Despite overarching benchmark discussion, I may:
  - Put extra focus on TPC-DI
  - Put on my Databricks hat for short segments

# Today's Scope

## Focus on LAKEHOUSE and AI Related Benchmarks

### Lakehouse/OLAP

### ML/Gen AI

- Primarily TPC but others will be mentioned
- Suggestion: view the following benchmarks through the lens of a **full end-to-end Lakehouse architecture**
- How can an organization get a "full picture" of an end-to-end TCO?

- Highlights and challenges
- Focus on Gen AI
- **What** to benchmark?
- Lessons learned from Mosaic

# Why Benchmark?

## Level-setting on the Value and Limitations of Benchmarks

Lies, Damned Lies, and Benchmarks

This is Part 1 of a 5-part series on DW/Analytics Performance Benchmarking and its use in **non-traditional** ways, including DW/Analytics Engineering and Maintenance.

Part 1. Background

Riffing on an old quote attributable (though apparently not originally) to Mark Twain [1] ...

There are three kinds of lies: lies, damned lies, and benchmarks [sic]

We've all seen one-upmanship and "fastest" claims from vendors in the computer technology space using industry-standard benchmarks. From the fastest computers, to the fastest storage systems, to the fastest database systems, there are numerous blogs and papers touting the dominance of one platform or approach over another. In some cases, usually when it suits them, vendors adhere to benchmark specifications and standard reporting requirements. In others, only portions of benchmarks are run and often selectively referenced and co-opted for the desired message. Such is the state of our industry and it will likely not change.

This is the case in the world of Data Warehousing and Analytics. The Transaction Processing Council (TPC) arose out of a need for standard performance metrics and well-defined configuration and pricing guidelines across a number of database vendors and platform offerings. Initially focused on OLTP systems, the organization broadened its work to data warehouse/analytics workloads and other related IT infrastructure areas.

https://rethinkio.com/lies-damned-lies-and-benchmarks/

# Why Benchmark?

## Level-setting on the Value and Limitations of Benchmarks

### Level playing field for all platforms

- Standardization and repeatability
  - To conform to the same practices
  - To conform to common industry operations, use cases, input/output, and scale
  - Industry "agreed upon" testing heuristics
- "Official" submissions

### Can be hard to believe any results

- Potential for:
  - Cheating
  - Bias
  - Abuse
- Slow pace of modernization

# Lakehouse Benchmarking

# TP(P)C - The Ubiquitous Standard

## Most prevalent and well-known

Transaction Processing Performance Council

- Formed in 1988
- Benchmarks across multiple domains
  - Decision Support (OLAP)
    - TPC-DI, TPC-H, TPC-DS
  - Transaction Processing (OLTP)
  - These are the only ones in scope today
  - "Big Data"
    - TPC-HS, TPC-BB
  - Virtualization
    - TPC-V, TPC-HCI
  - Internet of Things (TPC-IOT)
  - AI (TPC-AI)

### Active Benchmarks Per TPC Website:

| Active Benchmarks | |
| --- | --- |
| **Benchmark/Document** | **Current Version** |
| TPC-C | 5.11.0 |
| **TPC-DI** | **1.1.0** |
| **TPC-DS** | **3.2.0** |
| TPC-E | 1.14.0 |
| **TPC-H** | **3.0.1** |
| TPCX-AI | 1.0.3.1 |
| TPCX-BB | 1.6.2 |
| TPCX-HCI | 1.1.9 |
| TPCX-HS | 2.0.3 |
| TPCX-IOT | 2.1.0 |
| TPCX-V | 2.1.9 |

DATA+AI SUMMIT

# Other OLAP Benchmarks

## Are TPC Benchmarks The Only Game in Town?

- SSB (Star Schema Benchmark)
    - Designed to measure the performance of databases in a star schema setup
    - Simpler than TPC benchmarks but focused on specific aspects of OLAP querying

- ClickBench, The No-Join Benchmark
    - Focuses on workloads without joins
    - Simulates scenarios common in clickstream analytics

