

**DATA+AI**  
SUMMIT 2022

# Polars

An API for fast analytics



**Ritchie Vink**

Machine Learning Engineer, Xomnia

ORGANIZED BY  databricks

# About



## Ritchie Vink

- Author of Polars
- Background in Machine Learning and Software development
- Part-time open source developer @ Xomnia



# Agenda

- Why polars?
- Foundations
- Performance
- Expression API
- Small example



# Current DataFrame implementations

\*Python oriented

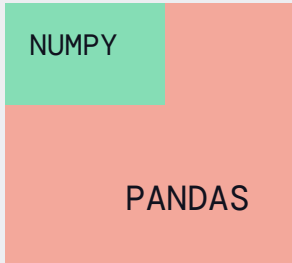
50+ years of RDMS design are not really applied in the data science community

- Almost all implementation are eager → no query optimization
- Huge wasteful materializations
- Responsibility on fast/memory efficient compute on user (most users are no OLAP experts)
- Terrible memory representation of string data == terrible performance



# Current DataFrame implementations

\*Python oriented



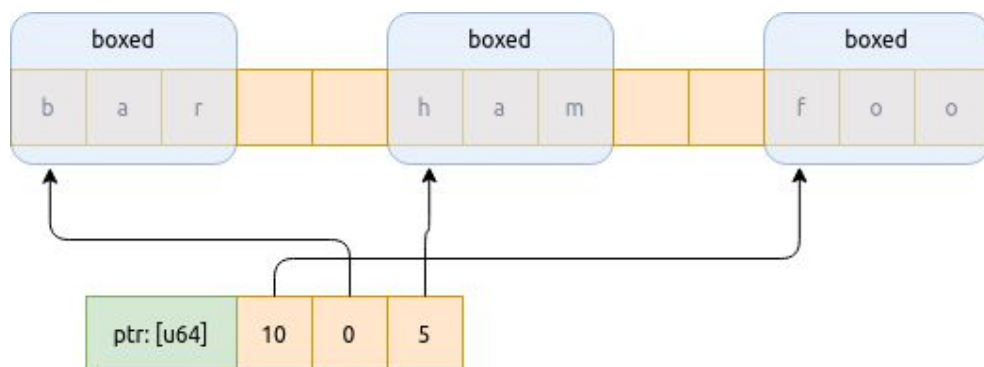
Problem: Numerical compute != Relational compute

- no missing data -> ambiguous use of NaN
- no string data -> boxed python strings -> ptr chasing
- eager evaluation -> no holistic query optimizations
- huge intermediate materializations



# Example: string data

## python strings



## arrow strings

data: [str]	f	o	o	b	a	r	h	a	m
offsets: [i64]	0	2	5	8					
validity bits	0	1	1	0	1	1	0	1	0



# Current DataFrame implementations

\*Python oriented

PANDAS

DASK/  
MODIN/ ETC.

Problem: Single threaded

- "simply" throw more CPU power to the problem
  - instead of rethinking the root problem
- All problems inherited from pandas
  - bad memory representation
  - high memory usage
  - can you really control memory?



# Introducing Polars



DataFrame backend/ query engine written in Rust

- Front end: {Python, Rust, NodeJS}
- Abstraction over arrow memory\*
- Vectorized parallel query engine
- Query planning/optimizations
- Powerful expression API\*
- COW++ semantics\*
- Fast native IO reader/writers

