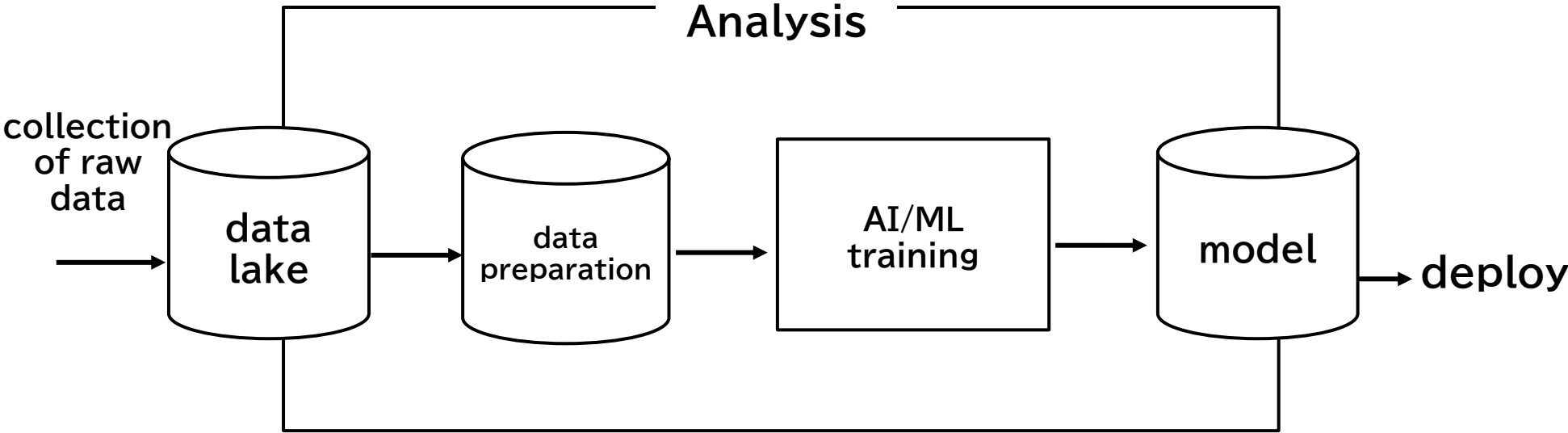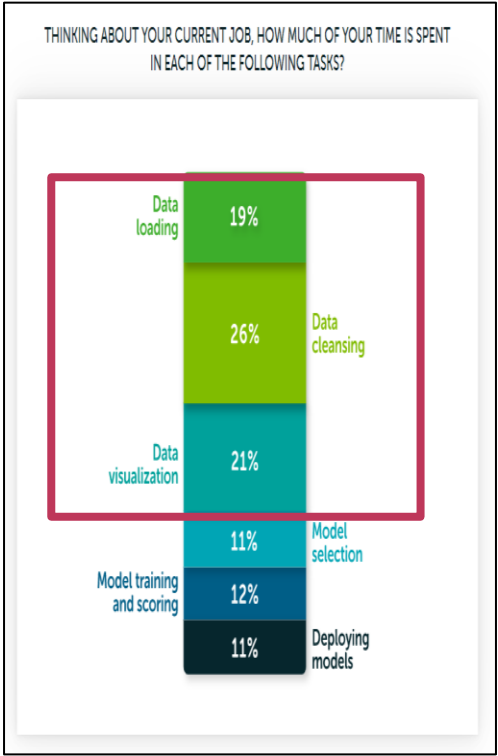# FireDucks: Pandas Accelerator using MLIR

September 28, 2024

Sourav Saha (NEC), Kazuhisa Ishizaka (NEC), Ashu Thakur (NEC)

