

WRANGLING DATA WITH DUCKDB

Will Angel

2024-10-29

KEY TAKEAWAYS

1. DuckDB is a *very* fast in-process SQL database.
2. Duckplyr is a package for using Dplyr with DuckDB.
3. You can use Duckplyr as a *drop in replacement* for dplyr.
2-30x speedup for dplyr code on large datasets.
4. Big data is smaller than it used to be.

AGENDA

- What is DuckDB
- Why should you care about DuckDB
- When should you use DuckDB
- How can you use DuckDB
- What is Duckplyr
- Data processing performance and profit
- A brief tangent on the shrinking of big

DUCKDB



DUC

WHAT IS DUCKDB?

DuckDB is an open source fast in-process analytical database

- Open Source: Free & Open
- Fast: Performant. **Quickly and efficiently runs analytics**
- In-process: **Runs locally** without a server
- Analytical: DuckDB is optimized for **aggregations** and a support online analytical processing (OLAP). DuckDB supports OLAP but is not as fast for online transaction processing (OLTP)
- Database: DuckDB can be used to efficiently store relational data

WHY SHOULD YOU CARE

- DuckDB is a great tool for local SQL and you can do SQL locally!
- DuckDB is starting to power the next generation of embedded analytical tools, so expect to see data filtering tools to get more power

WHY SHOULD YOU CARE (TECHNICAL)

- DuckDB is like SQLite but for data processing
crunch significant amounts of data locally
without spinning up a full database / data warehouse
which may create significant time/cost overhead
simplify system design
- DuckDB is versatile and fast for data processing
Competitive with spark/polars in benchmarks
- DuckDB is portable with zero dependencies

WHEN SHOULD YOU USE

- You should use DuckDB for SQL data processing
- You have medium-largish data
- You don't want the hassle of procuring resources.

HOW CAN YOU USE DUCK

- Directly use the DuckDB package
- use DBplyr to connect to a local DuckD
- use Duckplyr!

WHAT IS DUCKPLYR

- Duckplyr is a drop in replacement for Dplyr
- Duckplyr uses DuckDB's "relational" API and directly construct logical query plans
- This means you can speed up your Dplyr queries by simply changing the package!
- Unsupported operations will fall back to base R code will always run!
- Duckplyr *overwrites* Dplyr methods, so you need to load the package.

```
library("duckplyr")
```

DUCKPLYR CAVEATS

- Unsupported operations will fall back to Python code will always run!
 - This may mean you don't get a performance boost
 - If you're doing out of memory processing, it can cause issues
- Duckplyr is still under active development, so there are some surprising unsupported operations
 - "message": "No relational implementation for `group_by()`."

BENCHMARKING & SPE



duckplyr

DPLYR VERSUS DUCKPLYR

To compare Dplyr and Duckplyr performance, we used the Global Lake area, Climate, and Population (GLCP). This includes almost 80 million records of temperature data for lakes, for a 3.9GB CSV file.