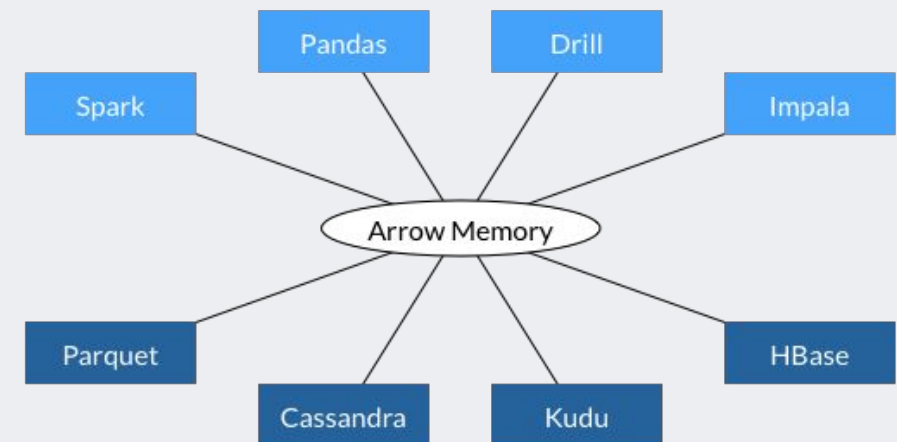




# Foundations: Arrow

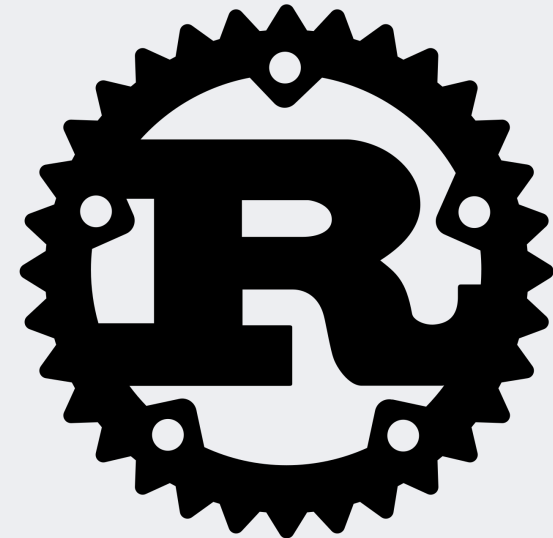
## columnar in-memory standard

- The future of (large) data communication
- << serialization/deserialization cost
- Within same process, free ptr sharing
  - (partial) zero copy interop with:
    - pyarrow
    - ray datasets
    - duckdb
    - dremio
    - ...
- Arrow2: native Rust implementation



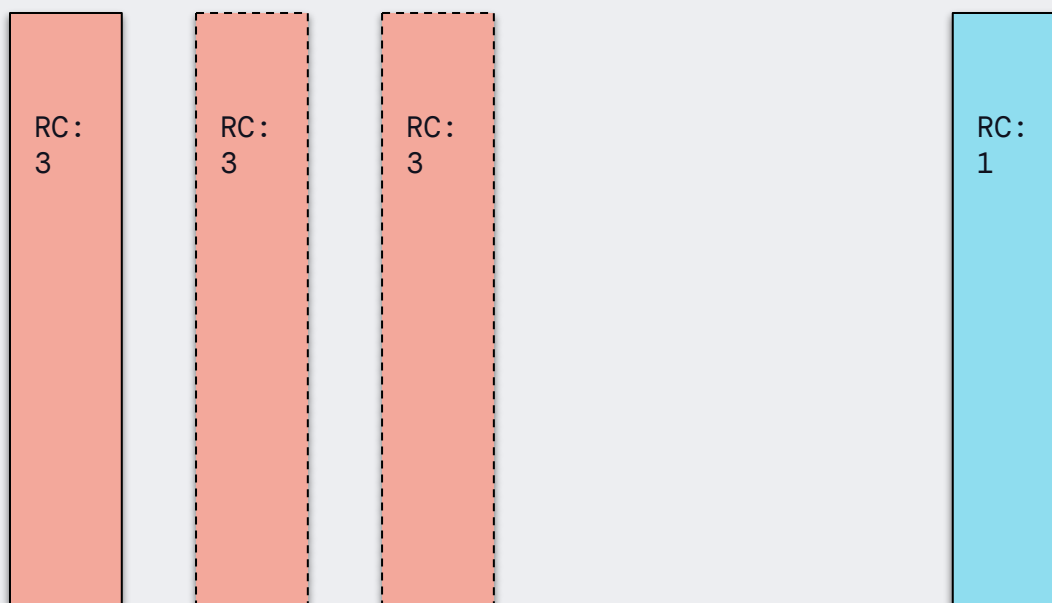
# Foundations: Rust

- Fast C/C++
- Borrow checker/ Ownership
  - memory safety guarantees (without losing control)
  - fearless concurrency