# Workflow of a Data Scientist
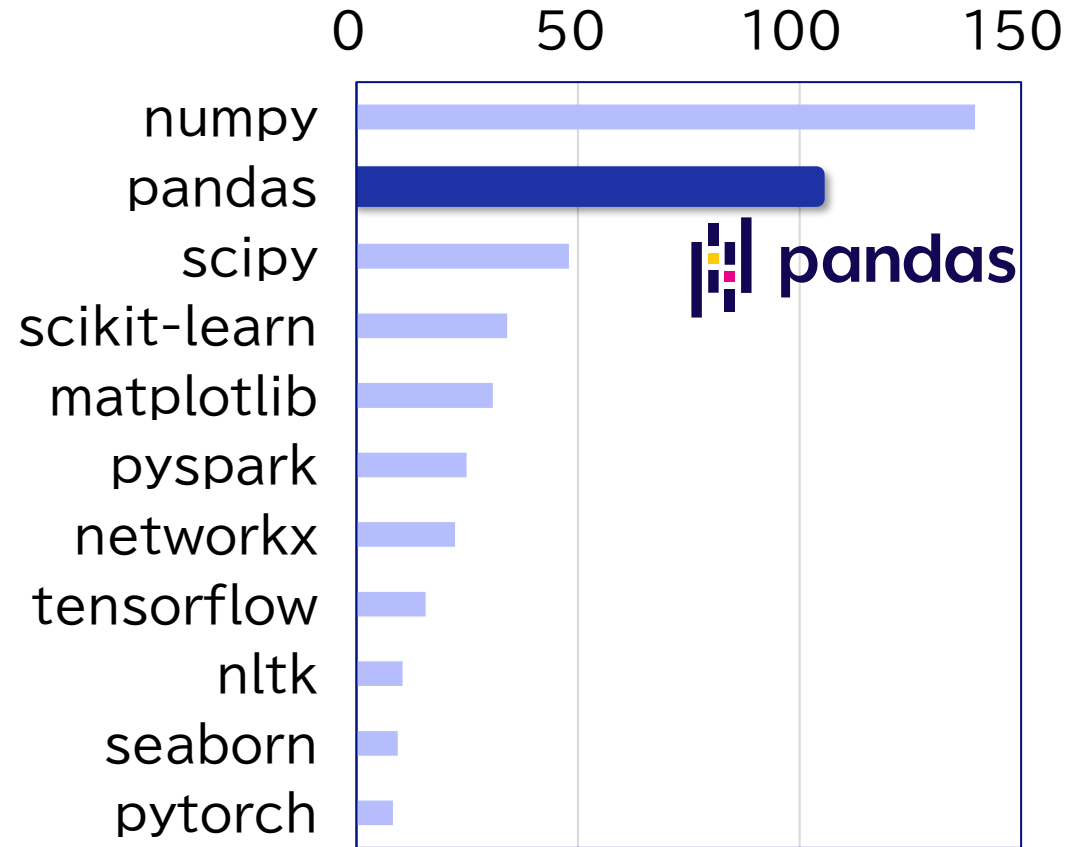
**almost 75% efforts of a Data Scientist spent on data preparation**



THINKING ABOUT YOUR CURRENT JOB, HOW MUCH OF YOUR TIME IS SPENT IN EACH OF THE FOLLOWING TASKS?

| | |
|---|---|
| Data loading | 19% |
| Data cleansing | 26% |
| Data visualization | 21% |
| Model selection | 11% |
| Model training and scoring | 12% |
| Deploying models | 11% |

Anaconda:
The State of Data Science 2020



Analysis

collection of raw data → data lake → data preparation → AI/ML training → model → deploy

\Orchestrating a brighter world    NEC

# Pandas: Its Pros and Cons

◆ **Most popular Python library for data analytics.**

Monthly download from pypi.org
(Data Analytics Libraries)

■ **pandas drawbacks:**

- It (mostly) doesn't support parallel computation.
- The choice of API heavily impacts the performance of a pandas application.
- Very slow execution reduces the efficiency of a data analyst.
- Long-running execution
  - produces higher cloud costs
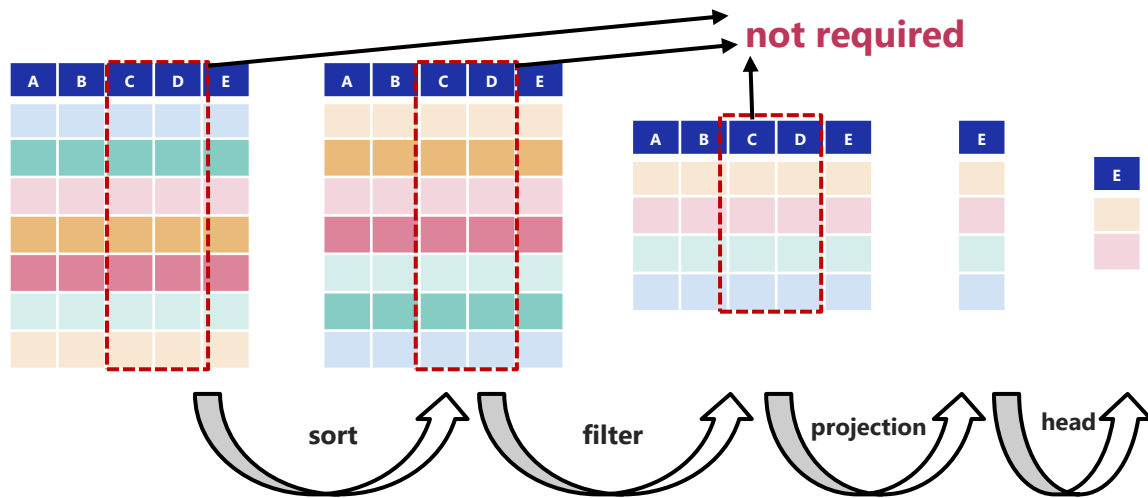  - attributes to higher $CO_2$ emission

