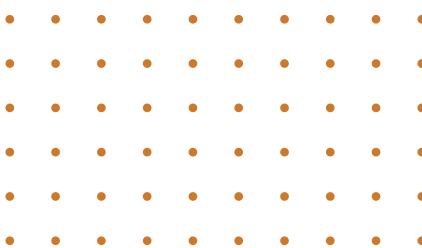
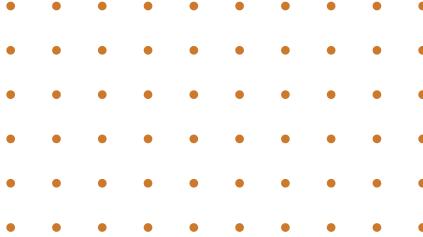
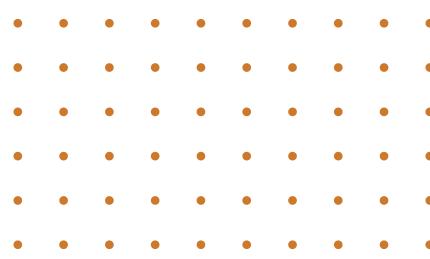
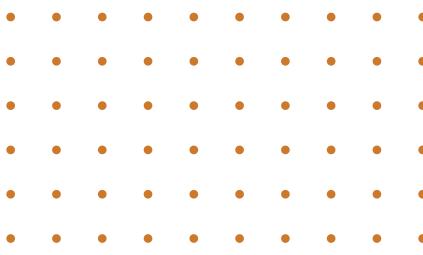


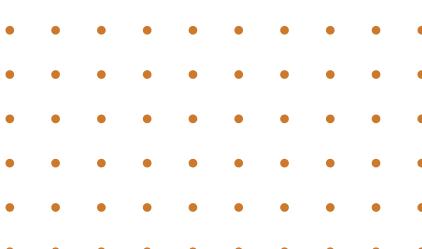
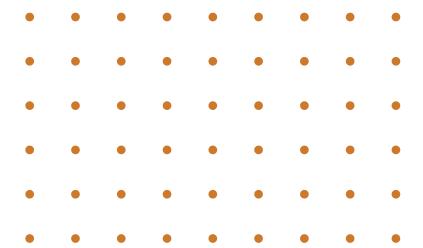
May 18, 2024

# Intelligent Tomato group Data Analysis



Name any  
**dataframe library**  
you have used?





**The largest number of rows**  
you have ever processed  
using Pandas library in Python?

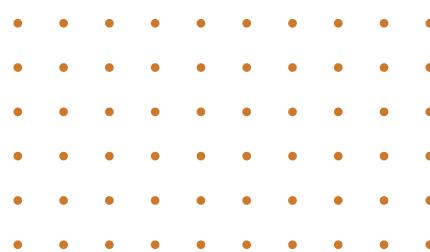
# How about... one billion rows?

## 1 Billion Row Challenge

Calculate the min, max, and average of 1 billion measurements

[Accept the challenge](#)

[Original blog post](#)



# Normal Pandas read\_csv

```
1 import pandas as pd
2
3 # Read the file into a DataFrame
4 df = pd.read_csv("measurements.txt", names=['station', 'measure'])
5
6 # Group by station and calculate min, max, and mean
7 result = df.groupby('station')['measure'].agg(['min', 'max', 'mean'])
8
9 print(result)
10
```

⊗

The Kernel crashed while executing code in the current cell or a previous cell.  
Please review the code in the cell(s) to identify a possible cause of the failure.  
Click [here](#) for more info.  
View Jupyter [log](#) for further details.

The kernel crashes because  
the file size exceeds the available RAM!

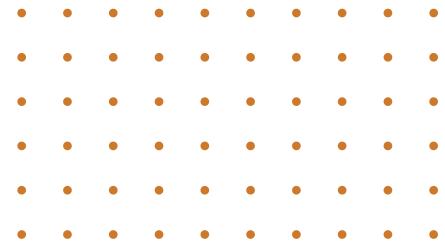
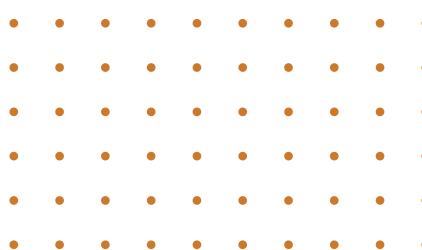
# Pandas read\_csv in chunks

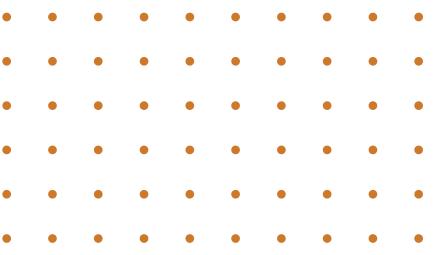
```
1 def process_chunk(chunk):
2     # Aggregates the data within the chunk using Pandas
3     aggregated = chunk.groupby('station')['measure'].agg(['min', 'max', 'mean']).reset_index()
4     return aggregated
5
6 def create_df_with_pandas(filename, total_linhas, chunksize=chunksize):
7     total_chunks = total_linhas // chunksize + (1 if total_linhas % chunksize else 0)
8     results = []
9
10    with pd.read_csv(filename, sep=';', header=None,
11                      names=['station', 'measure'],
12                      chunksize=chunksize) as reader:
13        # Wrapping the iterator with tqdm to visualize progress
14        with Pool(CONCURRENCY) as pool:
15            for chunk in tqdm(reader, total=total_chunks, desc="Processing"):
16                # Processes each chunk in parallel
17                result = pool.apply_async(process_chunk, (chunk,))
18                results.append(result)
19
20            results = [result.get() for result in results]
21
22    final_df = pd.concat(results, ignore_index=True)
23
24    final_aggregated_df = final_df.groupby('station').agg({
25        'min': 'min',
26        'max': 'max',
27        'mean': 'mean'
28    }).reset_index().sort_values('station')
29
30    return final_aggregated_df
```

Pandas allows us to read a large DataFrame in chunks and process them in parallel.

# GUESS

Does Pandas chunking helps?





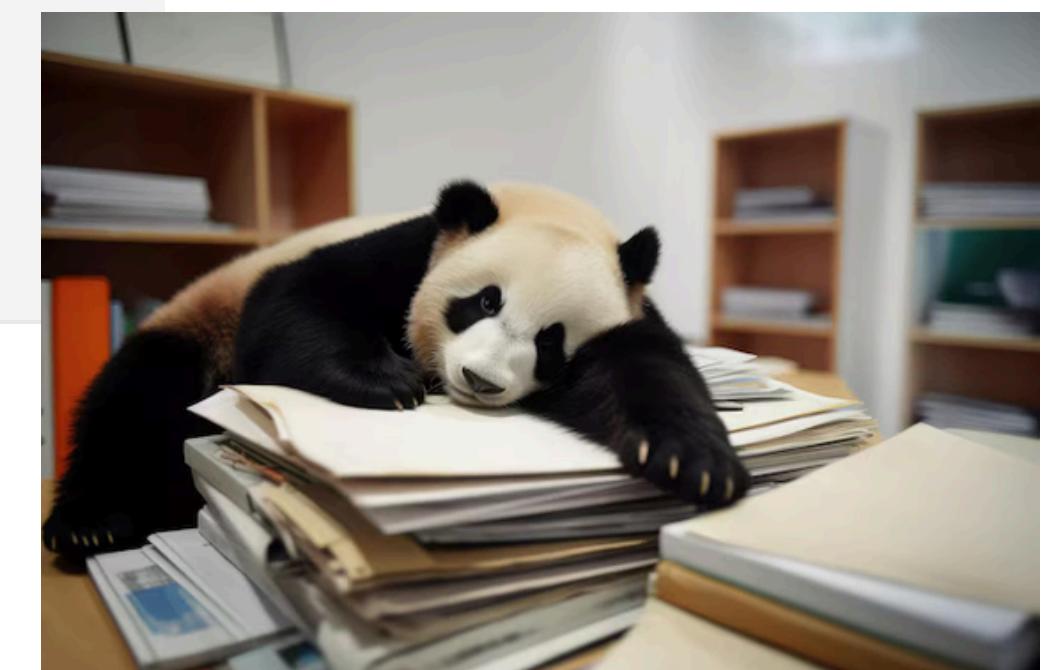
# 6 hours pass and our Pandas still keep working so hard...

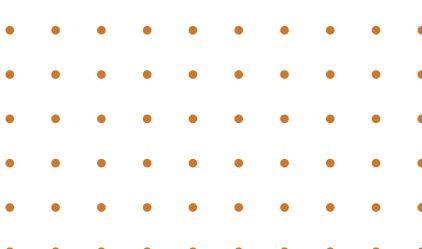
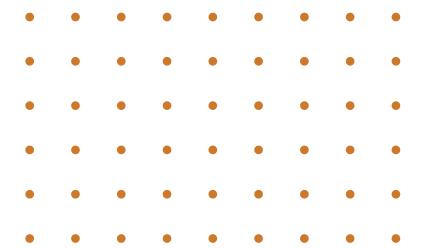
```
start_time = time.time()
df = create_df_with_pandas(filename, total_linhas, chunksize)
took = time.time() - start_time
```

```
print(df.head())
print(f"Pandas took: {took:.2f} sec")
```

⌚ 367m 58.6s

Processando: 100% |██████████| 10/10 [04:40<00:00, 28.02s/it]





**Don't be frustrated...**  
**Let's find a solution!**

# Google

⌚ fast dataframe library python





fast dataframe library python



All

Videos

Images

Shopping

News

More

Tools

Example

Pdf

Github

About 911,000 results (0.21 seconds)

**Polars** is a blazingly fast DataFrame library for manipulating structured data. The core is written in Rust, and available for Python, R and NodeJS.



Polars

<https://docs.pola.rs>



[Polars user guide: Index](#)

[About featured snippets](#) • [Feedback](#)

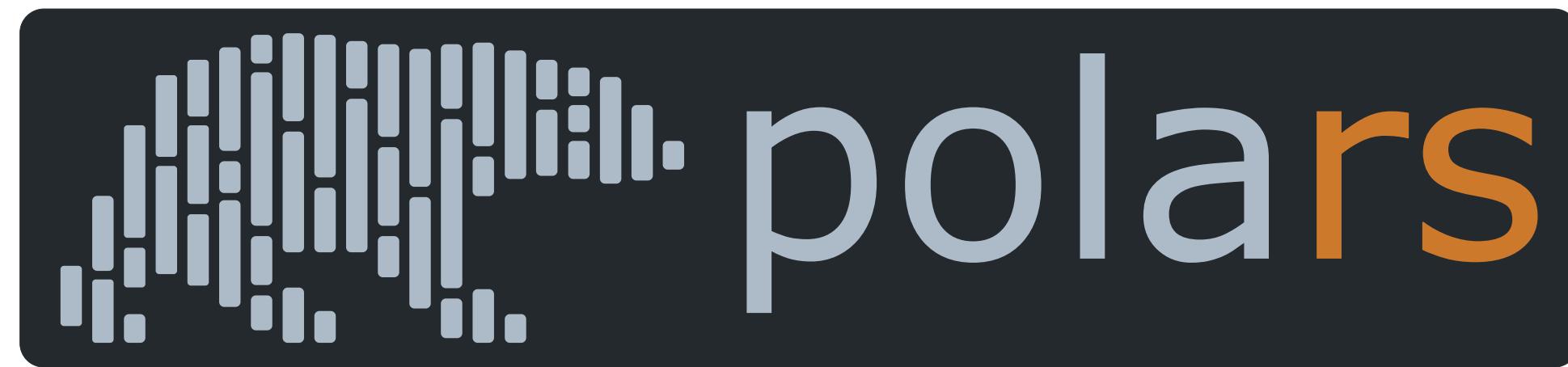


• • • • •

# Demonstrating...

**~30 seconds  
to find min, max, mean of  
1 billion rows**

INTELLIGENT DATA ANALYSIS



A **Lightning-Fast**  
DataFrame Library

# TOMATO TEAM MEMBERS



21120275  
Huỳnh Cao Khôi



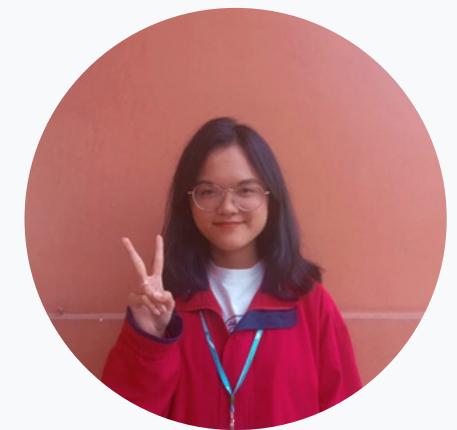
21120290  
Hoàng Trung Nam



21120302  
Huỳnh Trí Nhân



21120308  
Phạm Lê Tú Nhi



21120533  
Lê Thị Minh Phương



21120593  
Võ Hoàng Hoa Viên

# CONTENTS

## POLARS: A LIGHTNING-FAST DATAFRAME LIBRARY

01

STORY

02

KEY DESIGNS

05 CONCLUSION



03

COMPARISON

04

TUTORIAL



# 01 STORY



2020 - Coronavirus Pandemic

Vink's hobby project in

**RUST**

to learn:

- Query engines
- Apache arrow
- Rust

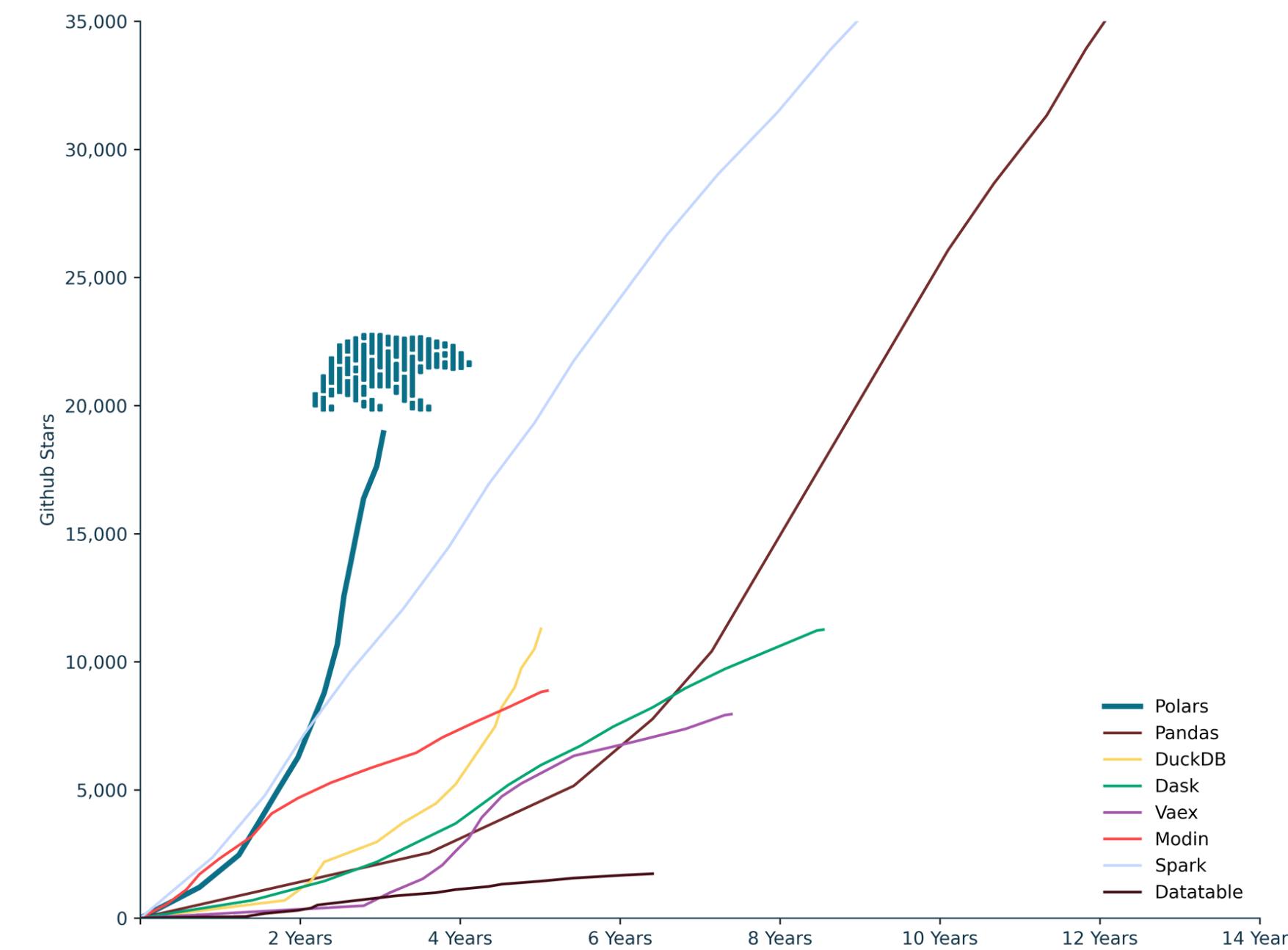
First commit: 23 June 2020

Richie Vink

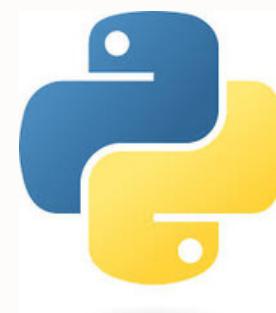
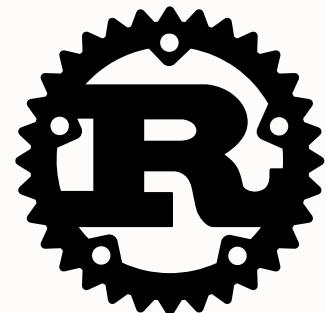


However...  
Polars grew  
beyond his expectation!!!