# Lakehouse Focus - 10k Foot View

## TPC fragments the Lakehouse Architecture into separate benchmarks

- There is no SQL consumption in TPC-DI

- No transformations in TPC-H or TPC-DS

  - Most "unofficial" results even skip the data loading step altogether

DATA AI SUMMIT

# ETL: TPC-DI

# TPC-DI: Data Integration

## The Ingestion and ETL One

- **ZERO** official submissions
  - Was Databricks first to code it?
- I originally presented **completed** benchmark at DAIS 2022
  - Not **submitted** (not for lack of trying)
- Extremely short TL;DR:
  - **ZERO** code given
  - **Ingest**: TXT, CSV, XML
  - **Transform**: based upon 100+ pages of business rules
  - **Load**: all 3 medallion layers

---

### TPC-DI

**TPC-DI is a benchmark for Data Integration**
Historically, the process of synchronizing a decision support s
ETL tools. Recently, ETL was replaced by the more comprehe
a unified data model representation and loading it into a data
and loads it into a data warehouse. The source and destinatio

**Specification**
- The current TPC-DI specification can be found on the **TPC**

**More Information**
- **TPC-DI: the first industry benchmark for Data Integration**

**Results**

    There are no TPC-DI results published yet.

---

# Is the TPC-DI Valuable?

A Frustrating Benchmark that hides some real valuable insight



The best official ETL benchmark available



The worst official ETL benchmark available

DATA'AI SUMMIT

14

# Is the TPC-DI Valuable?

## A Frustrating Benchmark that hides some real valuable insight



The **best** official ETL benchmark available

- Robust even though built for legacy DWs

- Business rules make for realistic test

  - Though it suffers DQ issues with data generator at higher scale factors

- Flexibility in how rules are coded

  - Allows practitioners to optimize to their platform



The **worst** official ETL benchmark available

- Is anyone aware of another "official" ETL benchmark?

- No "official" submittals

- Scoring metrics are confusing and do not even allow for cloud platforms

- No provided code means it is extremely frustrating to attempt this benchmark

  - Made worse by long, confusing business rules

# Excuse Me While I Digress…

I will speak longer on TPC-DI than originally planned

## Why?

A) Because me

2) Because everybody

D) Because Joe Abandoned Us

# 2 Years of TPC-DI on Databricks

## From Initial Implementation to Scorching Performance Today

This is a slide from the **DAIS 2022** Session in which we announced the TPC-DI had been finally implemented

Photon price per billion rows:

$1.51

## Traditional Notebook Workflow Results

Performance Dashboard

- The TPC-DI has a rather confusing benchmark algorithm
- Simplified: TCO approach based on cost per row processed

These were the best performing combinations with On-Demand Pricing:

| Run Time (minutes) | Worker | Total Costs | Price per Billion Rows | Photon | Graviton |
|---|---|---|---|---|---|
| 36.4 | m6gd.8xlarge | $23.28 | $1.44 | No | Yes |
| 24.0 | m6gd.4xlarge | $24.47 | $1.51 | Yes | Yes |

DATA+AI SUMMIT 2022

- **SPOT** instances drops this price to as low as **85 CENTS!**

25

DATA+AI SUMMIT

# 2 Years of TPC-DI on Databricks

## From Initial Implementation to Scorching Performance Today

In April 2023, we published a blog, How We Performed ETL on One Billion Records For Under a Dollar, to tout the power and TCO of Delta Live Tables on this benchmark.

Photon price per billion rows:

$0.96



Delta Live Tables Improves Performance and TCO

Databricks Baseline — 1 min 42 sec — $1.82

Delta Live Tables Optimized Performance — 54 sec — $0.96

Cost and Time to Process 1 Billion Rows in TPC-DI (Complex Enterprise ETL Benchmark)

DATA AI SUMMIT

# 2 Years of TPC-DI on Databricks

## From Initial Implementation to Scorching Performance Today

This video compiled in **September 2023** compares a **dbt** implementation against CDW competitors

Photon price per billion rows:
**$0.73**

DATA+AI SUMMIT

# A Prominent CDW Was missing…

The Truth is Out There…

# A Prominent CDW Was missing…

Some can handle large file sizes, others can't

- We tried benchmarking the other CDW but found it wholly intractable at larger scale factors since it is the only one that is unable to split raw files natively