The way of implementing a query in pandas-like library (that does not support query optimization) heavily impacts its performance!!

\Orchestrating a brighter world **NEC**

# Execution order matters to boost the performance of a data analysis tool

```
df.sort_values("A")
    .query("B > 1")["E"]
    .head(2)
```

※ sort-order: yellow->red->green->blue



**not required**

sort    filter    projection    head

**SAMPLE QUERY**

```
df.loc[:, ["A", "B", "E"]]
    .query("B > 1")
    .sort_values("A")["E"]
    .head(2)
```

projection    filter    sort    projection    head

reduction in the number of columns

reduction in the number of rows
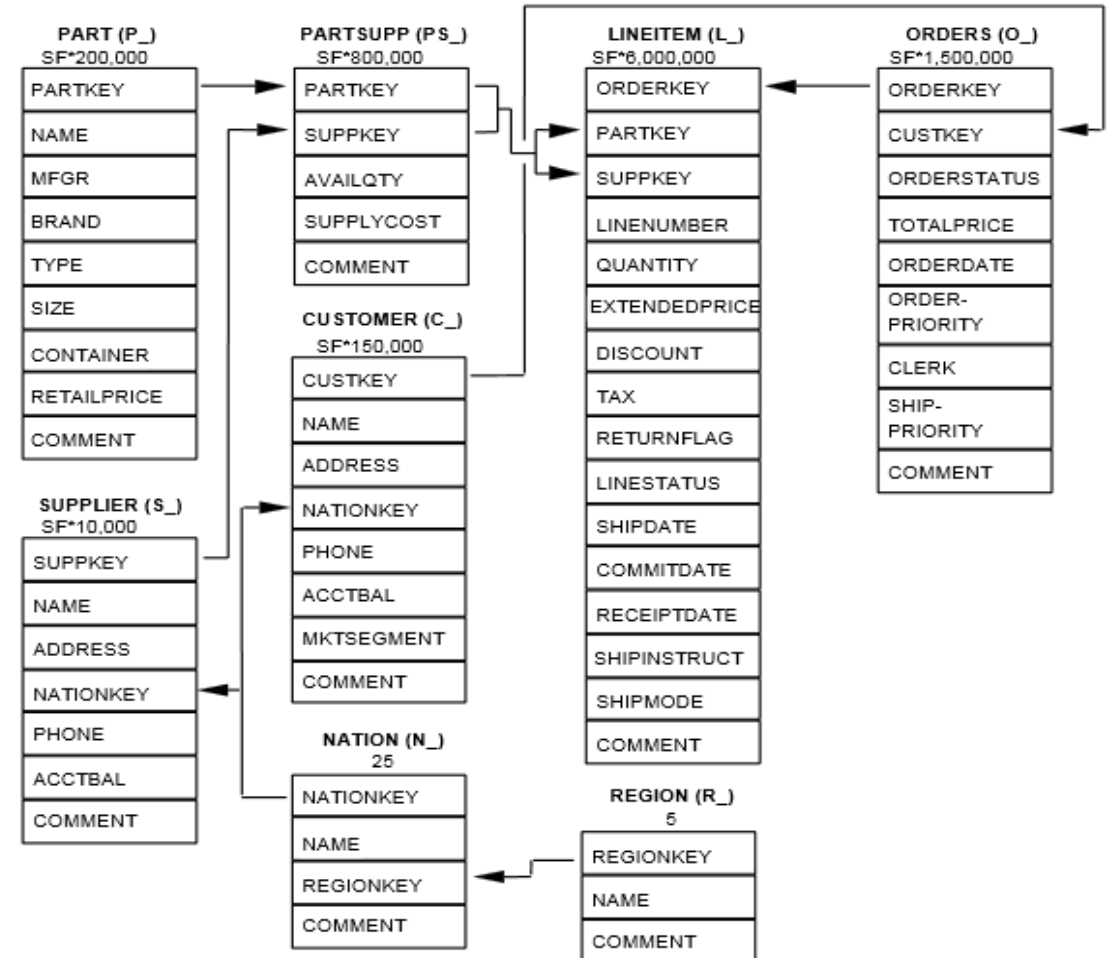
**OPTIMIZED QUERY**