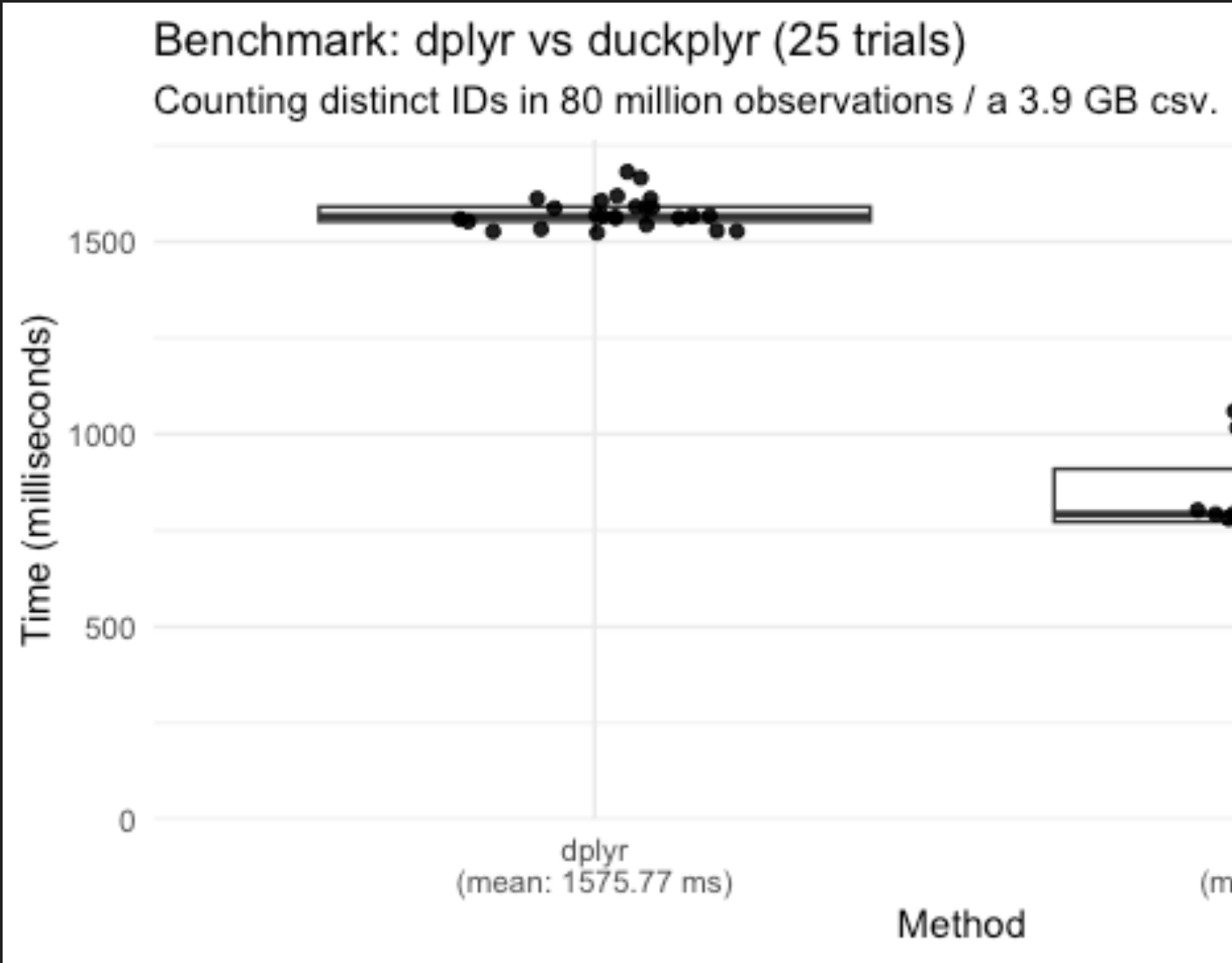
Dplyr:

```
data |>
  distinct(HYBAS_ID) |>
  nrow()
```

Duckplyr:

```
data |>
  duckplyr::as_duckplyr_tibble() |>
  distinct(HYBAS_ID) |>
  nrow()
```