- We weren't the only ones to notice

# 2 Years of TPC-DI on Databricks

## From Initial Implementation to Scorching Performance Today

- Last month (May 2024) we expanded the benchmark to test non-DWs– which required moving from dbt

- Since AWS has been improving EMR over the last few years this became the obvious first choice for non-dbt tests



| Task name | 3:51:57 PM | 3:57:28 PM |
|---|---|---|
| ingest_BatchDate | 14.2s | |
| ingest_DimDate | 13.7s | |
| Silver_DimBroker | 20s | |
| ingest_DimTime | 14.2s | |
| ingest_FinWire | 1m 20s | |
| ingest_ProspectIncremental | 1m 5s | |
| ingest_StatusType | 13.9s | |
| ingest_TaxRate | 13.9s | |
| Silver_DimCustomer | 28.2s | |
| Silver_DimAccount | 40.5s | |
| Gold_FactCashBalances | 2m 3s | |
| Silver_Prospect | 40.9s | |
| ingest_TradeType | 13.8s | |
| ingest_industry | 14.2s | |
| Silver_DimCompany | 10.6s | |
| Silver_DimSecurity | 26s | |
| Gold_FactWatches | 3m 5s | |
| Silver_DimTrade | 5m 2s | |
| Silver_FactHoldings | 1m 16s | |
| Silver_Financial_CIK | 46.6s | |
| Silver_Financial_CONAME | 49s | |
| Gold_FactMarketHistory | 8m 23s | |

Job ID
701634745375479

Job run ID
1039057252964353

Launched
Manually

Started
06/07/2024, 03:51:57 PM

Ended
06/07/2024, 04:02:40 PM

Duration ⓘ
10m 44s

Queue duration ⓘ
-

Status
✓ Succeeded

Lineage ⓘ
29 upstream tables, 21 downstream tables

Compute
● Shannon Barrow's Cluster

Driver: m7gd.2xlarge · Workers: m7gd.4xlarge · 9 workers · On-demand · DBR: 15.2.x-photon-scala2.12 · auto

View details | Spark UI | Logs | Metrics

Job parameters ⓘ

**9 Graviton 16-core workers = 144 cores**

# 2 Years of TPC-DI on Databricks

## From Initial Implementation to Scorching Performance Today

- **2.2x faster on ¼ the cores** as 2 years ago!
  - 24 minutes down to 10.75 minutes
  - 576 cores down to 144 cores

- Improvements from:
  - PHOTON shifting into *overdrive*
  - Gradual code and orchestration improvements
    - No code is provided - optimize code to match the platform
  - Newer generation VMs
  - Other platform enhancements

- 1 year ago: "A billion rows for under a dollar"

- Today: as low as **20¢** on spot (27¢ on-

### Price per Billion Rows

Bar chart "Price per Billion Rows" with Y-axis from $- to $1.60:
- Jun-22: ~$1.52
- Apr-23: ~$0.96
- Sep-23: ~$0.73
- May-24: ~$0.23

*{ Databricks Digression over }*

DATA AI SUMMIT

# SQL:
# TPC-H & TPC-DS

| | Persona |
|---|---|
| Ingestion | |
| Transformation and Curation | Data Engineer |
| Consumption | SQL / Data Analyst |

DATA AI SUMMIT

# So TPC-H… this still a thing?

The "OG" OLAP benchmark

- Released in April 1999 to "fix" issues with TPC-D

- However, the following year TPC moved to develop a new decision support benchmark to better reflect modern OLAP implementations

- In January 2012, the TPC-DS was released – begging the question why the TPC-H is still used by organizations

# How is TPC-H constructed?

With the "Easy Button"

- Inmon-style DW model with 1 very large table (lineitem) and 7 smaller tables

- All tables contain DATE and STRING columns that are joined using numeric business keys

- Low Query Complexity: 22 queries with only 1 LEFT JOIN, simple aggregates and subqueries, no nested CTEs,  and predicates applied directly to large tables

- Easier Tuning Complexity: often "super-tuned" and each query gets a perfectly covering index