Benchmark  
Graph



# SUPPORTED LANGUAGES



officially



on the way



## 02 KEY DESIGNS

To understand Polars and its key design decisions,  
we start with the familiar dataframe library



Pandas was build in April, 2008  
by **Wes McKinney**



In 2017, Wes McKinney published a post named “Apache Arrow and the 10 Things I Hate About Pandas”.



“

*I didn't know much about software engineering or even how to use Python's scientific computing stack well back then. My code was ugly and slow.*

”

Wes McKinney,  
Apache Arrow and the 10 things i hate about pandas

# The 10 things I hate about Pandas

by Wes McKinney

|  |  |
|--|--|
| <b>Internals too far from “the metal”</b>                  | <b>Complex groupby operations awkward and slow</b>           |
| <b>No support for memory-mapped datasets</b>               | <b>Appending data to a DataFrame tedious and very costly</b> |
| <b>Poor performance in database and file ingest/export</b> | <b>Limited, non-extensible type metadata</b>                 |
| <b>Warty missing data support</b>                          | <b>Eager evaluation model, no query planning</b>             |
| <b>Weak support for categorical data</b>                   | <b>“Slow”, limited multicore algorithms</b>                  |

[Source](#)

“

*I strongly feel that **Arrow** is a key technology  
for the next generation of data science tools.*

*[...] building a faster, cleaner core pandas  
implementation, which we may call **Pandas 2**.*

”

Wes McKinney,  
Apache Arrow and the 10 things i hate about pandas

And **Apache Arrow**,  
is actually the key technology  
for both **Pandas 2** and **Polars**

How can  optimize query engines?

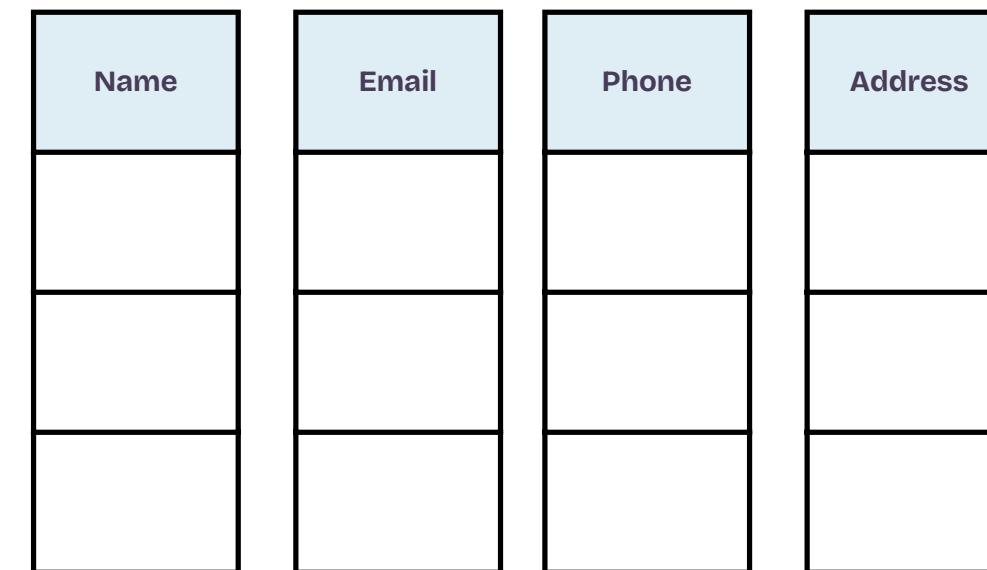
# Apache Arrow

## A language-independent **columnar** memory format

**Row storage**

| Name | Email | Phone | Address |
|------|-------|-------|---------|
|      |       |       |         |
|      |       |       |         |
|      |       |       |         |

**Columnar storage**



Data is stored row by row

Data is organized by columns

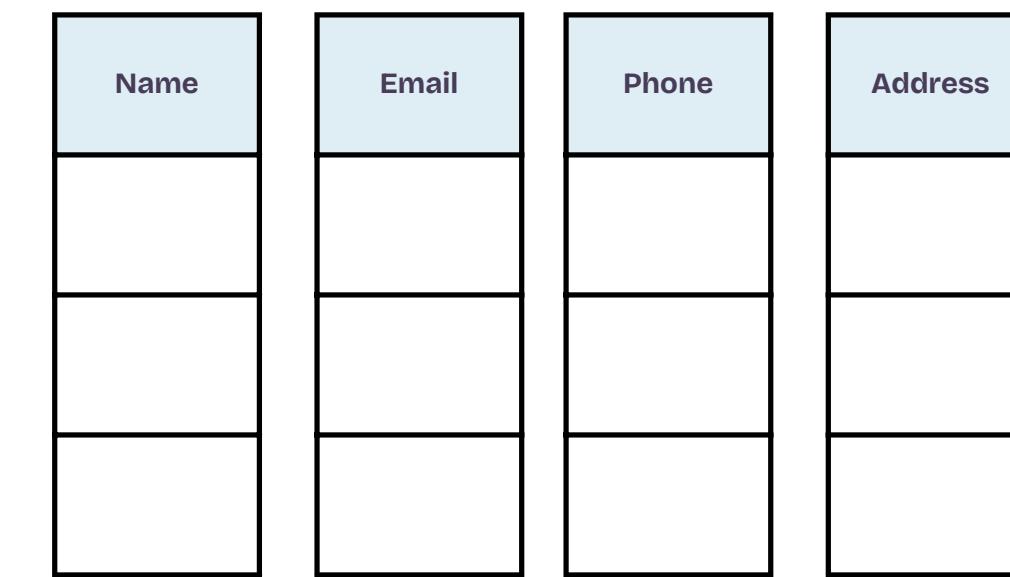
# Apache Arrow

## A language-independent **columnar** memory format

**Row storage**

| Name | Email | Phone | Address |
|------|-------|-------|---------|
|      |       |       |         |
|      |       |       |         |
|      |       |       |         |

**Columnar storage**



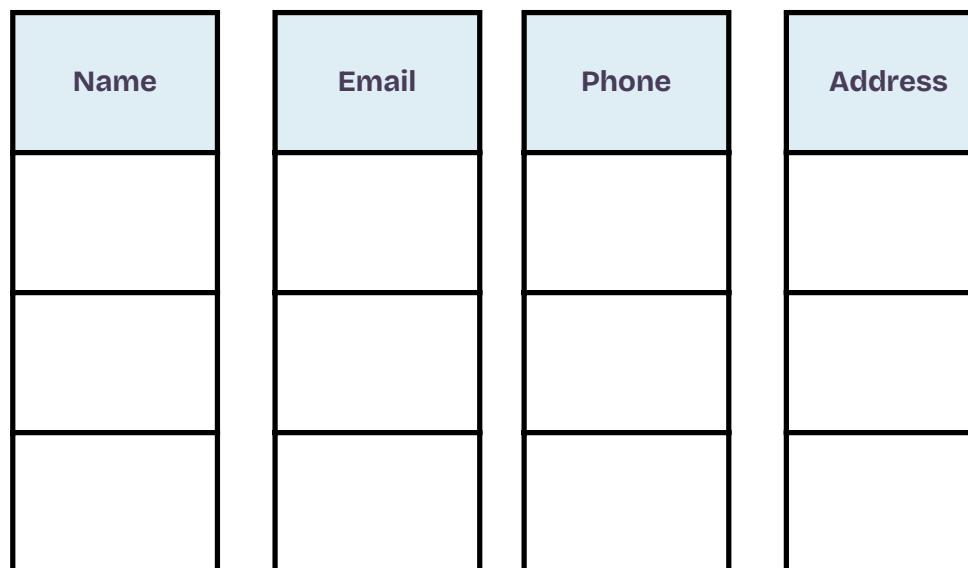
Efficient for **transactional operations**  
(selecting, inserting, updating, or deleting rows)

Efficient for **analytical queries** that scan  
large datasets but only access a few columns

# Apache Arrow

## A language-independent **columnar memory format**

### Columnar storage

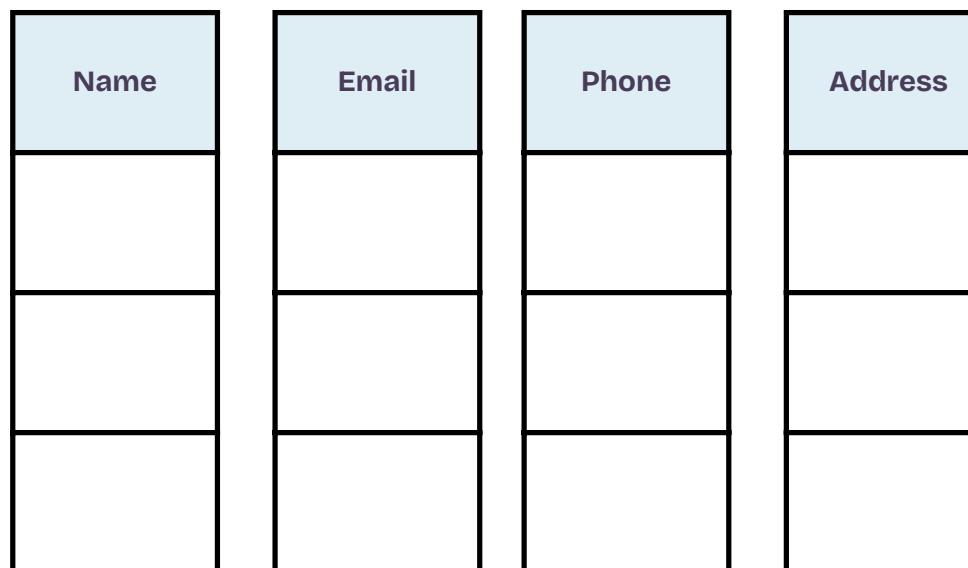


Why is **columnar format**  
such powerful?

# Apache Arrow

## A language-independent **columnar memory format**

### Columnar storage



Why is **columnar format**  
such powerful?

**DATA PARALLELISM**

# Data parallelism

Single Instruction Single Data

$$\begin{array}{rcl} a & + & a \\ \boxed{a} & & \boxed{a} \\ \hline & = & \boxed{2a} \end{array}$$
$$\begin{array}{rcl} b & + & b \\ \boxed{b} & & \boxed{b} \\ \hline & = & \boxed{2b} \end{array}$$
$$\begin{array}{rcl} c & + & c \\ \boxed{c} & & \boxed{c} \\ \hline & = & \boxed{2c} \end{array}$$
$$\begin{array}{rcl} d & + & d \\ \boxed{d} & & \boxed{d} \\ \hline & = & \boxed{2d} \end{array}$$

Single Instruction Multiple Data

$$\begin{array}{rcl} a & & a \\ \boxed{a} & & \boxed{a} \\ b & & b \\ \boxed{b} & & \boxed{b} \\ c & & c \\ \boxed{c} & & \boxed{c} \\ d & & d \\ \boxed{d} & & \boxed{d} \\ \hline & + & = \\ & & \end{array} \quad \begin{array}{l} 2a \\ 2b \\ 2c \\ 2d \end{array}$$

Perform a **Single Instruction** across **Multiple Data**  
at the same time (**SIMD**)

## Pandas 2 and Polars leverage Arrow Columnar Format to solve some Pandas problems

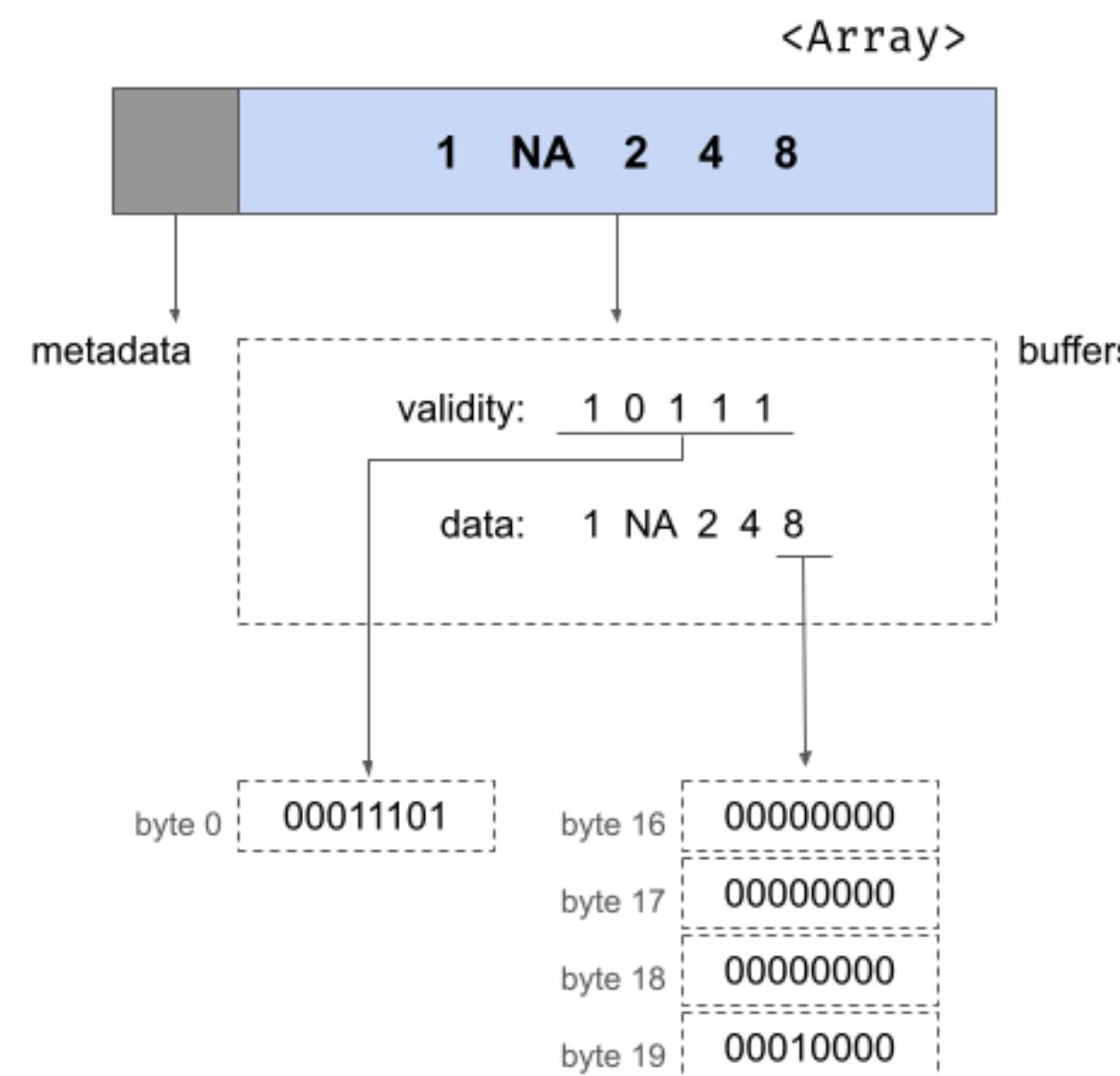
|   |  |
|---|--|
| Getting closer to the metal                         | Better groupby operations                    |
| No support for memory-mapped datasets               | DataFrame appending is more memory-efficient |
| Poor performance in database and file ingest/export | Limited, non-extensible type metadata        |
| Warty missing data support                          | Eager evaluation model, no query planning    |
| Weak support for categorical data                   | "Slow", limited multicore algorithms         |