DPLYR VERSUS DUCKPLYR



86% speed increase counting distinct IDs in 80 million row

DPLYR VERSUS DUCKPLYR

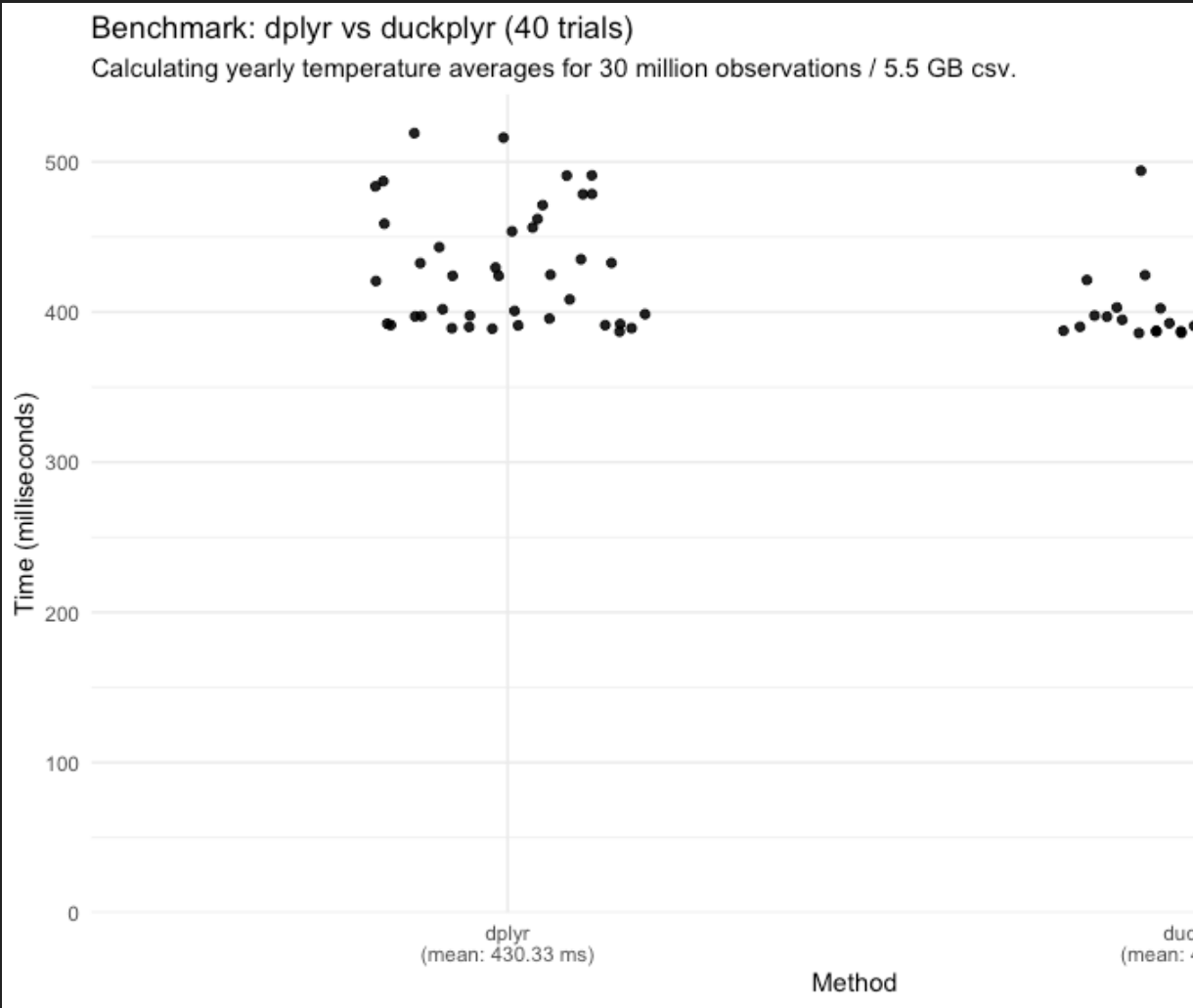
Now we'll take a look at the overall GLCP 10 million observations in a 5.5GB CSV file.

```
data |>
  summarise(
    avg_size = mean(total_km2), n = n(), .by = cou
```

```
library(duckplyr)
```

```
data |>
  summarise(
    avg_size = mean(total_km2), n = n(), .by = cou
```

DPLYR VERSUS DUCKPLYR



We ran into a limitation of Duckplyr and fell back to dplyr.
{"version": "0.4.1", "message": "No relational implementation"}