- Does not require a sophisticated optimizer: needs join reordering and predicate pushdown

DATA'AI SUMMIT

# TPC-DS: The Popular Kid in School

## The one we all fight over

- First published in 2012 to counter aging TPC-H's limitations and tackle current OLAP trends

- A benchmark that is often misused/misrepresented by skipping one or more of the 3 components designed to ensure operational considerations aren't forgotten for the sake of over-indexing on this benchmark's SQL queries
  - Load Test, Throughput Test, and Data Maintenance Test
  - A multi-cloud data warehouse platform even publishes a highly-tuned, preloaded TPC-DS dataset in all deployed warehouses for users to consume

- Still valuable in a vacuum when results can be trusted and validated

# How is TPC-DS constructed

Retailer selling goods via 3 different distribution channels: Store, Catalog, Internet

- Based on Kimball dimensional modeling (The Data Warehouse Toolkit)
  - Replaced TPC-H 3NF approach with hybrid approach between 3NF and star schema, or a "multiple snowflake schema"
- Significantly more complicated than TPC-H
  - Heavy on advanced SQL features/functions and lopsided filters
  - 99 queries compared to meager 22 in TPC-H
- 4 query classes:
  - pure reporting queries
  - pure ad-hoc queries
  - iterative OLAP queries
  - extraction or data mining queries

DATA AI SUMMIT

# Feature Comparison: TPC-H & TPC-DS

**Easier to consume cheat sheet**

| Feature | TPC-H | TPC-DS |
|---|---|---|
| Data Model | Simpler schema, uses Inmon style DW model | Complex schema, Kimball style dimensional model. |
| Schema | 1 very large table (lineitem)<br>7 smaller tables | 6 fact tables (3 _sales, 3 _returns)<br>18 dimension tables |
| Data Types | All tables contain DATE and STRING columns.<br>Tables are joined using numeric business keys | Fact tables use only INTEGER and NUMERIC columns.<br>Only dimension tables use TIMESTAMP and STRING<br>Tables are joined using numeric surrogate keys |
| Query Complexity | 22 queries: Low complexity<br>• Only uses 1 LEFT JOIN<br>• Only uses simple aggregates<br>• Subqueries are simple, no nested CTEs<br>• Predicates applied directly to large tables | 99 queries: High complexity<br>• 9 queries use LEFT JOIN, 3 use a cross join<br>• Complex aggregates, 15 queries use window functions<br>• Complex nested CTEs used in most queries<br>• Predicates applied <u>only</u> to dimension tables |
| Tuning Complexity | Easier for vendors to tune<br>• Often "super-tuned" with perfect indexes<br>• Does <u>not</u> require a sophisticated optimizer: | Harder for vendors to tune<br>• Optimizing a specific query can make others slower<br>• <u>Requires</u> a sophisticated query optimizer: must be |

DAT

# Is the TPC-H Valuable?

Best for: Ad -Hoc Manual Benchmarking



**Easy Peasy Man**

- Simpler schema, easier to understand and manage.

- Fewer benchmark queries and they are easy to understand

- Tables contain DATE and STRING columns that are used as predicates



**Too Simple and Easy to Shortcut.  Been replaced!**

- Less realistic, not representative of more complex modern data warehousing needs.

- Simple queries do no reflect hyper-complex real world queries from tools like Tableau and dbt

- Simple schema does not reflect best practices such as SCD type 2.

# Is the TPC-DS Valuable?

## Best for: Vendor Supported POC Evaluations

**Most modern of the Common SQL Benchmarks**

- Complex, realistic schema that better mimics enterprise data warehouses.

- Covers a broad spectrum of query types, SQL operators, and complex joins.

- Requires a sophisticated optimizer, testing more capabilities.

**Complexity & Popularity Result in Missed Stages**

- Higher complexity in setup and longer time to implement and tune.

- Many complex queries can make the results hard to evaluate.

- Can require significant resources to fully utilize and understand performance implications.