Apache Arrow uses **Validity bitmap buffer**



**Validity bitmap buffer**

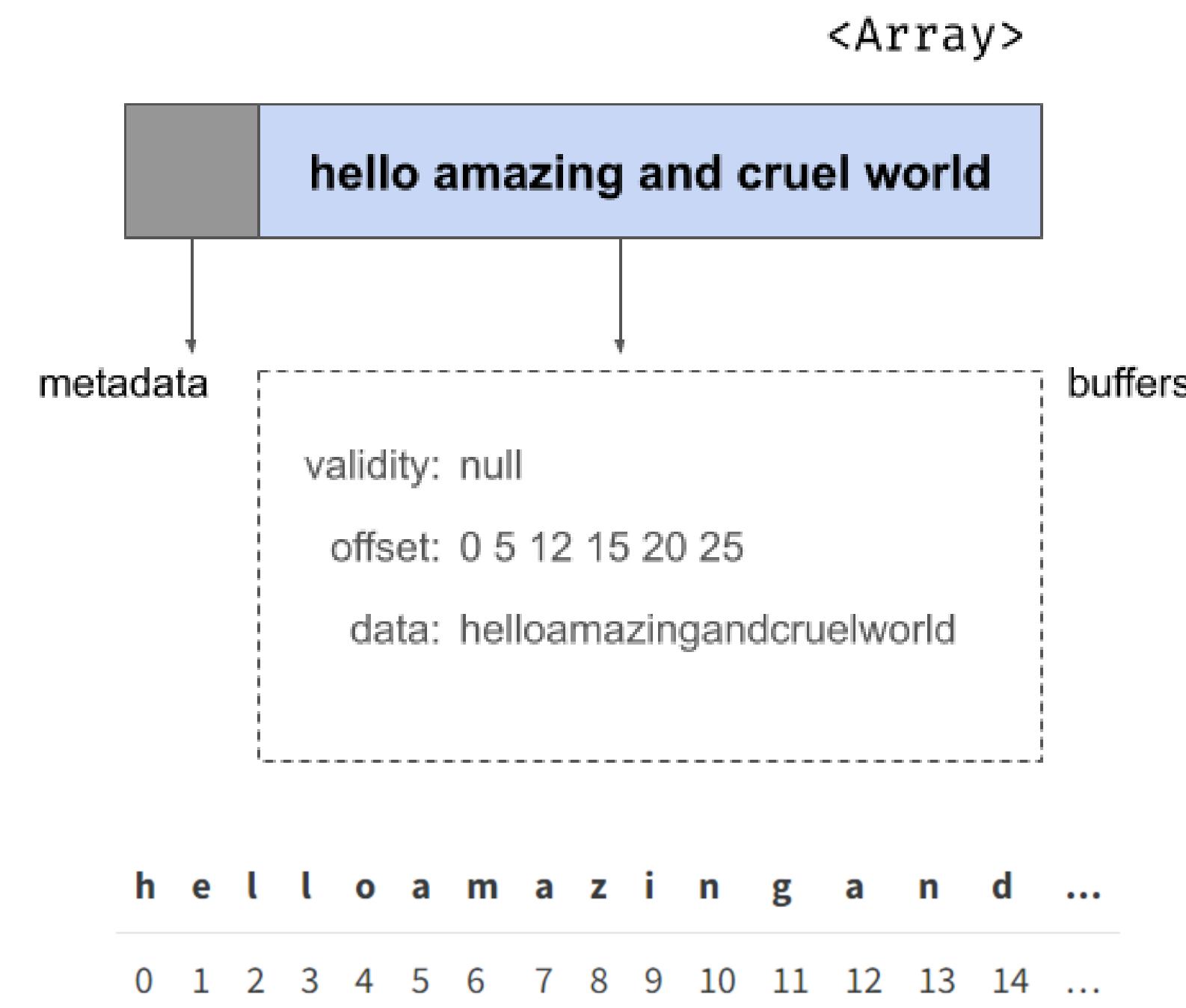
## Apache Arrow uses Validity bitmap buffer



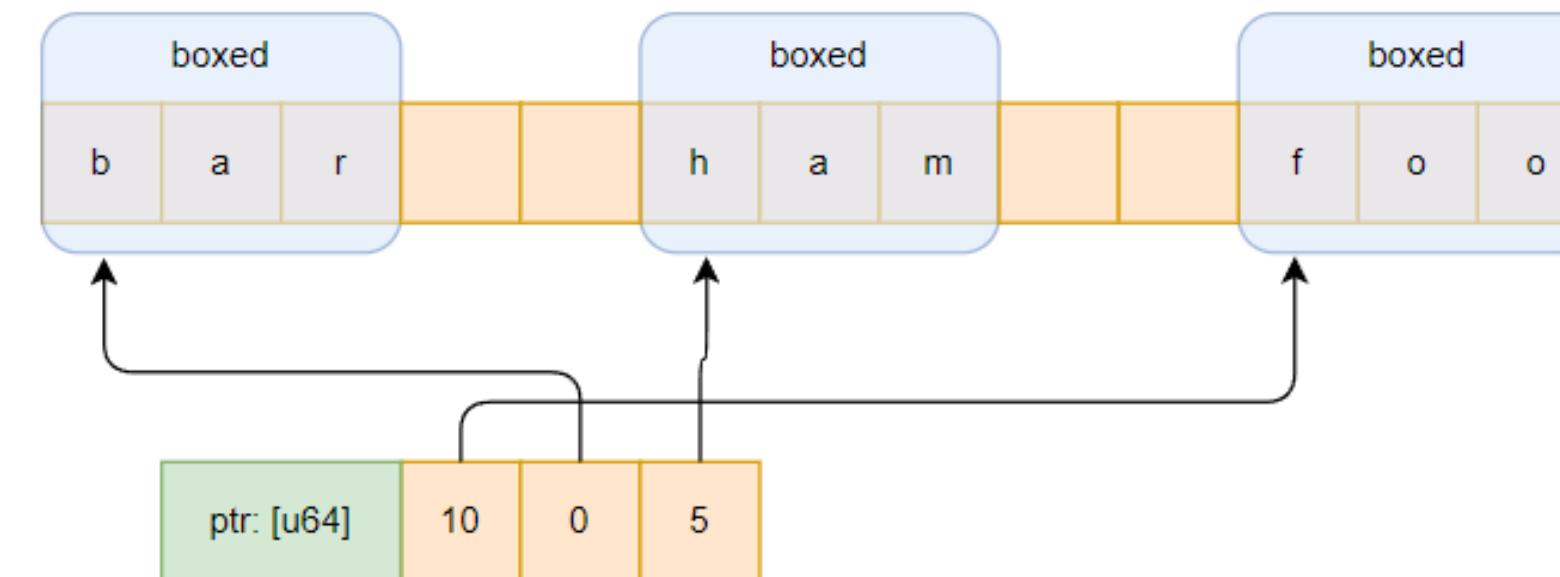
## Pandas 2 and Polars leverage Arrow Validity bitmap buffer to solve some Pandas problems

|   |  |
|---|--|
| Getting closer to the metal                         | Better groupby operations                    |
| No support for memory-mapped datasets               | DataFrame appending is more memory-efficient |
| Poor performance in database and file ingest/export | Extensible type metadata                     |
| Doing missing data right                            | Eager evaluation model, no query planning    |
| Weak support for categorical data                   | "Slow", limited multicore algorithms         |

## Arrow string array is stored sequentially



Pandas string array  
(Python object array)



Polars string array  
(Arrow LargeString array)

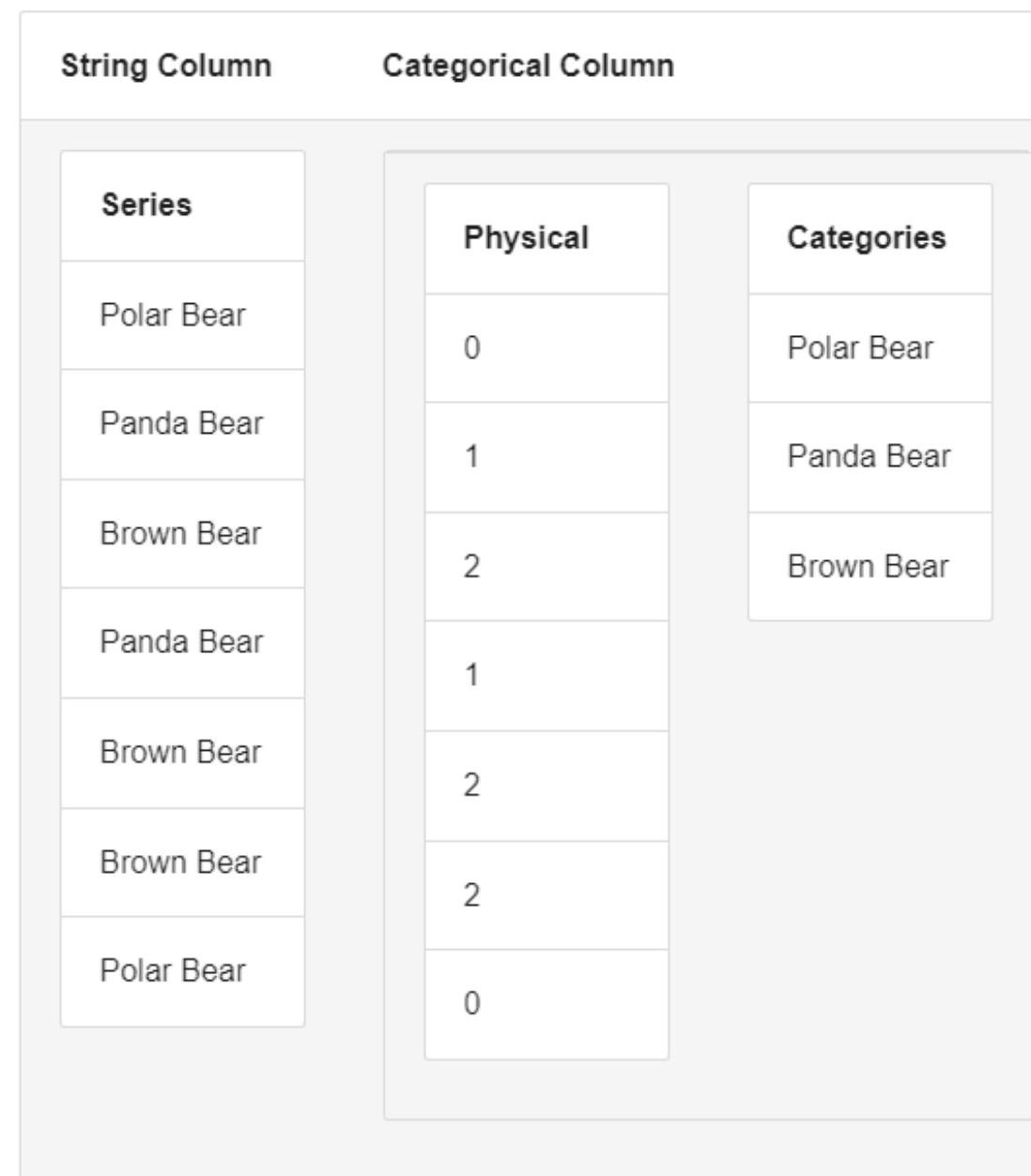
|                |      |   |   |   |   |   |   |   |   |
|----------------|------|---|---|---|---|---|---|---|---|
| data: [str]    | f    | o | o | b | a | r | h | a | m |
| offsets: [i64] | 0    | 2 | 5 | 8 |   |   |   |   |   |
| validity bits  | null |   |   |   |   |   |   |   |   |

## Pandas 2 and Polars leverage Arrow String to solve one more Pandas problem

|   |  |
|---|--|
| Getting closer to the metal                         | Better groupby operations                    |
| No support for memory-mapped datasets               | DataFrame appending is more memory-efficient |
| Poor performance in database and file ingest/export | Extensible type metadata                     |
| Doing missing data right                            | Eager evaluation model, no query planning    |
| Weak support for categorical data                   | "Slow", limited multicore algorithms         |

# Apache Arrow

## Supports categorical data well with dictionary



Memory-efficient  
Performant

## Pandas 2 and Polars leverage Arrow Dictionary Type to solve one more Pandas problem

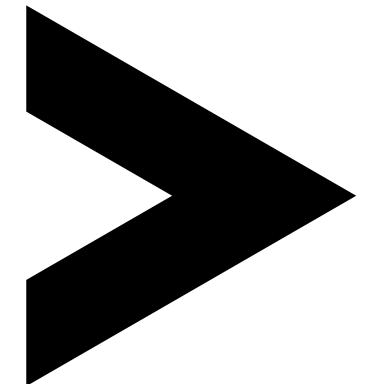
|   |  |
|---|--|
| Getting closer to the metal                         | Better groupby operations                    |
| No support for memory-mapped datasets               | DataFrame appending is more memory-efficient |
| Poor performance in database and file ingest/export | Extensible type metadata                     |
| Doing missing data right                            | Eager evaluation model, no query planning    |
| Better support for categorical data                 | "Slow", limited multicore algorithms         |