# Example COW++

- Atomically reference counted
- Rust: No mutable aliases



# Polars' performance

# Polars' performance

- Fast cache friendly data-structures and algorithms
- Work stealing parallelism
- In-node parallelism
  - radix groupby
  - partitioned groupby
  - radix join
- SIMD operations
- Query optimizations
- native IO



# H2OAI's db-benchmark

## groupby + join benchmarks



### join summary

#### basic questions

Input table: 100,000,000 rows x 7 columns ( 5 GB )

Library	Version	Date	Time	Status
Polars	0.8.8	2021-06-30	43s	Second time
data.table	1.14.1	2021-06-30	92s	Second time
ClickHouse	21.3.2.5	2021-05-12	159s	Second time
spark	3.1.2	2021-05-31	332s	Second time
DataFrames.jl	1.1.1	2021-06-03	349s	Second time
dplyr	1.0.7	2021-06-20	370s	Second time
(py)datatable	1.0.0a0	2021-06-30	500s	Second time
pandas	1.2.5	2021-06-30	628s	Second time
DuckDB	0.2.7	2021-06-15	630s	Second time
dask	2021.04.1	2021-05-09	internal error	internal error
cuDF*	0.19.2	2021-05-31	internal error	internal error
Arrow	4.0.1	2021-05-31	not yet implemented	not yet implemented
Modin		see README	pending	pending

### groupby summary

#### advanced questions

Input table: 100,000,000 rows x 9 columns ( 5 GB )

Library	Version	Date	Time	Status
Polars	0.8.8	2021-06-30	57s	Second time
ClickHouse	21.3.2.5	2021-05-12	69s	Second time
DataFrames.jl	1.1.1	2021-05-15	116s	Second time
data.table	1.14.1	2021-06-30	120s	Second time
DuckDB	0.2.7	2021-06-15	157s	Second time
(py)datatable	1.0.0a0	2021-06-30	323s	Second time
pandas	1.2.5	2021-06-30	1081s	Second time
Arrow	4.0.1	2021-05-31	4273s	Second time
dplyr	1.0.7	2021-06-20	4378s	Second time
spark	3.1.2	2021-05-31	not yet implemented	not yet implemented
dask	2021.04.1	2021-05-09	internal error	internal error
cuDF*	0.19.2	2021-05-31	out of memory	out of memory
Modin		see README	pending	pending