# ...TPC-LH?



Persona

Ingestion

Transformation and Curation

Consumption

Bronze          Silver          Gold

**How do we get to here?**
**Benchmark it all**

SQL

Data Engineer

Data Analyst

# The State of "Lakehouse" Benchmarks

## How can the industry do better?

### No Real "Official" Lakehouse Benchmark

- Each benchmark focuses on a portion of end-to-end Lakehouse platform
  - Favors bias and "shortcuts" to improve performance
- Reveals flaws in keeping benchmarks current
  - Example: TPC-DI has no way to calculate its benchmarked metric for cloud platforms

### LHBench: Berkeley white paper implementing a Lakehouse benchmark on EMR

- Composed of 4 tests:
  - TPC-DS
  - TPC-DS Refresh
  - Merge Microbenchmark
  - Large File Count
- Pattern appears sound and it's a great start - but not "official"
  - Community may want to move beyond TPC-DS as the core of the benchmark.

# Thought Experiment: Cluster TPC-DI

## How costly is optimizing and what can be learned to "balance" a Lakehouse benchmark?

- Modified for **OPTIMIZE** on all fact tables

- Adjusting for cluster start times, it takes **44% longer**
  - $6.42 OnDemand nodes ($4.92 on spot)

- Despite tuning the tables this price is still less than:
  - Half the price of EMR
  - 1/3 the price of Big Query
  - Over 15x cheaper than any other CDW we

# Thought Experiment: Cluster TPC-DI

## How costly is optimizing and what can be learned to "balance" a Lakehouse benchmark?

- Modified for OPTIMIZE on all fact tables
- Adjusting for cluster start times, it takes **44% longer**

### Is this worth it?

The answer is always the same: **if** consumption **savings** are **greater than** than the **costs** to optimize the data



DATA·AI SUMMIT

©2024 Databricks Inc. — All rights...

# Thought Experiment: Cluster TPC-DI

## How costly is optimizing and what can be learned to "balance" a Lakehouse benchmark?

- Autostats feature: Stats on Write!

- Leverages Liquid Clustering



STATS

NO STATS

DATA·AI SUMMIT

# Thought Experiment: Cluster TPC-DI

**How costly is optimizing and what can be learned to "balance" a Lakehouse benchmark?**

- Autostats feature: Stats on Write!

- Leverages Liquid Clustering

- Point lookup ad-hoc type query

```
1  set use_cached_result = false;
2  select * from factmarkethistory
3  where sk_securityid = 42949723972 and sk_dateid = 20170401
```

**Raw results** ⌄  +

| ¹²₃ sk_securityid | ¹²₃ sk_companyid | ¹²₃ sk_dateid |
|---|---|---|
| 1 | 42949723972 | 19700316303537 | 20170401 |

**Optimized**

Wall-clock duration ⓘ

| Total wall-clock duration | 350 ms |
|---|---|
| Scheduling ⓘ | 35 ms | 10% |
| Waiting for compute ⓘ | 0 ms |
| Waiting in queue ⓘ | 35 ms |
| Running ⓘ | 315 ms | 90% |
| Optimizing query & pruning files ⓘ | 155 ms |
| Executing ⓘ | 160 ms |

| Start time | 2024-06-09 21:44:17.640 -04:00 |
| End time | 2024-06-09 21:44:17.990 -04:00 |
| Result fetching by client ⓘ | 49 ms |

Aggregated task time ⓘ

| Tasks total time | 39 ms |
| Tasks time in Photon | 82 % |

IO

| Rows returned | 1 |
| Rows read | 1 |
| Bytes read | 94.19 MB |
| Bytes pruned | 170.86 GB |
| Bytes read from cache | 100 % |
| Bytes written | 0 bytes |

**.35s vs 10.2s**

Files & partitions

| Files read | 2 |
| Files pruned | 2,731 |
| Partitions read | 0 |

**30x improvement**

**Not Optimized**

Wall-clock duration ⓘ

| Total wall-clock duration | 10 s 235 ms |
|---|---|
| Scheduling ⓘ | 43 ms | 0% |
| Waiting for compute ⓘ | 0 ms |
| Waiting in queue ⓘ | 43 ms |
| Running ⓘ | 10 s 192 ms | 100% |
| Optimizing query & pruning files ⓘ | 139 ms |
| Executing ⓘ | 10 s 53 ms |

| Start time | 2024-06-09 21:43:35.606 -04:00 |
| End time | 2024-06-09 21:43:45.841 -04:00 |
| Result fetching by client ⓘ | 45 ms |