## Arrow array supports memory-mapped datasets

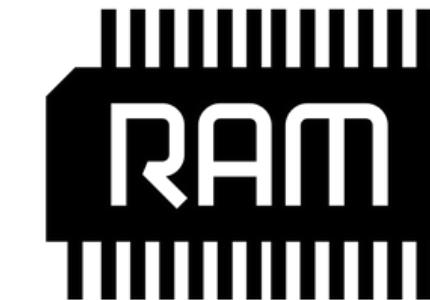
Dataset



50GB

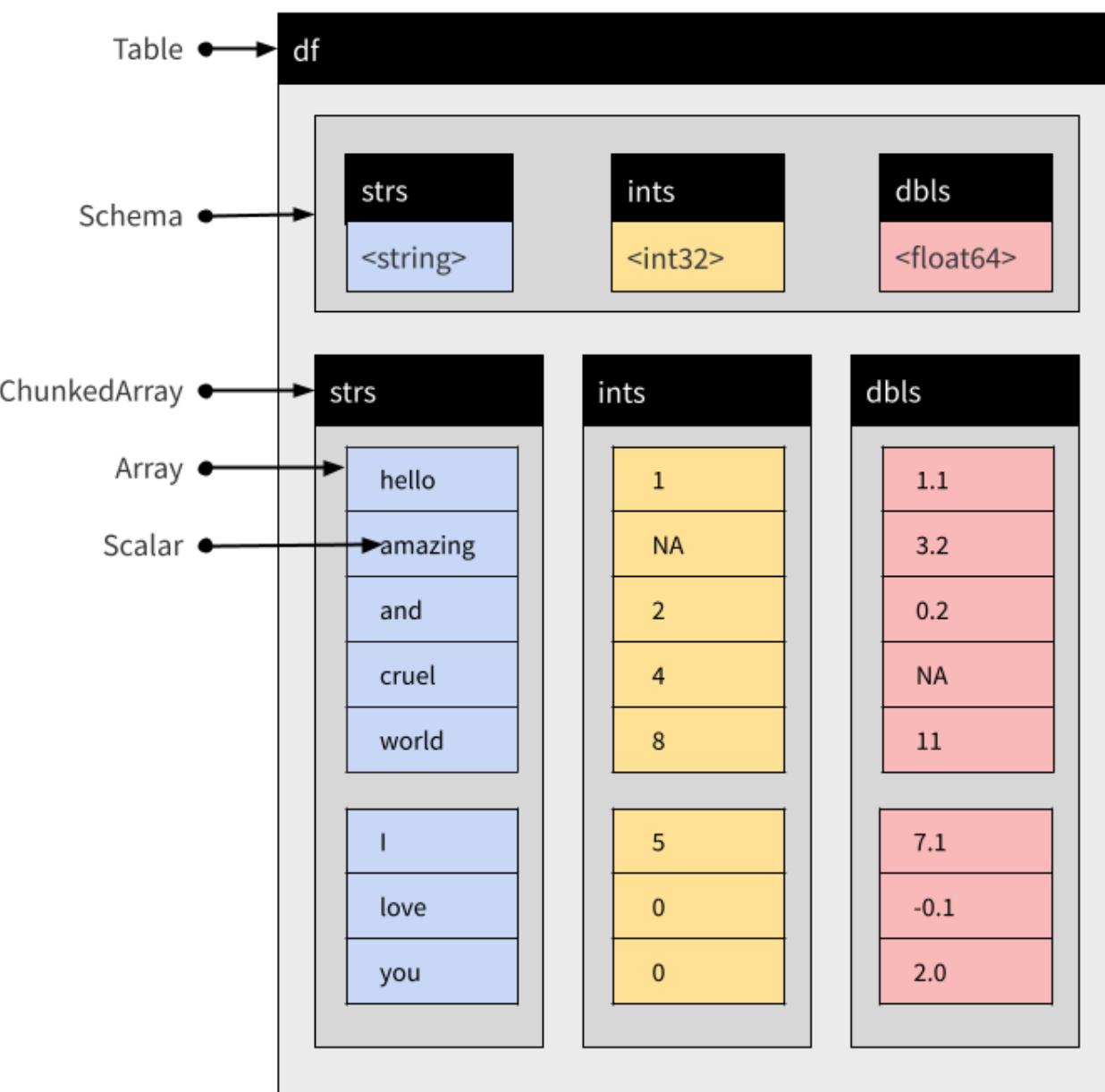


RAM

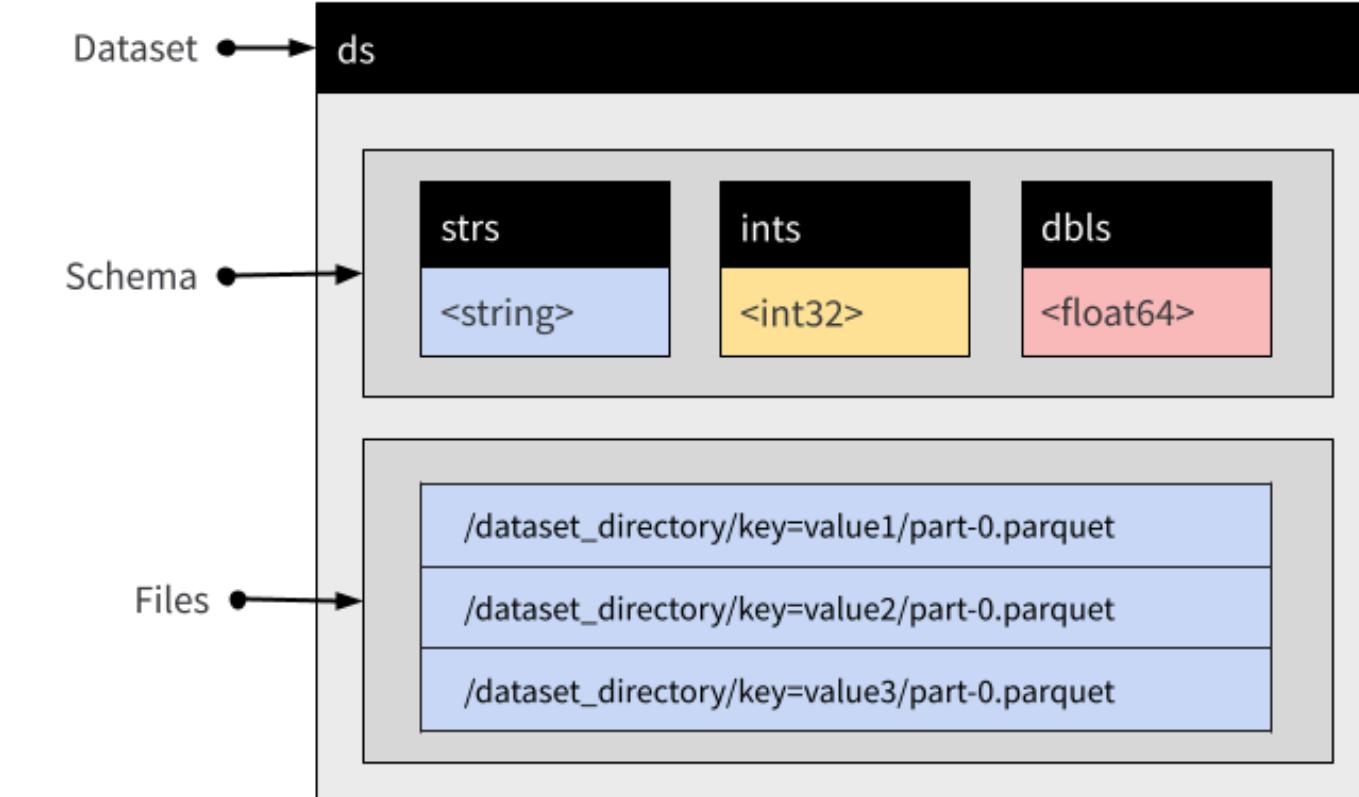


16GB

# Arrow array supports memory-mapped datasets

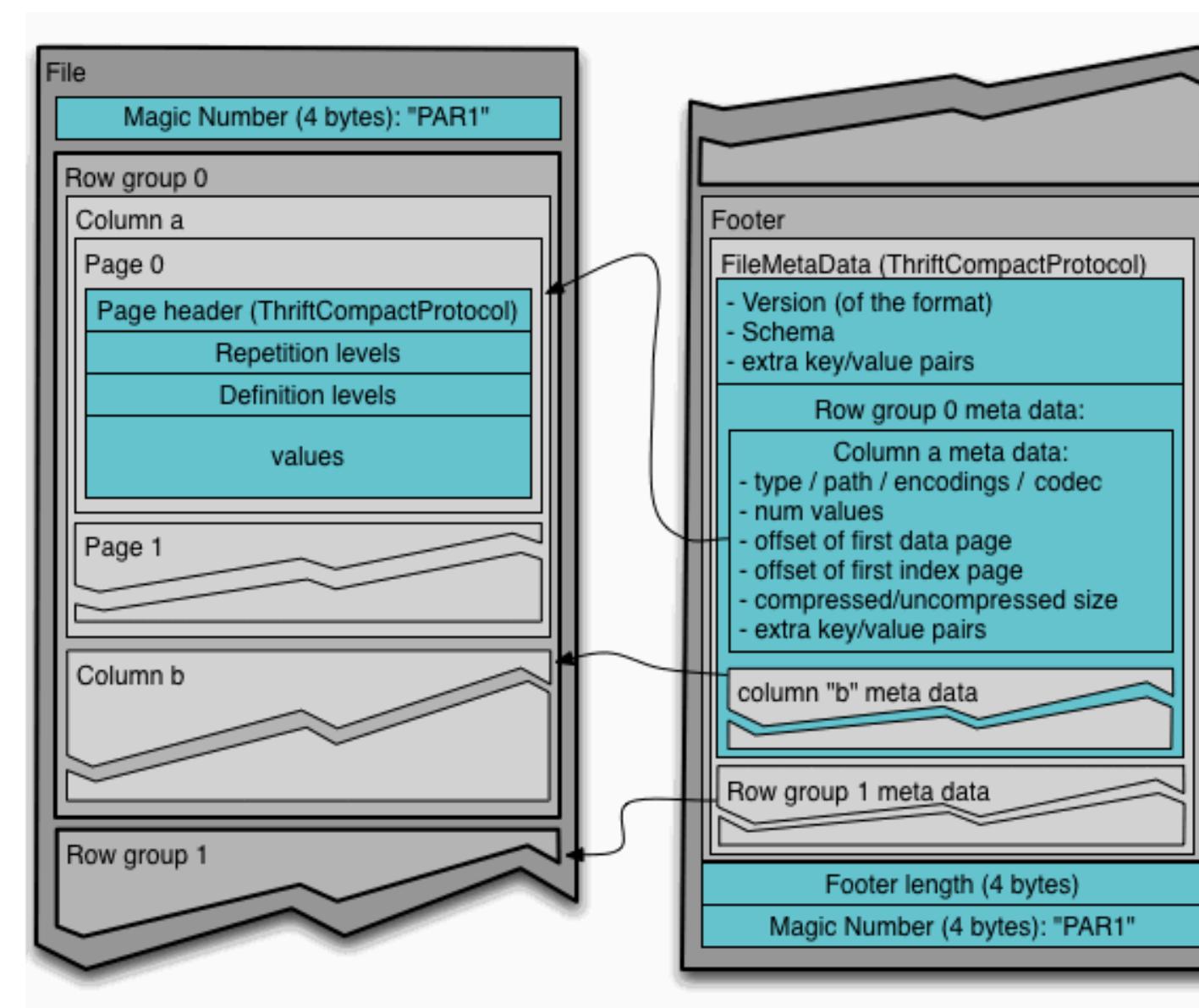


In-memory dataframe



On-disk dataset

## Arrow array supports memory-mapped datasets



Parquet file structure

## Pandas 2 and Polars leverage memory-mapped dataset support of Arrow to solve one more Pandas problem

|   |  |
|---|--|
| Getting closer to the metal               | Better groupby operations                        |
| <b>Support for memory-mapped datasets</b> | DataFrame appending is more memory-efficient     |
| High speed data ingest and export         | Extensible type metadata                         |
| Doing missing data right                  | <b>Eager evaluation model, no query planning</b> |
| Better support for categorical data       | <b>"Slow", limited multicore algorithms</b>      |

# Query planning

## Eager evaluation (imperative)

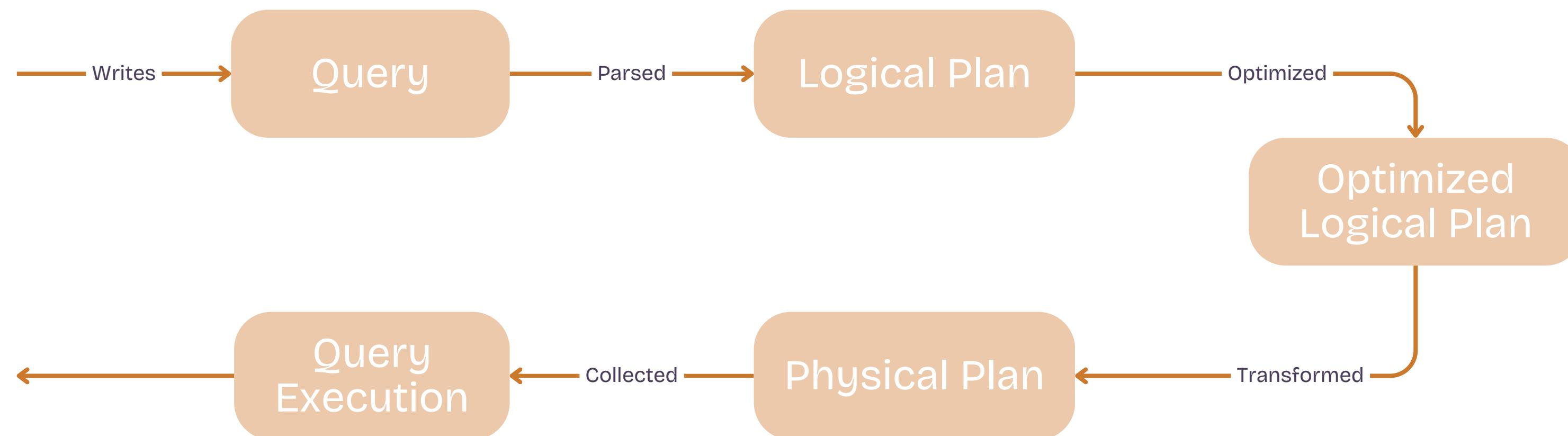
The expression  
is evaluated  
as soon as  
it is called.



## Lazy evaluation (declarative)

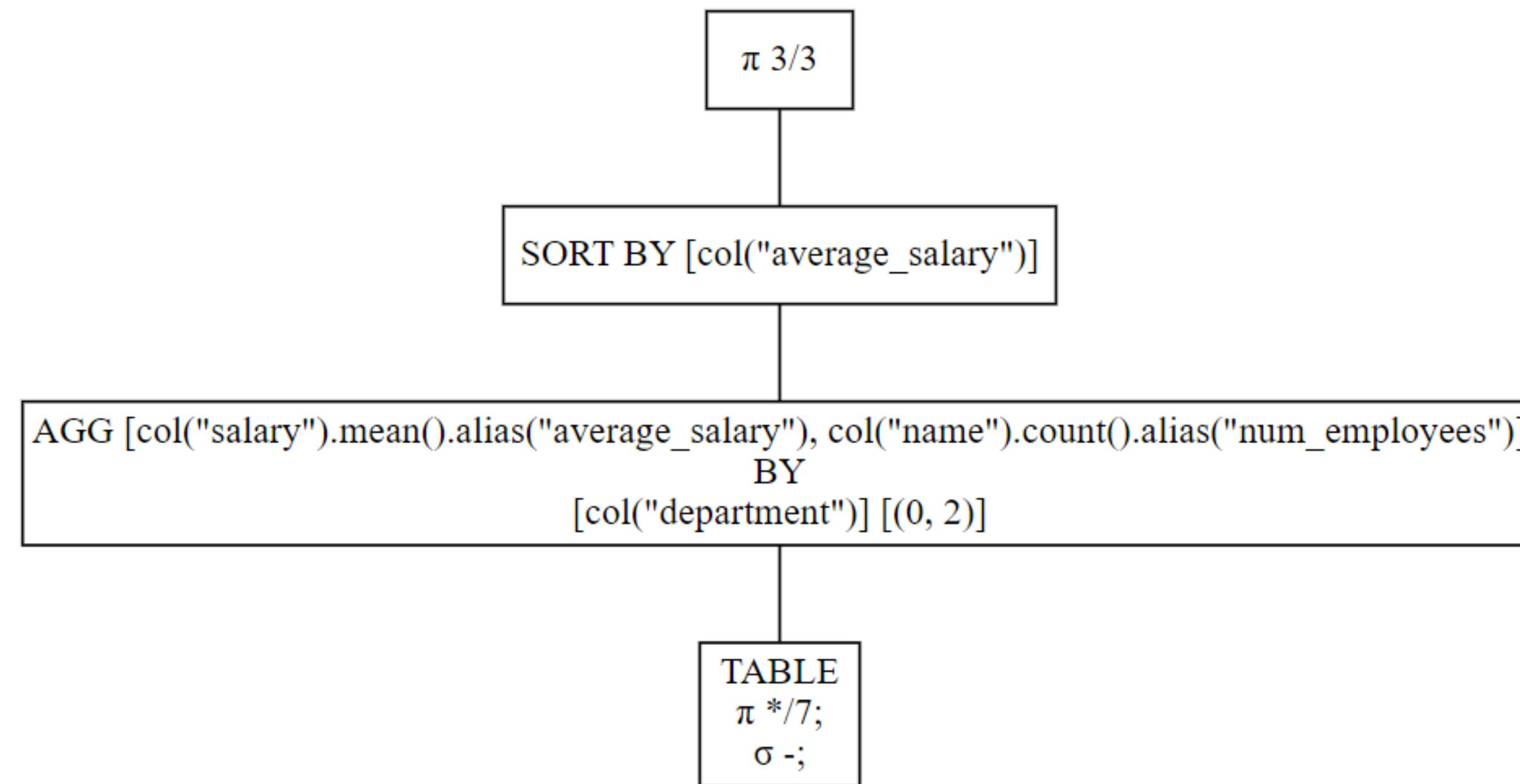
The expression  
is stored as  
computational  
graph to be  
evaluated  
when needed.

# Query planning

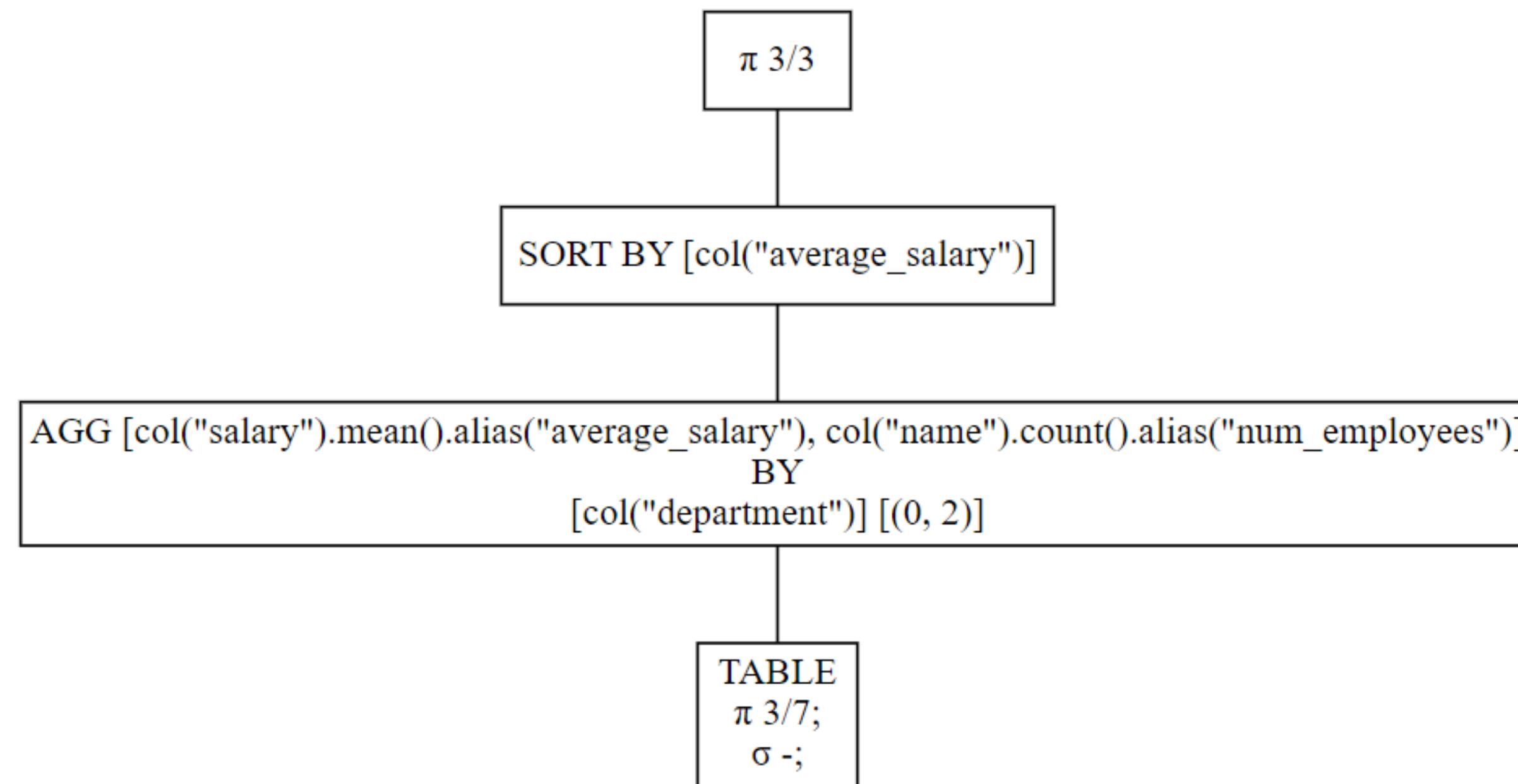


`query.collect()`

|                  |   |   |
|------------------|---|---|
|                  |  pandas |   |
| Eager evaluation |        |  |
| Lazy evaluation  |        |  |



```
query.show_graph(optimized=False)
```