#### basic questions

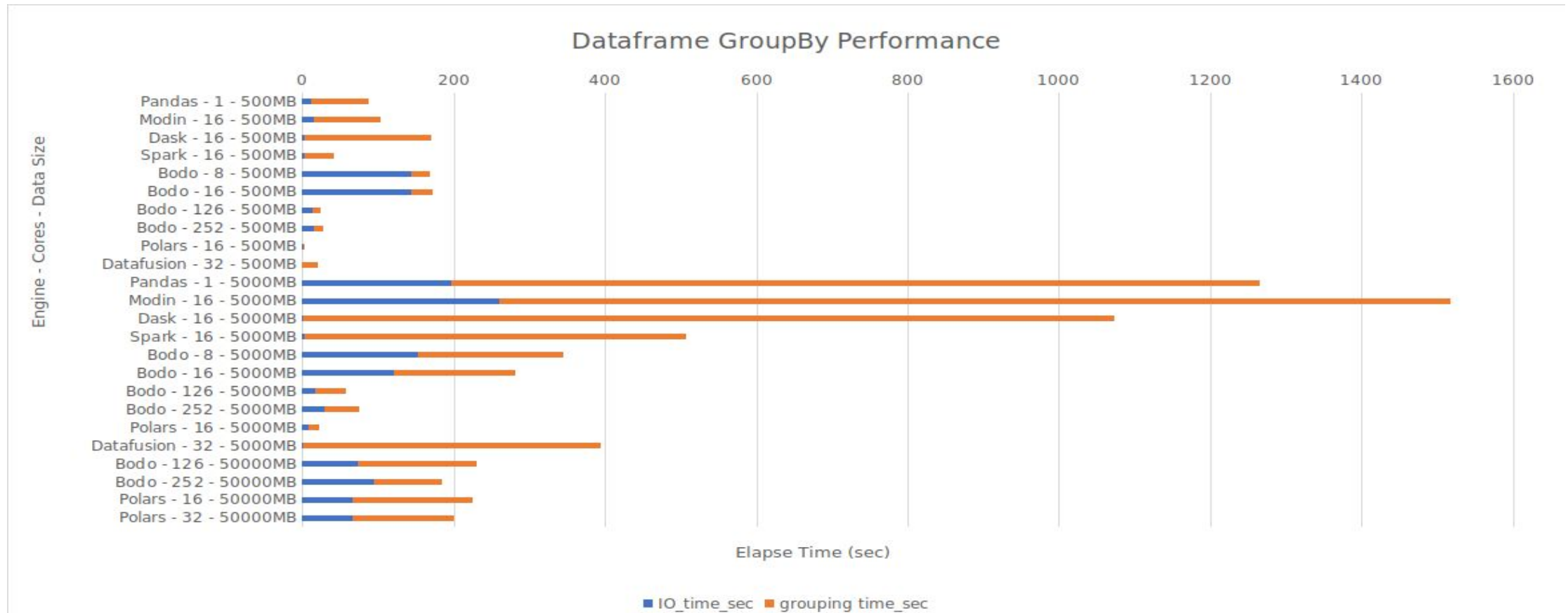
Input table: 1,000,000,000 rows x 9 columns ( 50 GB )

Library	Version	Date	Time	Status
Polars	0.8.8	2021-06-30	143s	Second time
data.table	1.14.1	2021-06-30	155s	Second time
DataFrames.jl	1.1.1	2021-05-15	200s	Second time
ClickHouse	21.3.2.5	2021-05-12	256s	Second time
cuDF*	0.19.2	2021-05-31	492s	Second time
spark	3.1.2	2021-05-31	568s	Second time
(py)datatable	1.0.0a0	2021-06-30	730s	Second time
dplyr	1.0.7	2021-06-20	internal error	internal error
pandas	1.2.5	2021-06-30	out of memory	out of memory
dask	2021.04.1	2021-05-09	out of memory	out of memory
Arrow	4.0.1	2021-05-31	internal error	internal error
DuckDB*	0.2.7	2021-06-15	out of memory	out of memory
Modin		see README	pending	pending



# H2OAI's db-benchmark

only groupby



# Polars' Expressions



# Expressions

- Strive for a small API surface
  - composable blocks
- Vocabulary of a programming language small AND powerful
- control AND performance



# Expressions

```
type Expr = Fn(Series) -> Series
where Series: <column> | <group> | <list el>
// f(g(Series)) = (f ° g)(Series)
// lazy
// optimizable
```



# Expressions

```
// very composable/flexible  
my_expr = pl.col("foo").sort_by("bar") /  
pl.col("ham").where(pl.col("spam") > 19).sum()
```



# Expressions

```
// very composable/flexible  
my_expr = pl.col("foo").sort_by("bar") /  
pl.col("ham").where(pl.col("spam") > 19).sum()
```



# Expressions

Contexts deal with the expression execution

- embarrassingly parallel
- cache and optimize expressions
- in-expression parallelism

```
df.select([
  pl.col("foo").rank(),
  pl.when(pl.col("bar") > 10).then(pl.col("foo").first()).otherwise("ham"),
  pl.col("foo").diff().alias("diff_foo")
])
```



# Example

**DATA+AI**  
SUMMIT 2022

# Thank you

Ritchie Vink

- github: <https://github.com/pola-rs/polars>
- discord: <https://discord.com/invite/4UfP5cfBE7>