Aggregated task time ⓘ

| Tasks total time | 3.59 m |
| Tasks time in Photon | 98 % |

IO

| Rows returned | 1 |
| Rows read | 1 |
| Bytes read | 2.57 GB |
| Bytes pruned | 716.18 MB |
| Bytes read from cache | 62 % |
| Bytes written | 0 bytes |

Files & partitions

| Files read | 448 |
| Files pruned | 155 |
| Partitions read | 0 |

# Thought Experiment: Cluster TPC-DI

How costly is optimizing and what can be learned to "balance" a Lakehouse benchmark?

- Autostats feature: Stats on Write!
- Leverages Liquid Clustering Point lookup ad-hoc type query
- BI-like Query using Dimensional filtering and dynamic file pruning



**Optimized**

**Not Optimized**

**20x task time improvement!**

# Thought Experiment: Cluster TPC-DI

## How costly is optimizing and what can be learned to "balance" a Lakehouse benchmark?

- How do we balance the Query Load based on the TPC-DI optimize added latency?

- Back of the napkin math...
  - Conservatively assume 2x performance gains for SQL
  - Assume 5 minutes longer to optimize ETL (at 10k scale factor)– need to save at 5 minutes in queries
  - $5 = \frac{NonOptimizedTablesQueryTime}{2}$
  - Therefore if there is 2x improvement in the SQL times and we need to make up 5 minutes, we need approximately 10 minutes of non-optimized tables query time
    - 5 minutes on optimized tables

DATA AI SUMMIT

# AI Benchmarking

# Why Benchmark in AI/ML?

- Standardized methods allow us to quantitatively know the capabilities of different models, software, and hardware enabling fair comparisons across different solutions.

- Allow ML developers to measure the inference time, memory usage, power consumption, and other metrics that characterize a system.

- Goals and Objectives:
  - Performance assessment
  - Resource evaluation
  - Validation and verification
  - Competitive analysis
  - Credibility
  - Regulation and Standardization

# What to Benchmark in AI/ML?

How does one benchmark something so subjective?

- 3 primary categories:
  - Hardware/System
  - Model
  - Data
- Granularity:
  - Micro
  - Macro
  - End to End
- Training vs Inference



https://harvard-edge.github.io/cs249r_book/contents/benchmarking/benchmarking.html

# In an LLM Not Far Away...

# Keeping Pace and Choosing Wisely

**Benchmarks are rapidly created and deprecated, what can Mosaic's Gauntlet teach us?**

- According to Stanford's 2024 AI INDEX REPORT
  - 15 benchmarks were deprecated in 2023 alone - many of which were less than 4 years old
  - 18 new benchmarks were added in 2023
- The "Mosaic Evaluation Gauntlet" (blog)
  - Evaluated 39 public benchmarks split across 6 core competencies
  - In order to prioritize the metrics that are most useful for research tasks across model scales, we tested the benchmarks using a series of increasingly advanced models



Group 1: Well-behaved metrics robust to few-shot settings

These benchmarks reliably ordered models by training scale and monotonically improved at any number of shots. We believe that these benchmarks can provide a reliable evaluation signal for models in this range.

Figure 2: Monotonically improving benchmarks. These benchmarks include popular tasks like Lambada, BoolQ, Arc, and Hellaswag.

DATA AI SUMMIT

# Challenges and Trends

## Human evaluation is "in"

- Practitioners are growing incredibly skeptical about Academic Benchmarks
    - Habitual issues overfitting models to existing benchmarks
    - MMLU, HumanEval, Hellaswag are bona fide benchmarks but model creators' game the system for models to do well on them

- Accordingly, practitioners today tend to prefer evaluating their LLM options by human preference in the real-world - like LMSYS
    - The HAI Stanford Report even points out "human evaluation is in" (Chapter 2)
    - LMSYS: Allows users to vote on the better response based on a prompt they provide to the LLMs - the user is blind to the choice of the models they're given).

# Q&A