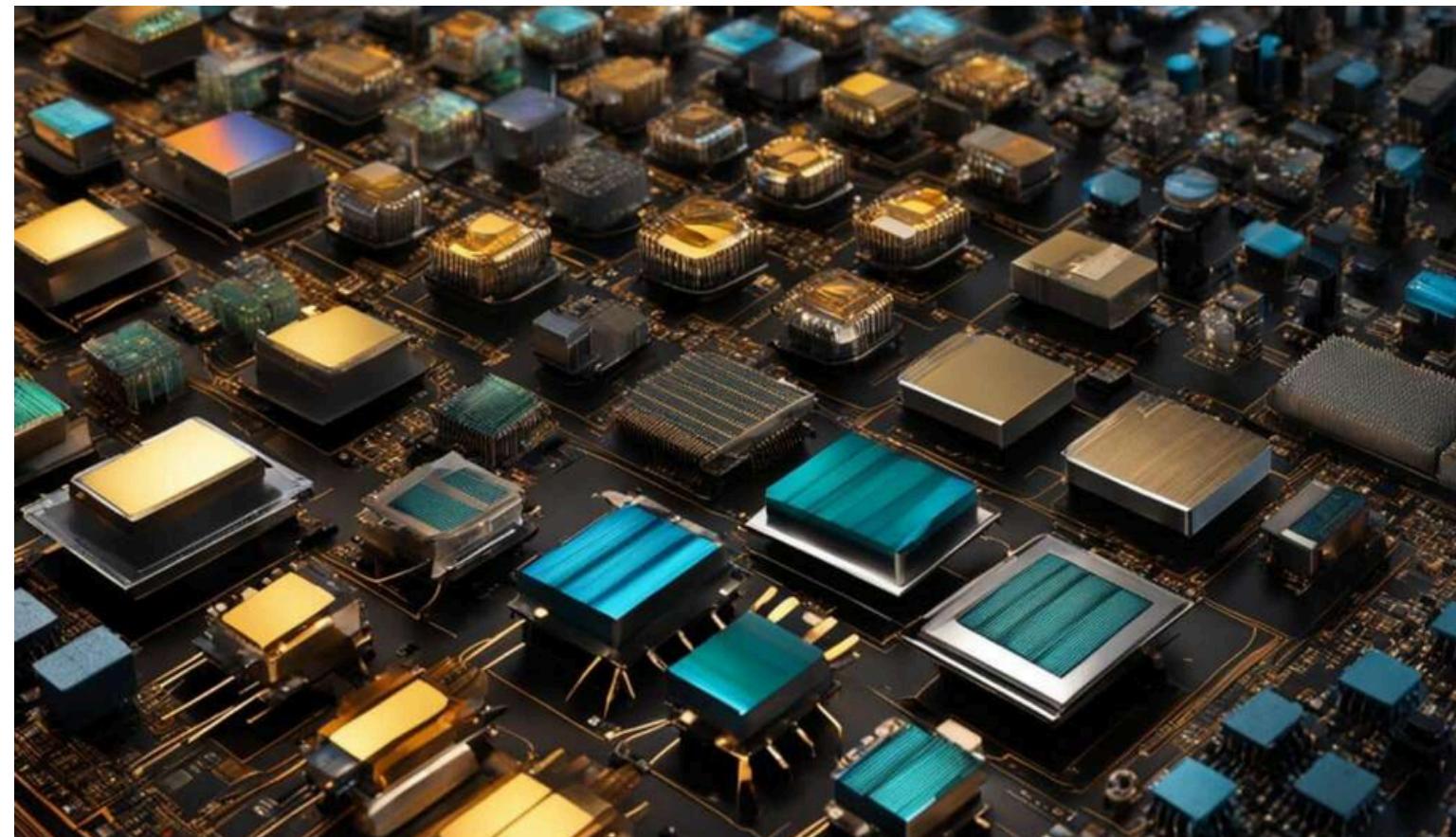


```
query.show_graph(optimized=True)
```

## Polars supports **Query Planning**, which is a problem of Pandas (even Pandas 2)

|                                     |  |
|-------------------------------------|--|
| Getting closer to the metal         | Better groupby operations                          |
| Support for memory-mapped datasets  | DataFrame appending is more memory-efficient       |
| High speed data ingest and export   | Extensible type metadata                           |
| Doing missing data right            | <b>Support query planning with Lazy Evaluation</b> |
| Better support for categorical data | <b>"Slow", limited multicore algorithms</b>        |

# Parallelism

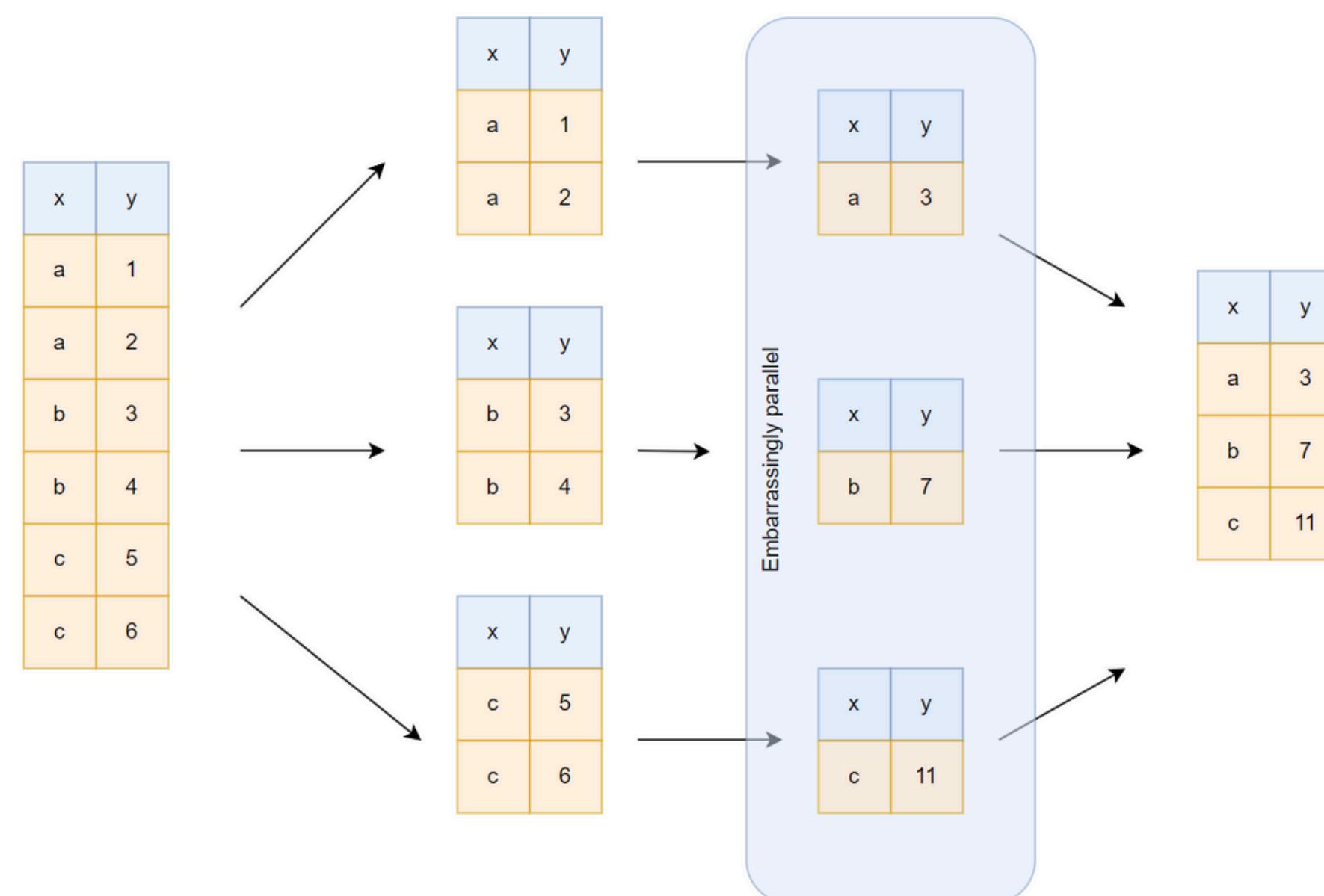


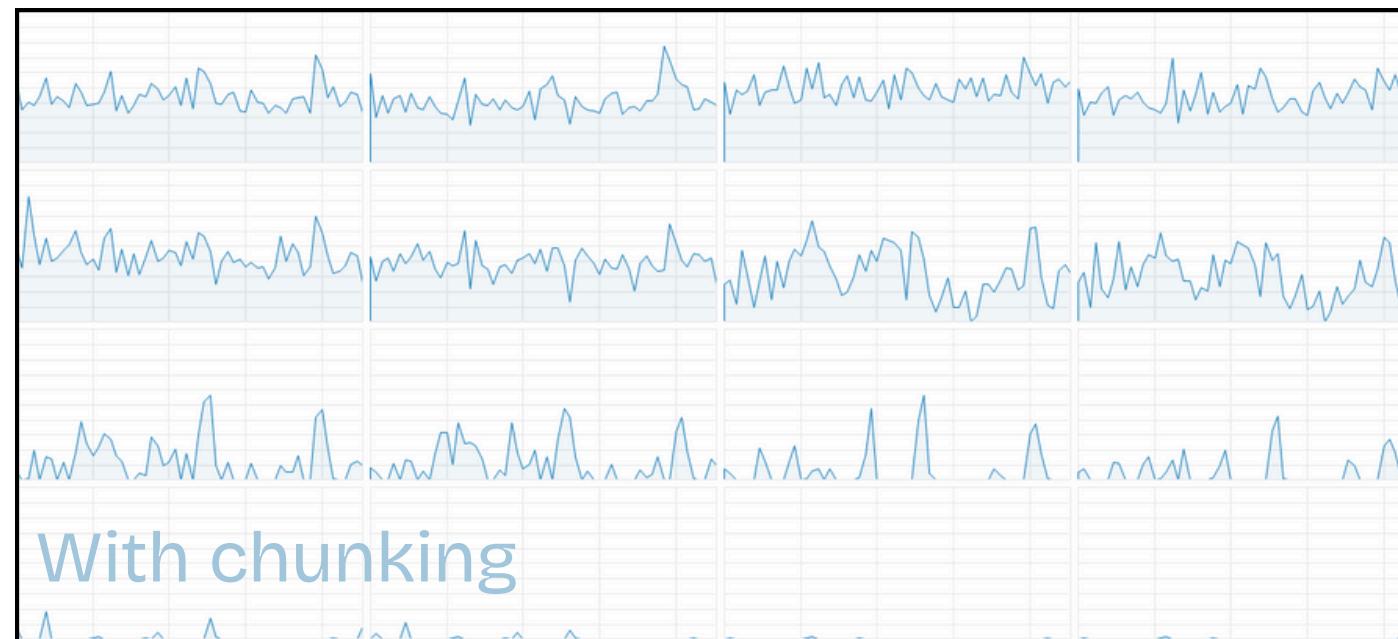
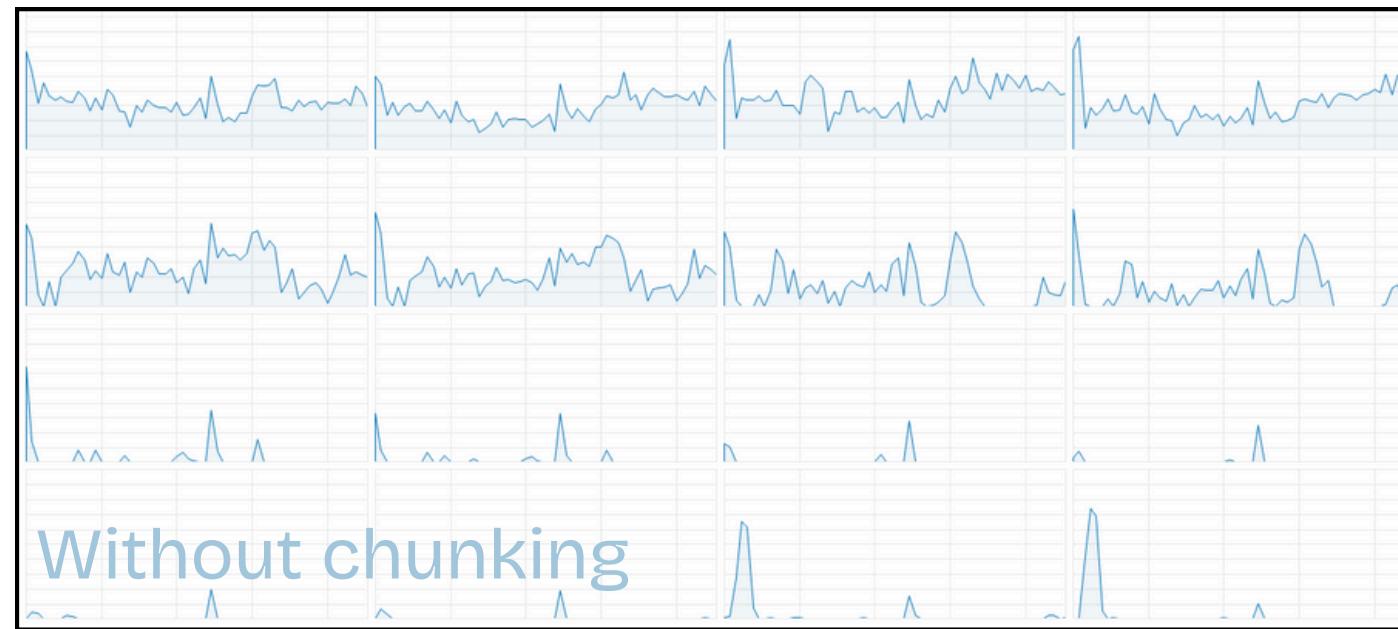
The advent of **multi-core processors** has revolutionized the world of computing with **parallel** processing.