DPLYR VERSUS ~~DUCKDB~~ DUCKDB

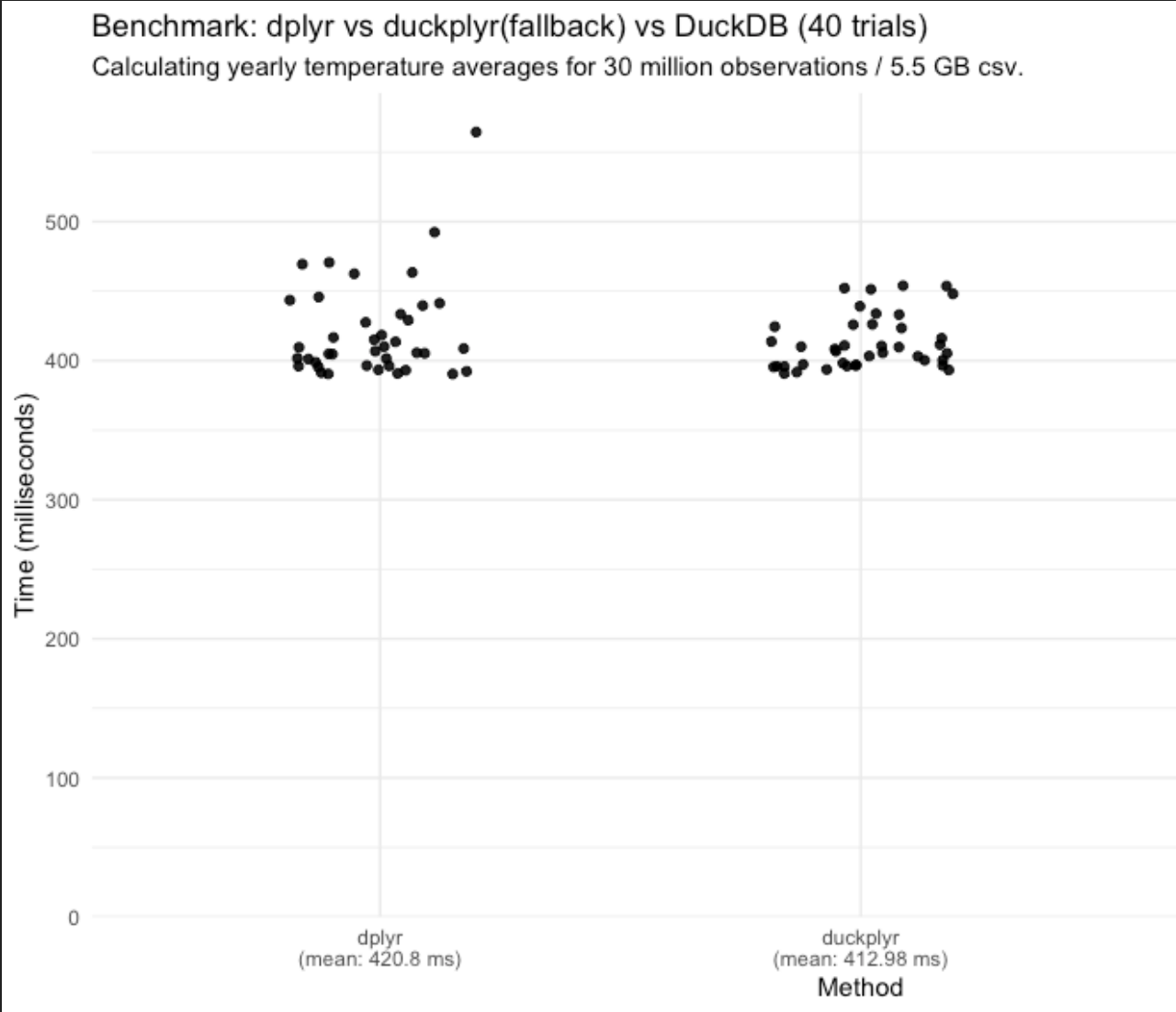
Let's switch to DuckDB and write some S

```
library(duckdb)

con <- dbConnect(duckdb())
duckdb_register(con, "data", data)