\Orchestrating a brighter world   NEC

# Exercise: Query #3 from TPC-H Benchmark (SQL -> pandas)

◆ query to retrieve the 10 unshipped orders with the highest value.

```sql
SELECT l_orderkey,
               sum(l_extendedprice * (1 - l_discount)) as revenue,
               o_orderdate,
               o_shippriority
FROM customer, orders, lineitem
WHERE
     c_mktsegment = 'BUILDING' AND
     c_custkey = o_custkey AND
     l_orderkey = o_orderkey AND
     o_orderdate < date '1995-03-15' AND
     l_shipdate > date '1995-03-15'
GROUP BY l_orderkey, o_orderdate, o_shippriority
ORDER BY revenue desc, o_orderdate
LIMIT 10;
```

```python
rescols = ["l_orderkey", "revenue", "o_orderdate", "o_shippriority"]
result = (
 customer.merge(orders, left_on="c_custkey", right_on="o_custkey")
   .merge(lineitem, left_on="o_orderkey", right_on="l_orderkey")
   .pipe(lambda df: df[df["c_mktsegment"] == "BUILDING"])
   .pipe(lambda df: df[df["o_orderdate"] < datetime(1995, 3, 15)])
   .pipe(lambda df: df[df["l_shipdate"] > datetime(1995, 3, 15)])
   .assign(revenue=lambda df: df["l_extendedprice"] * (1 - df["l_discount"]))
   .groupby(["l_orderkey", "o_orderdate", "o_shippriority"], as_index=False)
   .agg({"revenue": "sum"})[rescols]
   .sort_values(["revenue", "o_orderdate"], ascending=[False, True])
   .head(10)
)
```

Orchestrating a brighter world    NEC

# Exercise: Query #3 from TPC-H Benchmark (pandas -> optimized pandas)

```python
rescols = ["l_orderkey", "revenue", "o_orderdate", "o_shippriority"]
result = (
 customer.merge(orders, left_on="c_custkey", right_on="o_custkey")
  .merge(lineitem, left_on="o_orderkey", right_on="l_orderkey")
  .pipe(lambda df: df[df["c_mktsegment"] == "BUILDING"])
  .pipe(lambda df: df[df["o_orderdate"] < datetime(1995, 3, 15)])
  .pipe(lambda df: df[df["l_shipdate"] > datetime(1995, 3, 15)])
  .assign(revenue=lambda df: df["l_extendedprice"] * (1 - df["l_discount"]))
  .groupby(["l_orderkey", "o_orderdate", "o_shippriority"], as_index=False)
  .agg({"revenue": "sum"})[rescols]
  .sort_values(["revenue", "o_orderdate"], ascending=[False, True])
  .head(10)
)
```