**Pandas** does not take advantage of parallelism.

**But Polars does!**

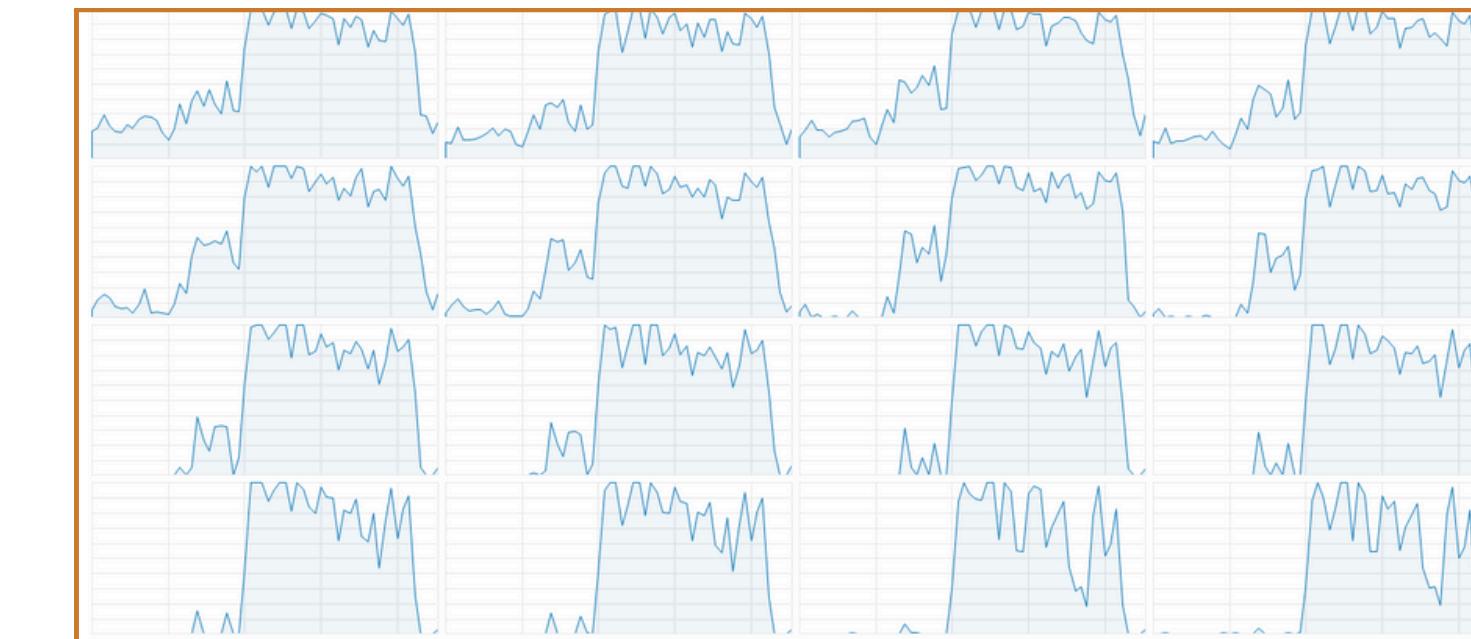






Pandas

# CPU utilization of **Pandas** and **Polars** in One Billion Row Challenge



Polars

## Polars supports **Parallelism very well**, which is a problem of Pandas (even Pandas 2)

|                                    |  |
|------------------------------------|--|
| Getting closer to the metal        | Better groupby operations                    |
| Support for memory-mapped datasets | DataFrame appending is more memory-efficient |
| High speed data ingest and export  | Extensible type metadata                     |
| Doing missing data right           | Support query planning                       |
| Weak support for categorical data  | <b>Leverage parallelism</b>                  |