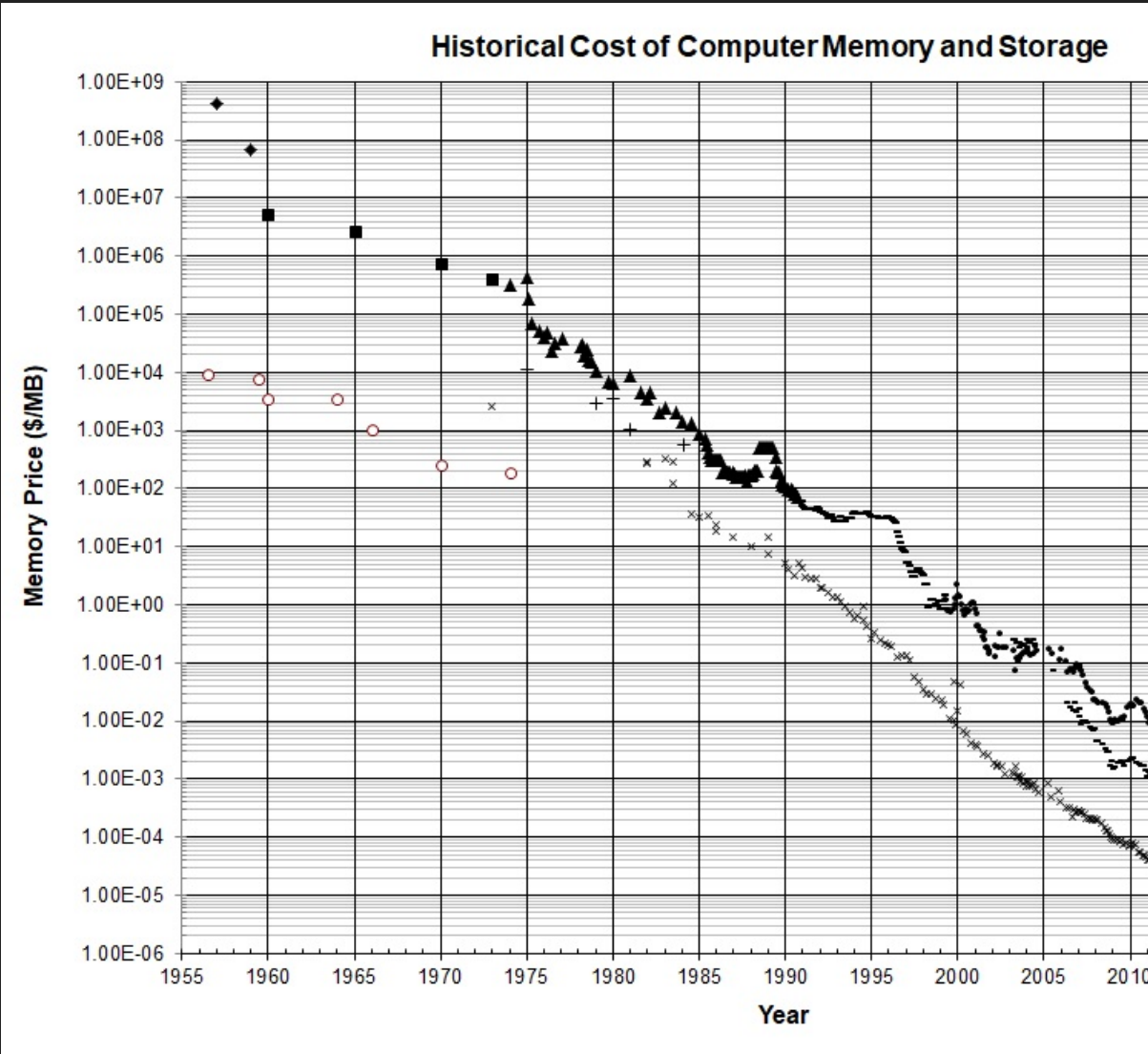
dbGetQuery(con, "
  SELECT
    country,
    AVG(total_km2) as avg_size,
    COUNT(*) as n
  FROM data
  GROUP BY country
")
```