**Exec-time: 68.55 s**

**Scale Factor: 10**

**6.5x**

**Exec-time: 10.33 s**

```python
# projection-filter: to reduce scope of "customer" table to be processed
cust = customer[["c_custkey", "c_mktsegment"]]
f_cust = cust[cust["c_mktsegment"] == "BUILDING"]

# projection-filter: to reduce scope of "orders" table to be processed
ord = orders[["o_custkey", "o_orderkey", "o_orderdate", "o_shippriority"]]
f_ord = ord[ord["o_orderdate"] < datetime(1995, 3, 15)]

# projection-filter: to reduce scope of "lineitem" table to be processed
litem = lineitem[["l_orderkey", "l_shipdate", "l_extendedprice", "l_discount"]]
f_litem = litem[litem["l_shipdate"] > datetime(1995, 3, 15)]

rescols = ["l_orderkey", "revenue", "o_orderdate", "o_shippriority"]
result = ( f_cust.merge(f_ord, left_on="c_custkey", right_on="o_custkey")
  .merge(f_litem, left_on="o_orderkey", right_on="l_orderkey")
  .assign(revenue=lambda df: df["l_extendedprice"] * (1 - df["l_discount"]))
  .pipe(lambda df: df[rescols])
  .groupby(["l_orderkey", "o_orderdate", "o_shippriority"], as_index=False)
  .agg({"revenue": "sum"})[rescols]
  .sort_values(["revenue", "o_orderdate"], ascending=[False, True])
  .head(10)
)
```

\Orchestrating a brighter world    NEC

# Idea #1

- **Can such optimization be automated?**
  - Yes, using LLVM/MLIR define-by-run mechanism we can build specialized intermediate representation for each pandas API.
  - The generated IRs can be parsed to implement different domain-specific optimizations, such as projection pushdown, predicate pushdown, etc.
  - the optimized IRs can be translated back to the pandas API.

```
df.sort_values("A")
   .query("B > 1")["E"]
   .head(2)
```

**pandas API to intermediate representation (IR)**

**MLIR**

**JIT optimization**

```
%v2 = "sort_values_op"(%v1, "A")
%v3 = "filter_op"(%v2, "B > 1")
%v4 = "project_op"(%v3, ["E"])
%v5 = "slice_op"(%v4, 2)
```

**input IR**

**optimized IR**

```
df.loc[:, ["A", "B", "E"]]
   .query("B > 1")
   .sort_values("A")["E"]
   .head(2)
```

**intermediate representation (IR) to pandas API**

```
%t1 = "project_op"(%v1, ["A", "B", "E"])
%t2 = "filter_op"(%t1, "B > 1")
%t3 = "sort_values_op"(%t2, "A")
%t4 = "project_op"(%t3, ["E"])
%t5 = "slice_op"(%t4, 2)
```

\Orchestrating a brighter world **NEC**

# Idea #2

- Pandas methods are slow due to poor memory utilization and single-core computation.
- But pandas is one of the most popular data manipulation tools.
- **How can we solve the core performance issue in pandas while keeping the same API for users?**
  - Well, we can
    - have a frontend with pandas API that generates IR.
    - develop our own library parallelizing the workload of DataFrame-related methods as a backend.
    - translate the optimized IRs to the **backend library API** (instead of pandas API).

```
df.sort_values("A")
   .query("B > 1")["E"]
   .head(2)
```
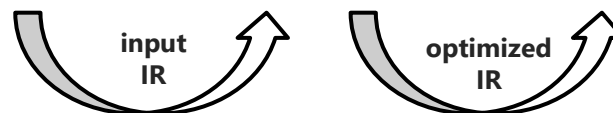
```
t1 = backend::project_columns(df, {"A", "B", "C"});
t2 = backend::filter_rows(t1, "B > 1");
t3 = backend::sort_values(t2, "A");
t4 = backend::project_columns(t3, {"E"});
t5 = backend::slice_rows(t4, 2);
```

**frontend pandas API to
intermediate representation (IR)**

**MLIR**

**intermediate representation (IR)
to backend API**

```
%v2 = "sort_values_op"(%v1, "A")
%v3 = "filter_op"(%v2, "B > 1")
%v4 = "project_op"(%v3, ["E"])
%v5 = "slice_op"(%v4, 2)
```

```
%t1 = "project_op"(%v1, ["A", "B", "E"])
%t2 = "filter_op"(%t1, "B > 1")
%t3 = "sort_values_op"(%t2, "A")
%t4 = "project_op"(%t3, ["E"])
%t5 = "slice_op"(%t4, 2)
```

**JIT optimization**

**input
IR**

**optimized
IR**

\Orchestrating a brighter world   **NEC**

# Introducing FireDucks

**FireDucks** (**F**lexible **IR E**ngine for DataFrame) is a high-performance compiler-accelerated DataFrame library with highly compatible pandas APIs.
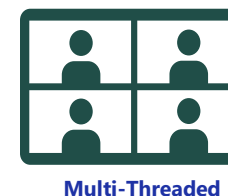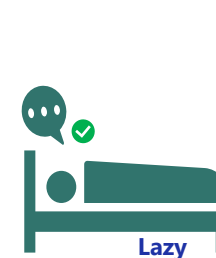
## Speed: significantly faster than pandas

- FireDucks is multithreaded to fully exploit the modern processor
- Lazy execution model with Just-In-Time optimization using a defined-by-run mechanism supported by MLIR (a subproject of LLVM).
  - supports <u>both lazy and non-lazy execution</u> models without modifying user programs (same API).