# Polars

Apache Arrow

Query planning

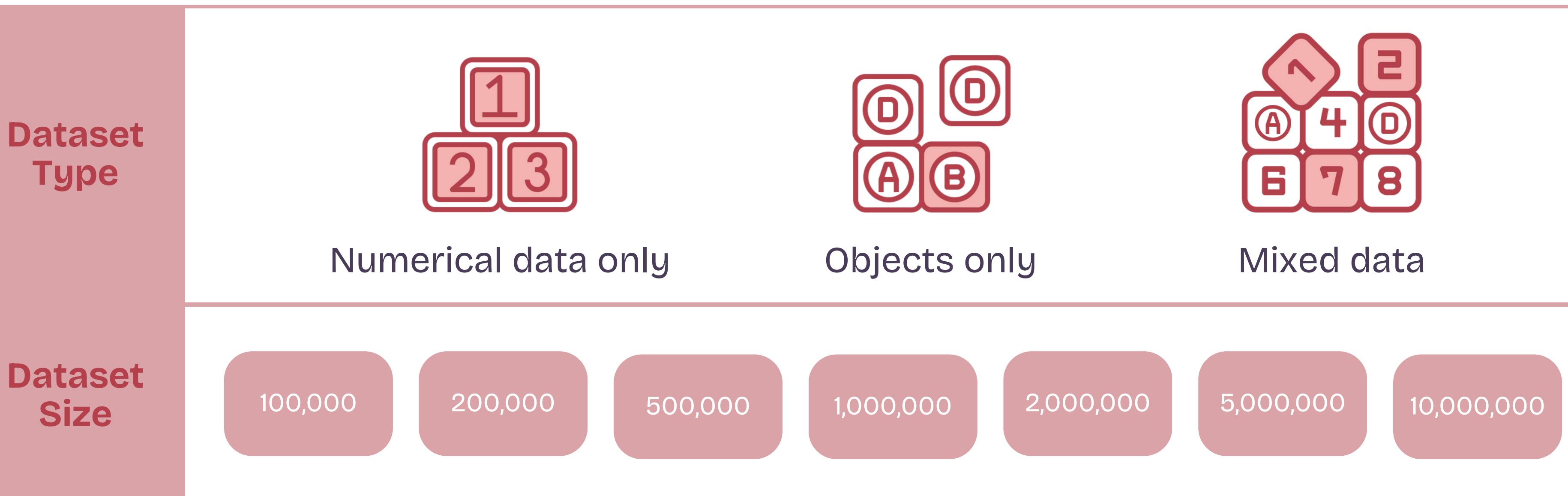
Parallelism

Columnar format and SIMD, Validity buffer, Memory mapping, Categorical Type



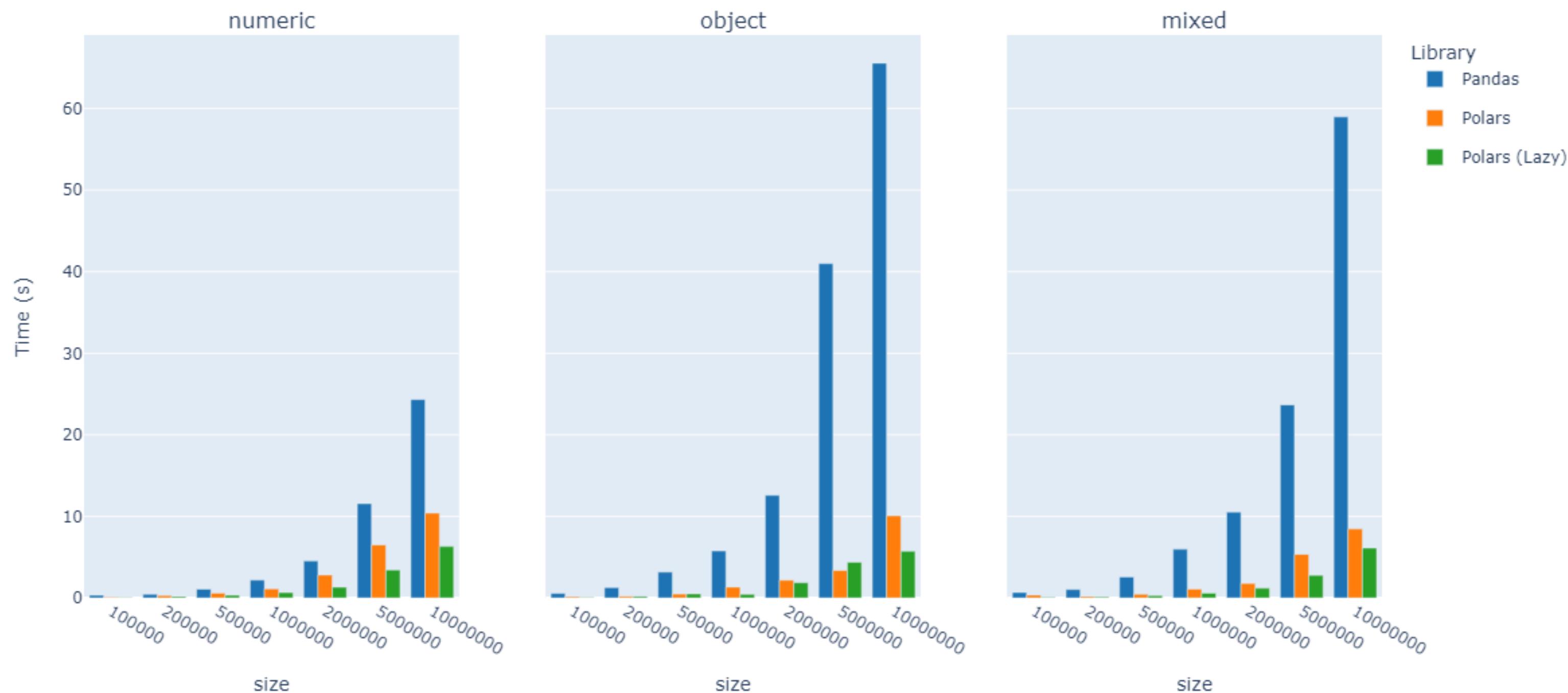
# 03 COMPARISON

## Compare run time of Pandas 1 and Polars on datasets



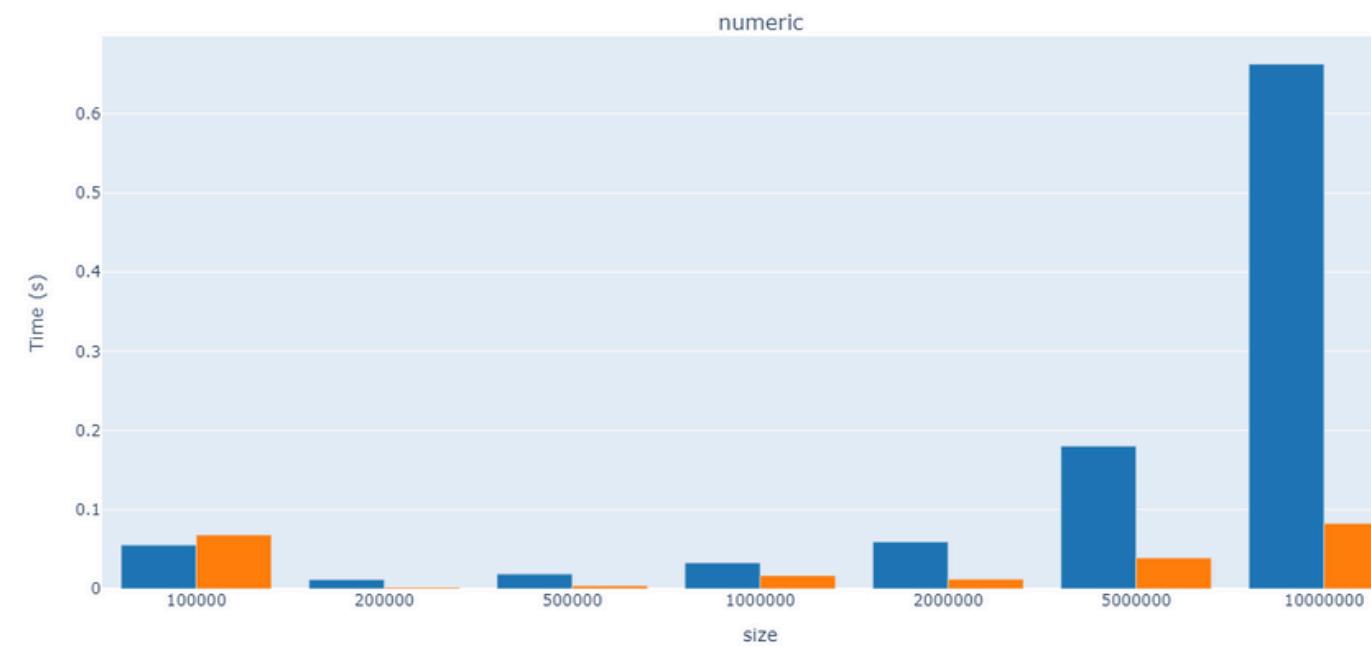
# Read csv file

Comparison of Pandas and Polars Read Time

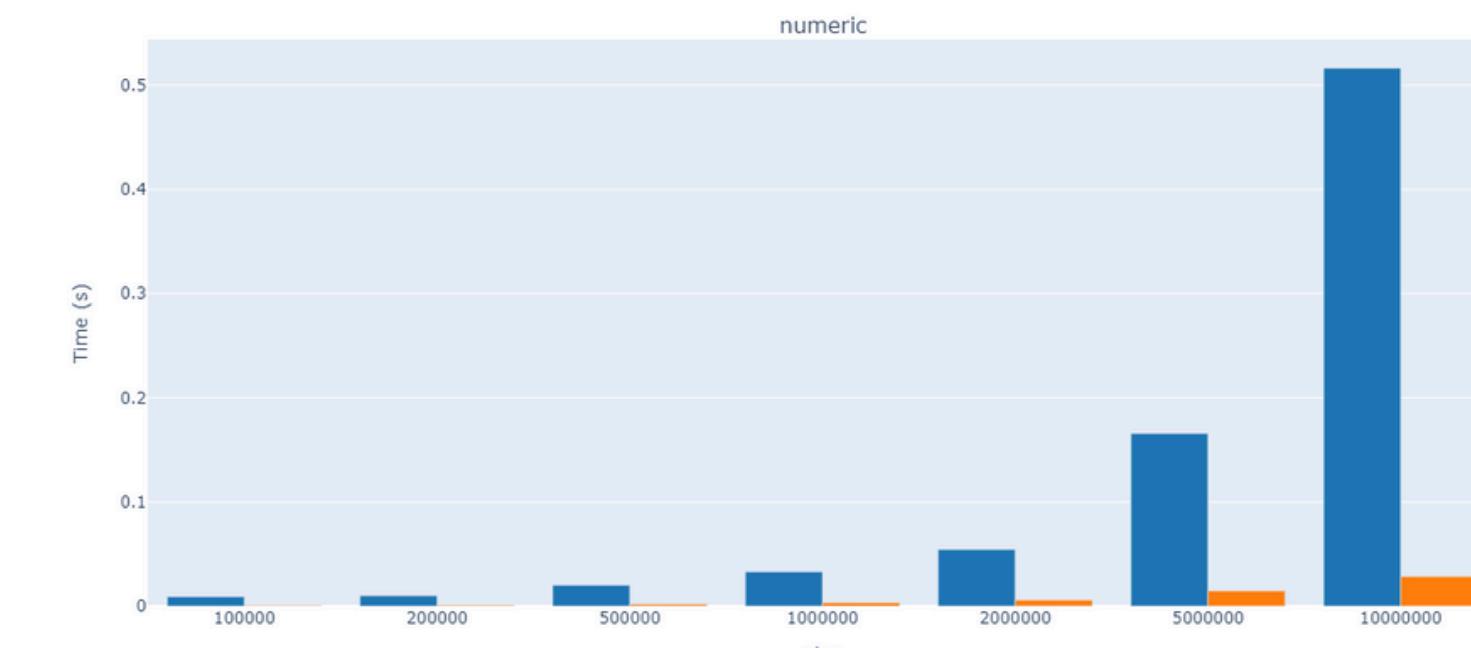


# Aggregation

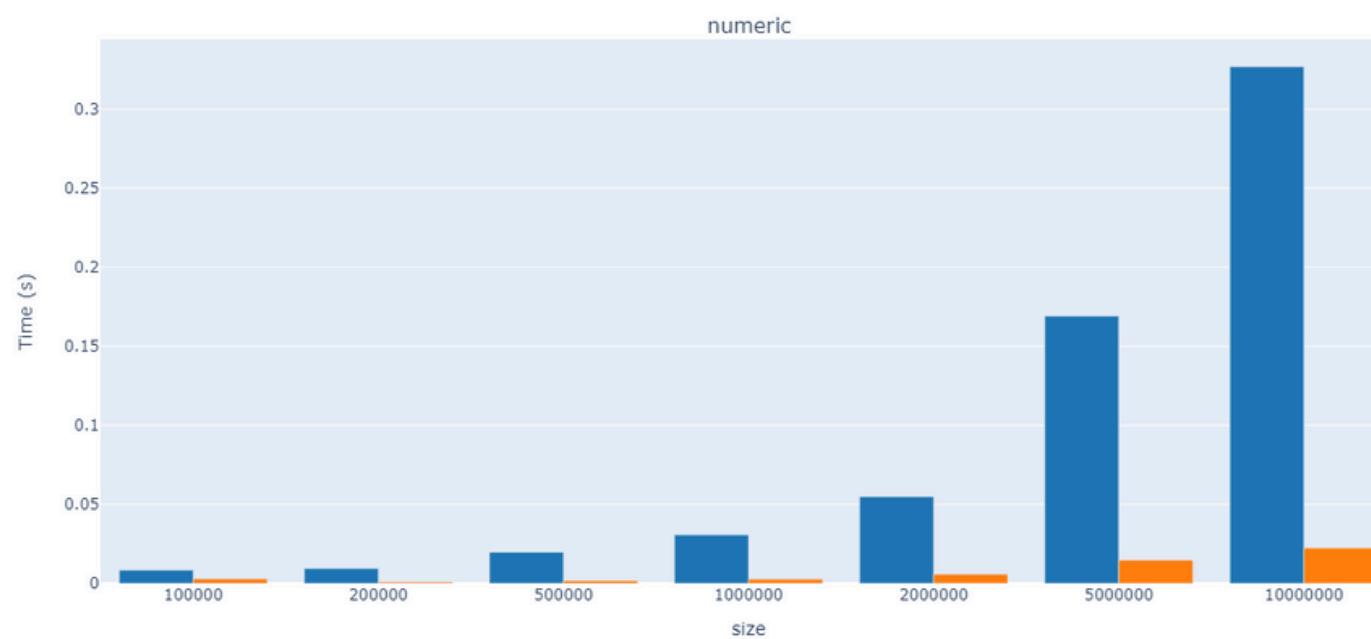
Min



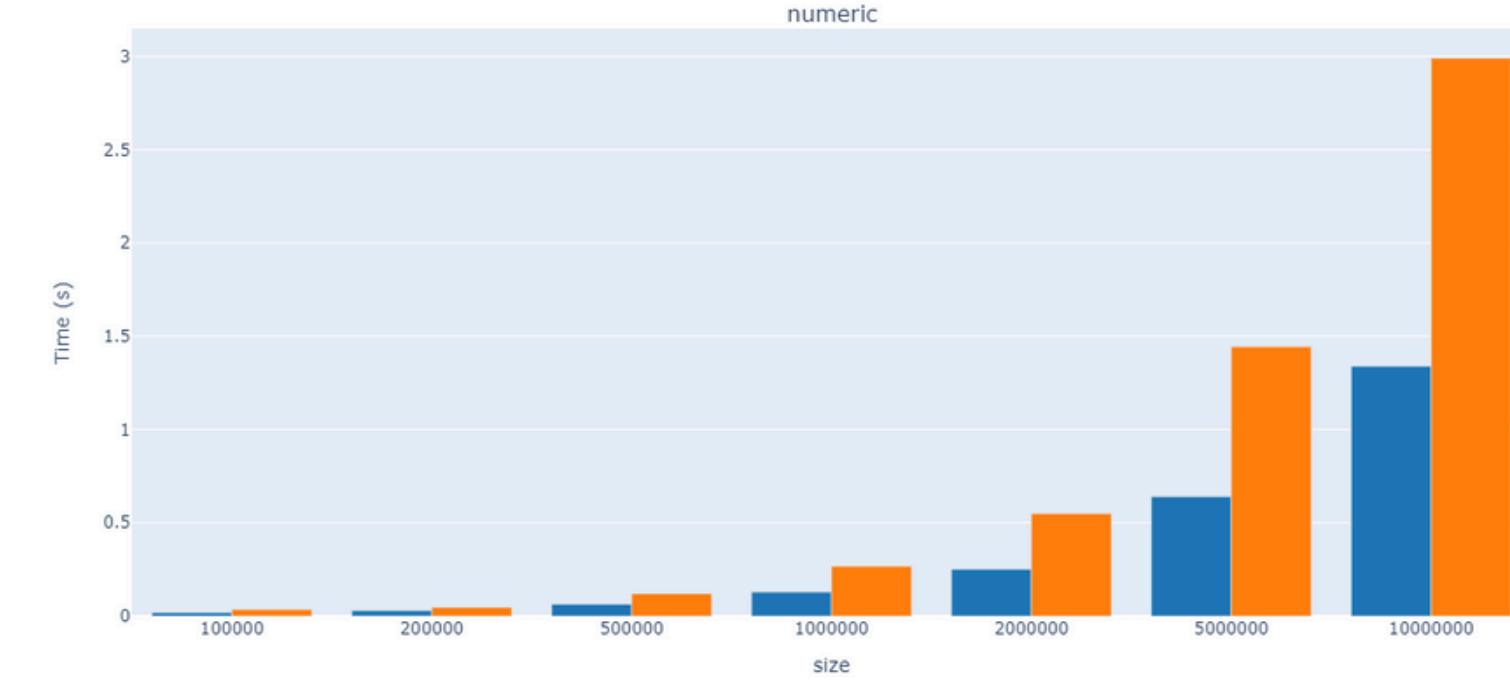
Max



Mean

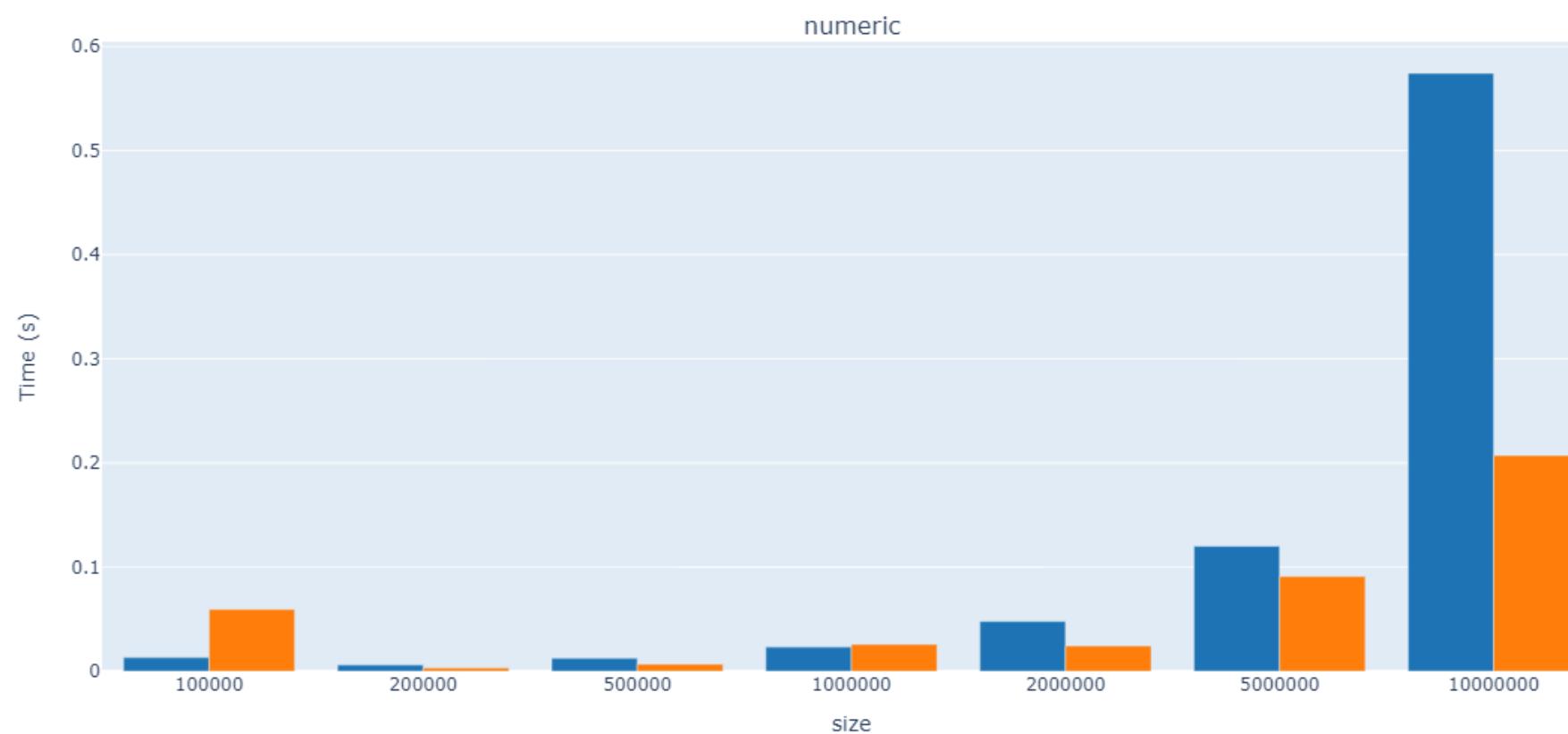


Median



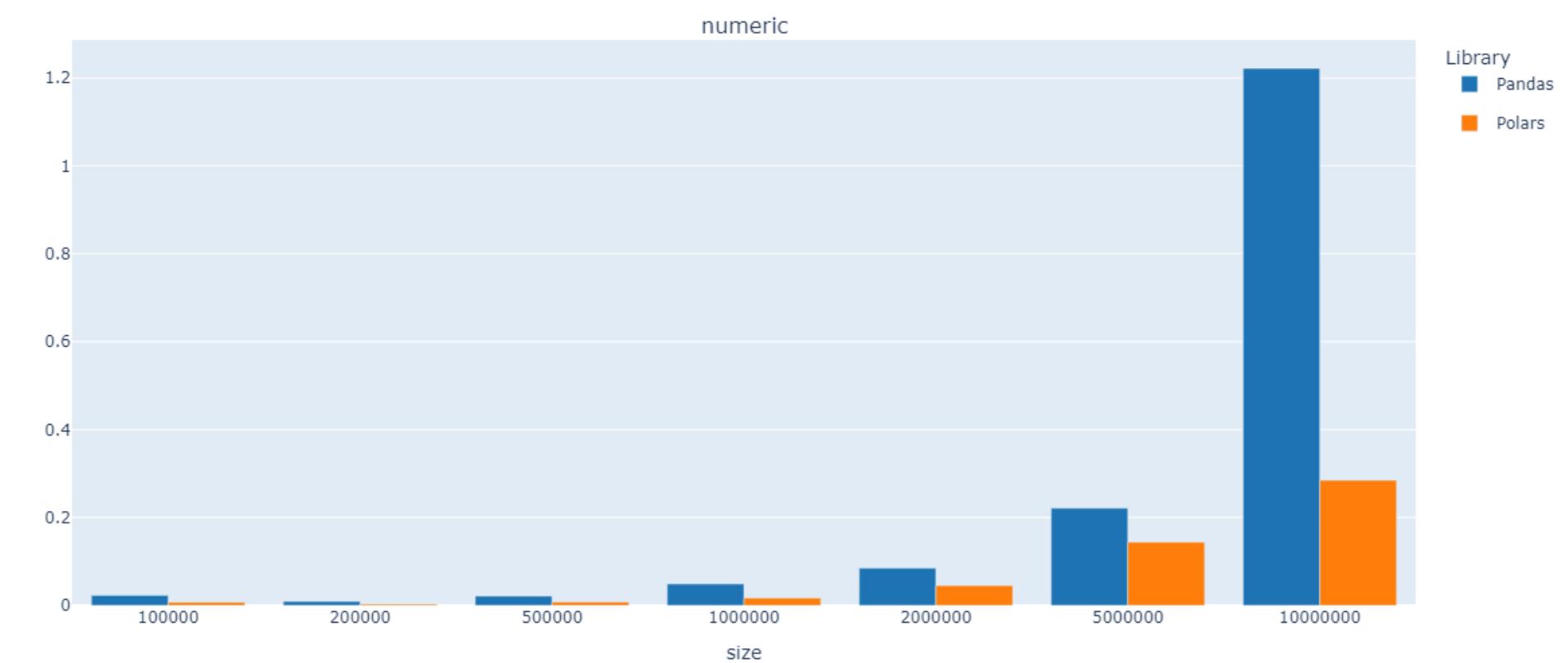
# Query

Comparison of Pandas and Polars Query 1 Time



Query 1: Filter

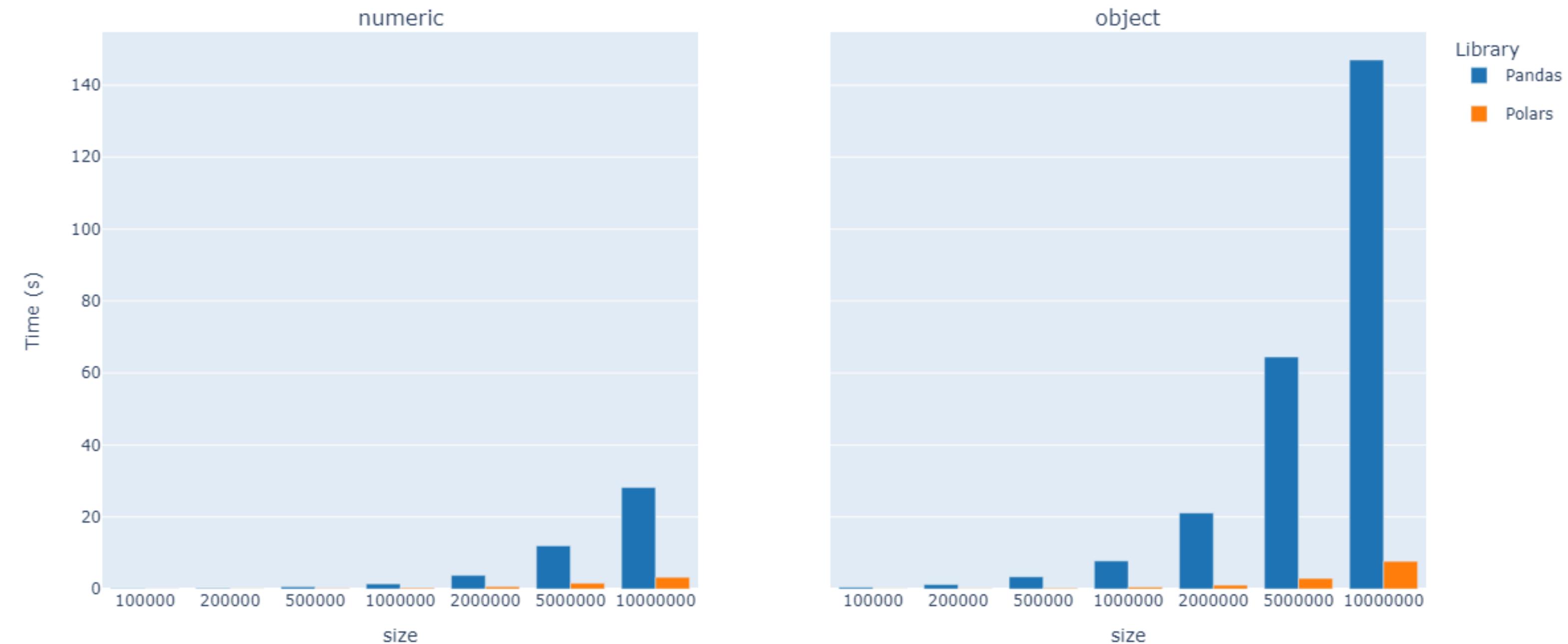
Comparison of Pandas and Polars Query 2 Time