DPLYR VERSUS DUCKPI DUCKDB

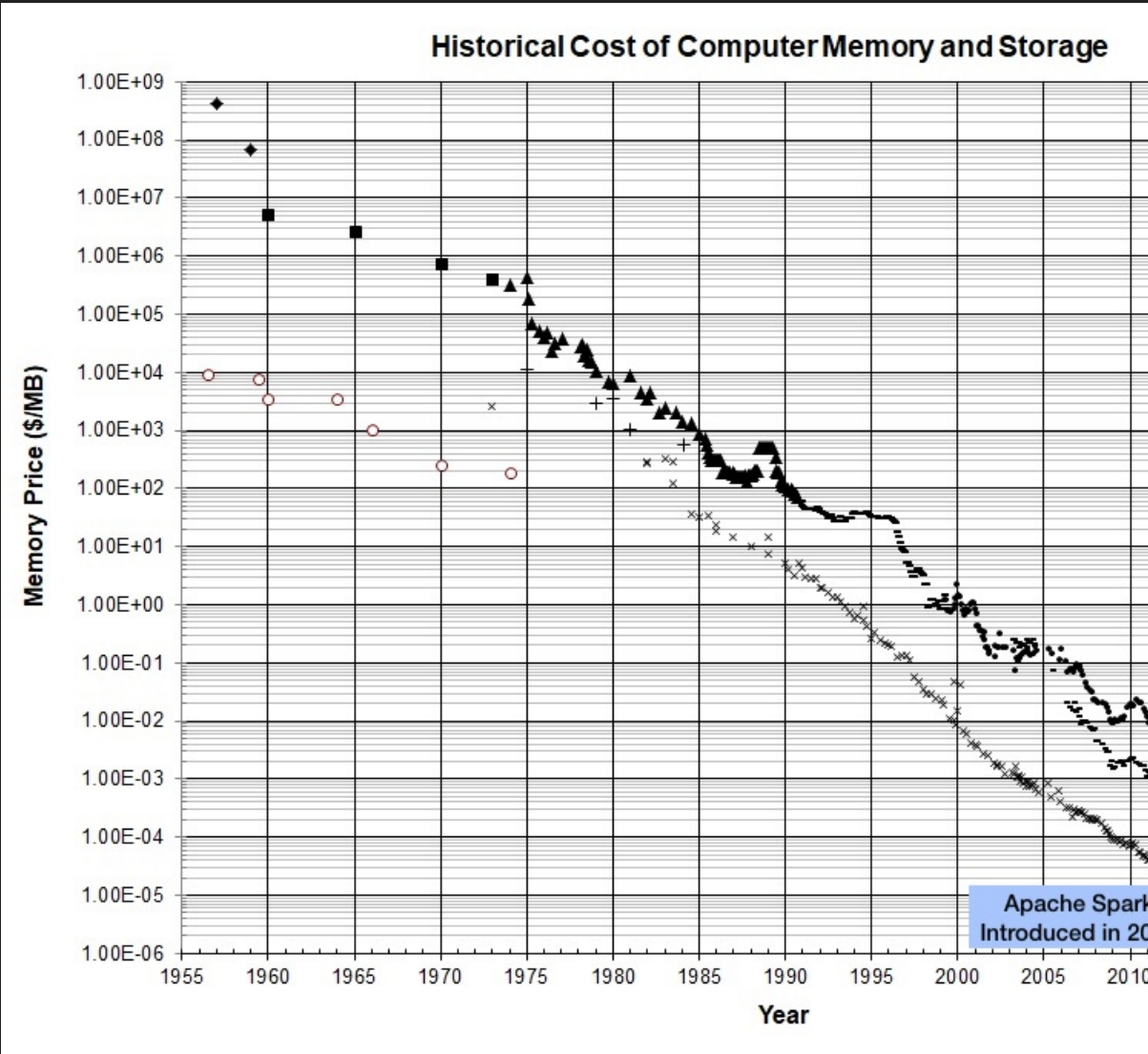


**BIG DATA
SMALLER
THAN IT U
TO BE**

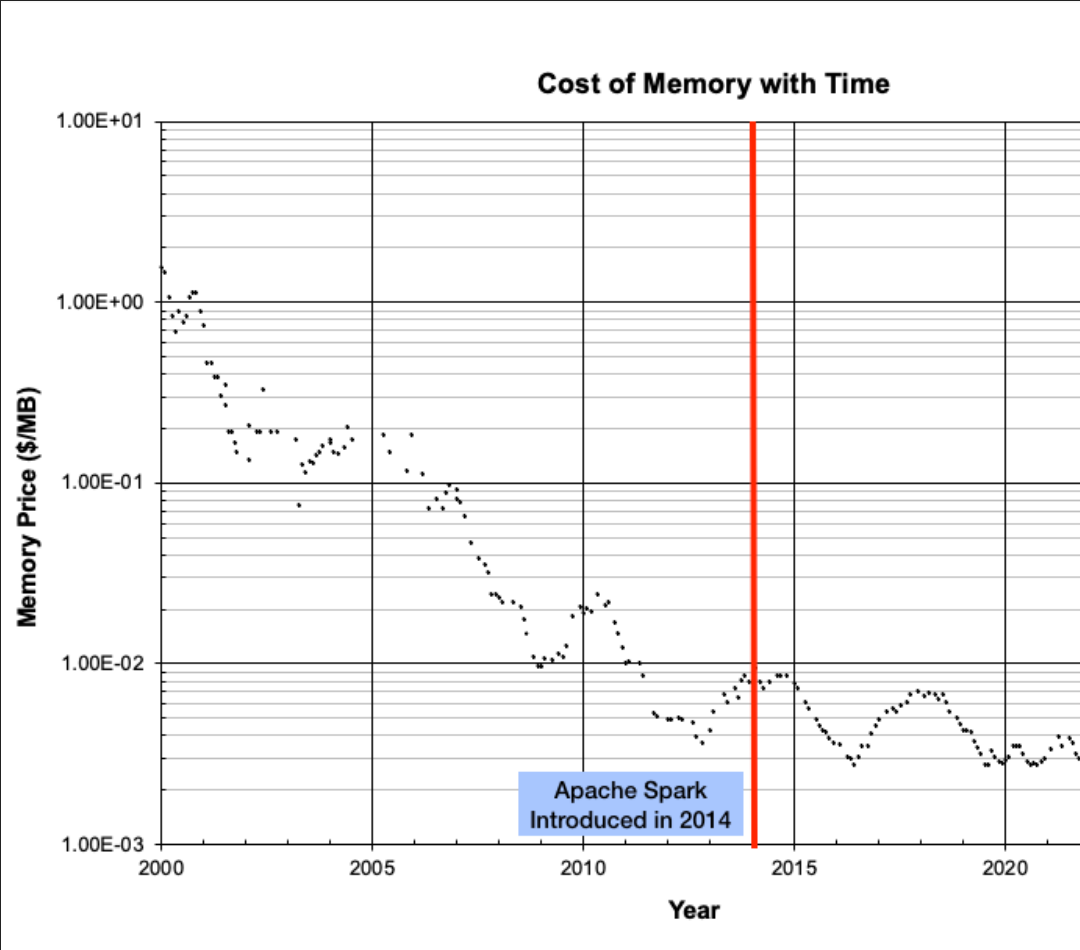
COMPUTERS ARE GETTING BIGGER



COMPUTERS ARE GETTING BIGGER

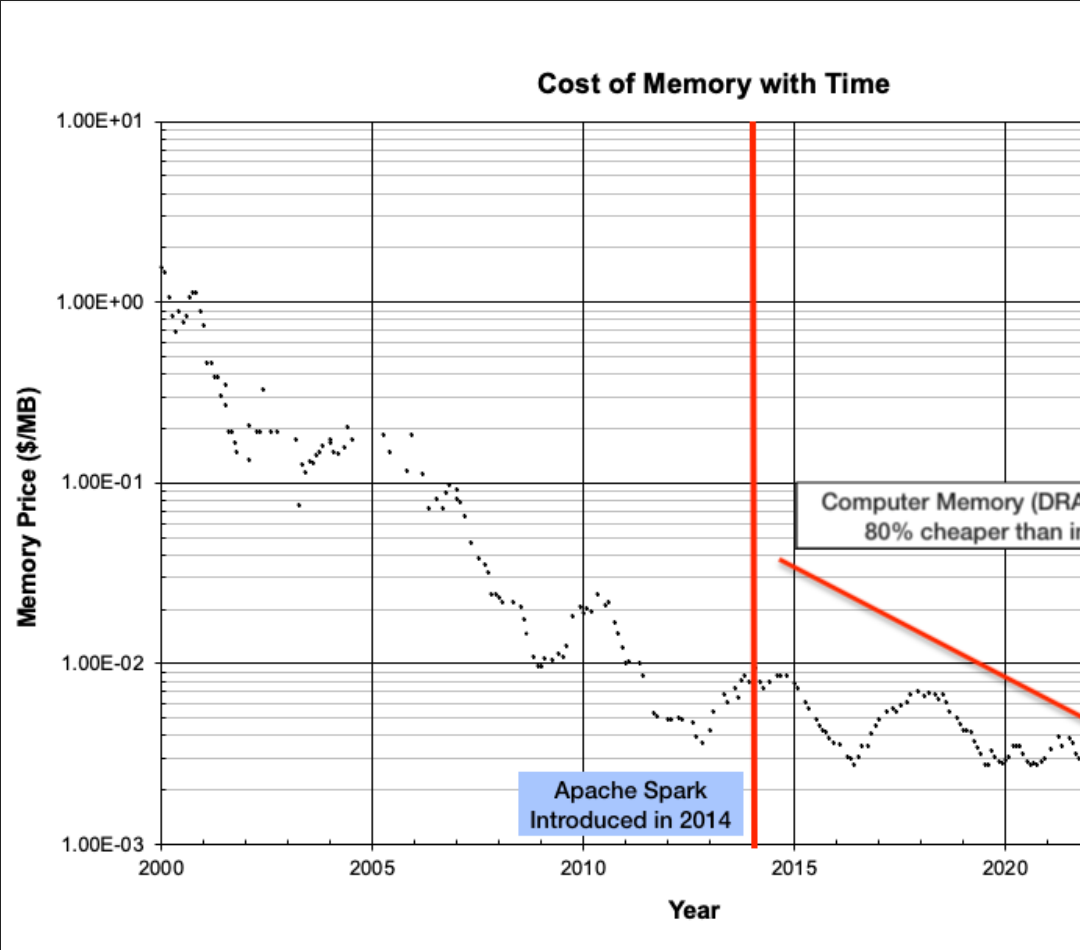


MEMORY HAS GOTTEN SIGNIFICANTLY CHEAPER



<https://jcmit.net/memoryprice.htm>

MEMORY HAS GOTTEN SIGNIFICANTLY CHEAPER



<https://jcmit.net/memoryprice.htm>

IN CONCLUSION:

BIG DATA IS SMALLER THAN IT USED TO BE

With more access to cheaper computing tools (DuckDB/Polars(TidyPolars)/DataFusion) to process data more efficiently:

1. We'll end up with more "Annoyingly Massive" files that are too big to open in Excel, fits in RAM on a desktop.
2. We can defer investing in large scale distributed systems like Spark more often, and wait longer before writing distributed code.

THANK YOU!

- [/in/william-angel/](#)
- [@datadrivenangel](#)
- www.williamangel.net
- yt@williamangel.net



All ducks to