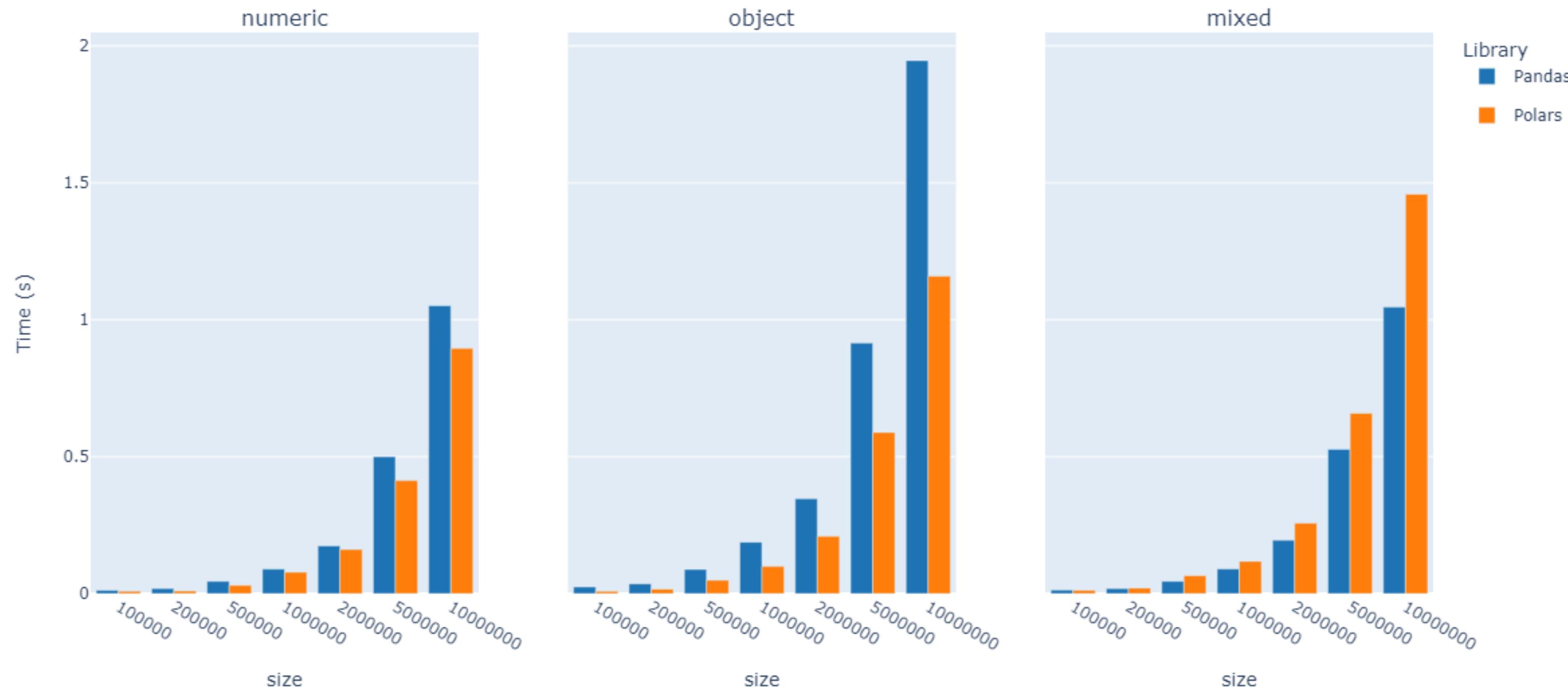
Query 2: Filter + Select

# Sort

Comparison of Pandas and Polars Sort Time

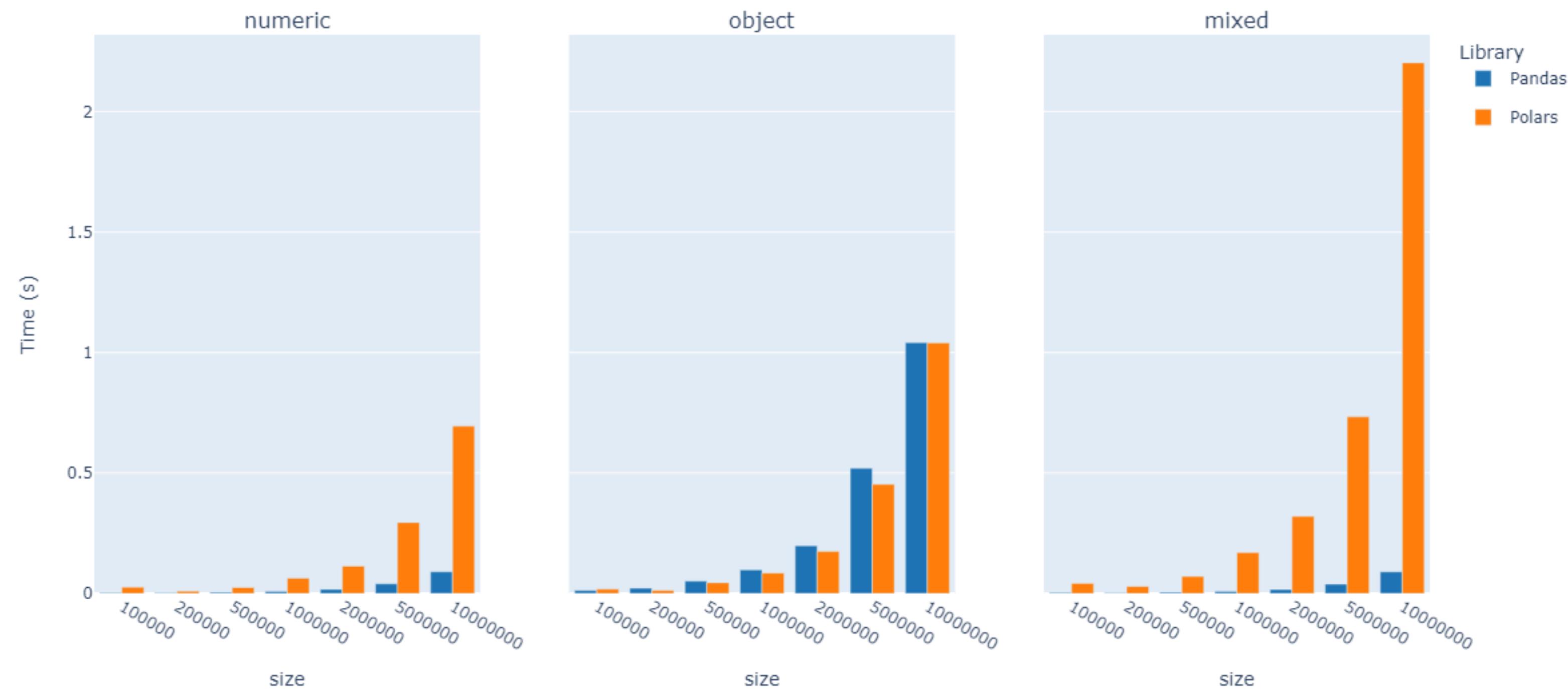


Comparison of Pandas and Polars Grouping String Time



## Groupby on string column

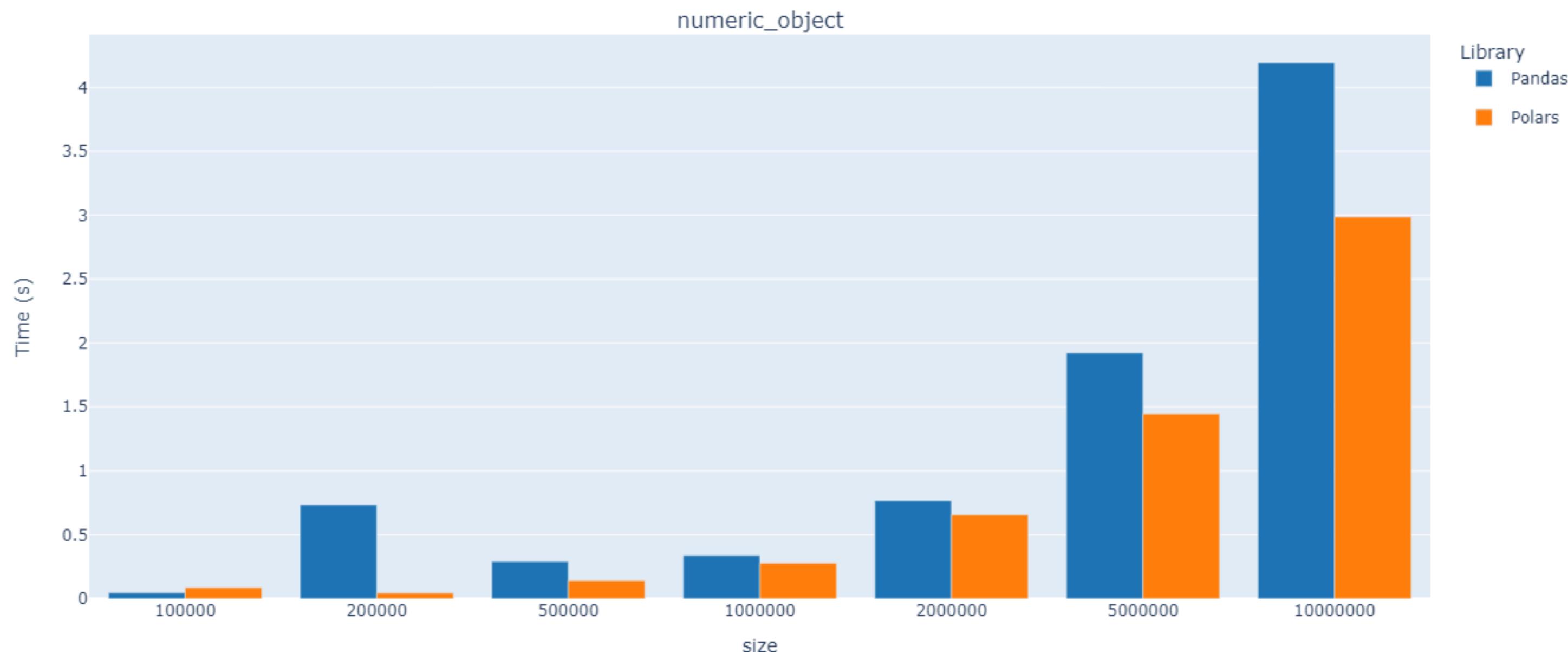
Comparison of Pandas and Polars Grouping Category Time



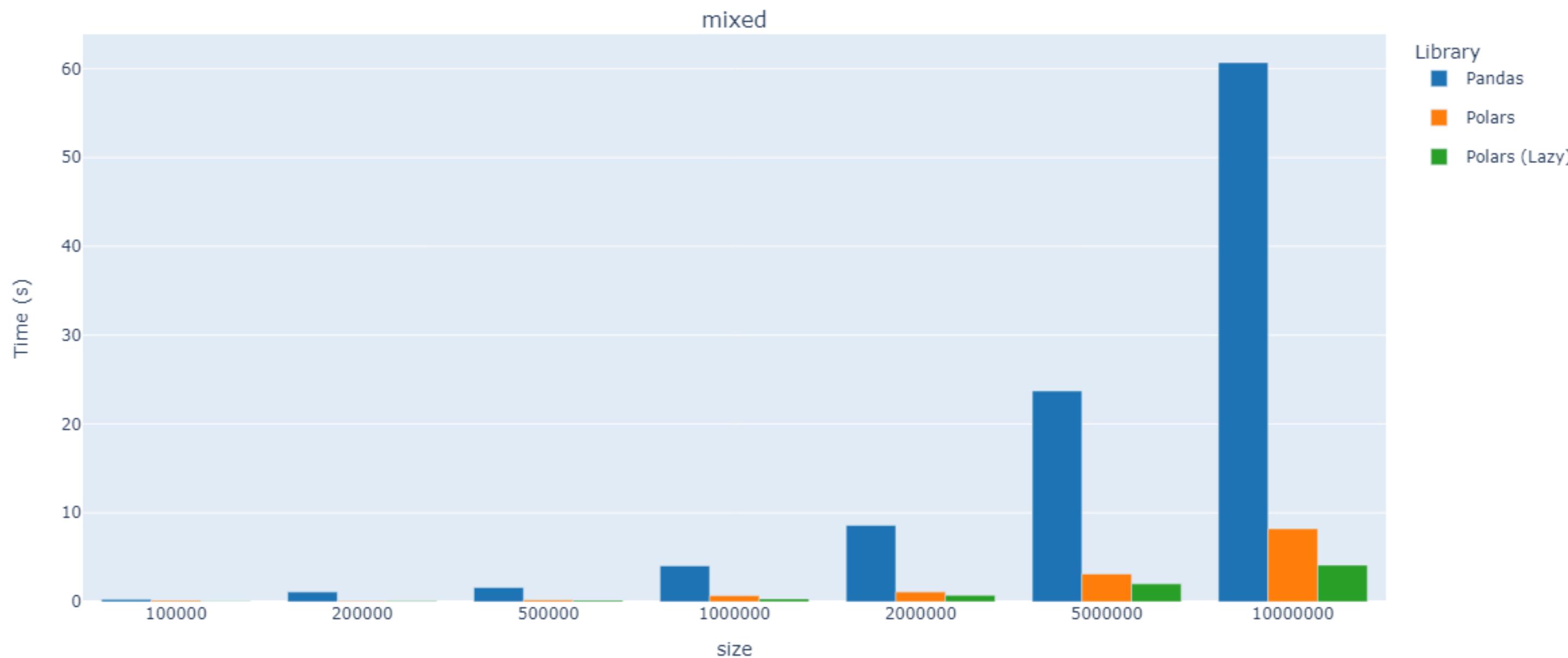
## Groupby on categorical column

# Join

Comparison of Pandas and Polars Join Time



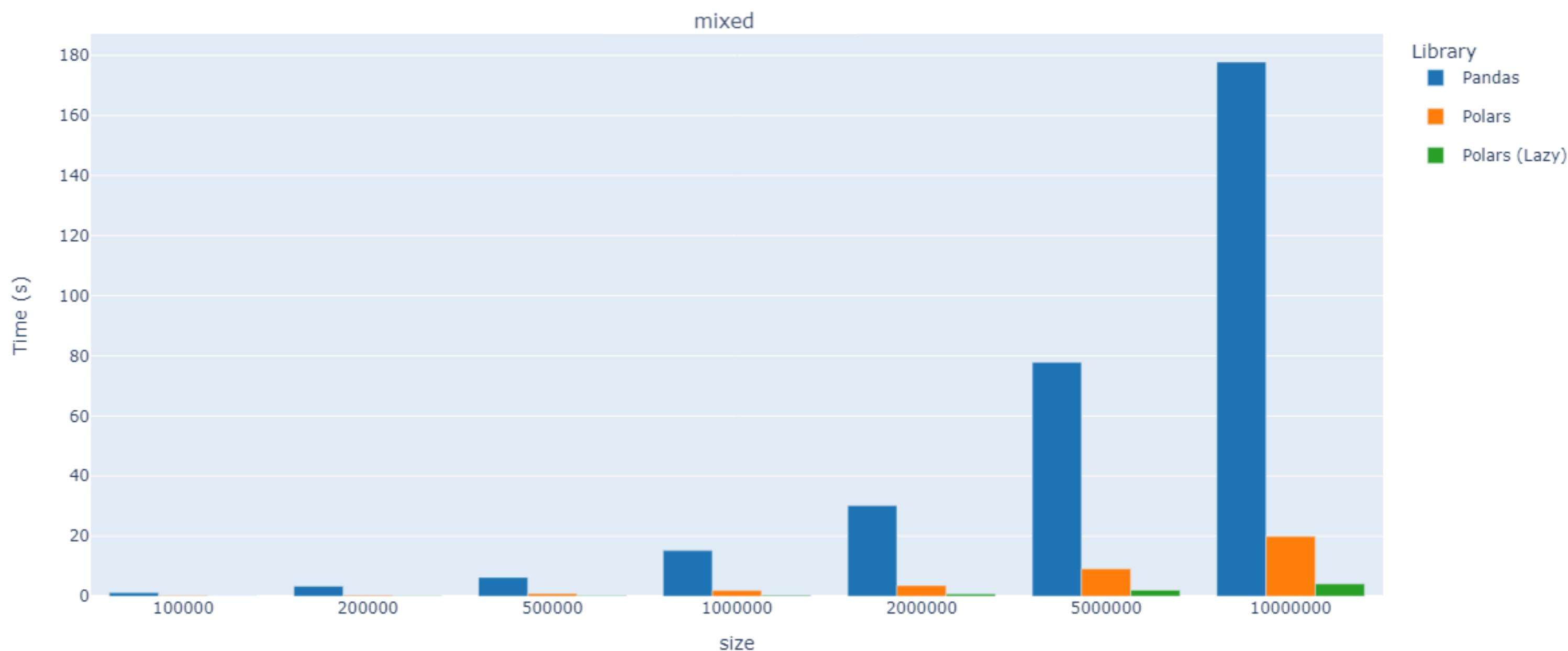
Comparison of Pandas and Polars Manipulation Time



## Compound Manipulation

## Read + Manipulation

Comparison of Pandas and Polars Read and Manipulation Time





04

# TUTORIAL

# 05 CONCLUSION

# Pros | Polars over Pandas 1

|  |   |
|--|---|
| <b>Close to the metal</b>                  | <b>Better groupby operations</b>                    |
| <b>Support for memory-mapped datasets</b>  | <b>DataFrame appending is more memory-efficient</b> |
| <b>High speed data ingest and export</b>   | <b>Extensible type metadata</b>                     |
| <b>Doing missing data right</b>            | <b>Support query planning</b>                       |
| <b>Better support for categorical data</b> | <b>Levarage parallelism</b>                         |

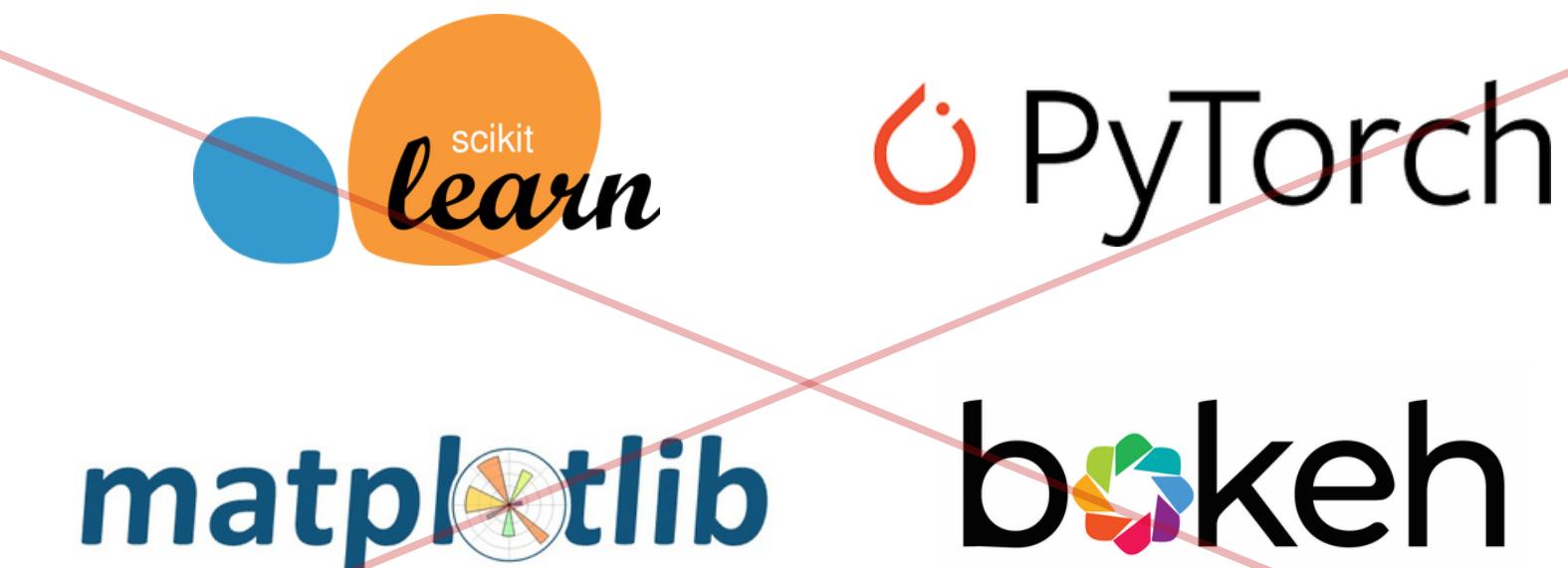
## Pros | Polars over Pandas 2

|                                     |  |
|-------------------------------------|--|
| Close to the metal                  | Better groupby operations                    |
| Support for memory-mapped datasets  | DataFrame appending is more memory-efficient |
| High-speed data ingest and export   | Extensible type metadata                     |
| Doing missing data right            | <b>Support query planning</b>                |
| Better support for categorical data | <b>Leverage parallelism</b>                  |

## Cons | Polars has limited ecosystem

 **plotly**

Only compatible with



Not compatible yet

`to_numpy()`

Thank you  
for your attention!