**MLIR**

## Ease of use: drop-in replacement of pandas

- FireDucks is highly compatible with pandas API
  - <u>seamless integration is possible</u> not only for an existing pandas program but also for any external libraries (like seaborn, scikit-learn, etc.) that internally use pandas dataframes.
- No extra learning is required
- No code modification is required

Lazy

JIT optimization

Cloud-friendly

Multi-Threaded

No new learning

Eco-friendly

data analysis

lightning-fast

\Orchestrating a brighter world    **NEC**

# How does FireDucks work?

User Program

IR Builder

MLIR

Generated IR-OPs → Optimization Passes

Multi-core Kernel Executor

```
sorted = df.sort_values("b")
result = sorted["a"]
```

↓

```
%v2 = "fireducks.sort_values"(%v1,"b")
%v3 = "fireducks.project"(%v2,["a"])
```

↓ **print (result)**

```
%v11 = "fireducks.project"(%v1,["a","b"])
%v2 = "fireducks.sort_values"(%v11,"b")
%v3 = "fireducks.project"(%v2,["a"])
```

↓

```
tmp = df[["a","b"]]
sorted = tmp.sort_values("b")
result = sorted["a"]
```

**Primary Objective: <u>Write Once, Execute Anywhere</u>**

\Orchestrating a brighter world    **NEC**

# Let's Have a Quick Demo!
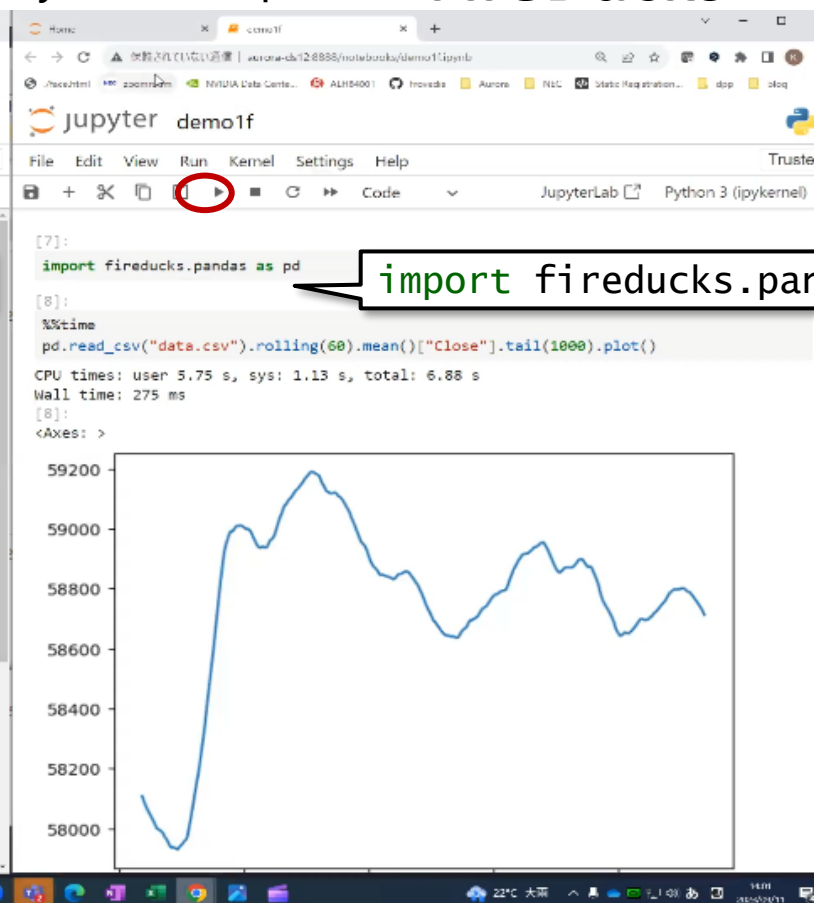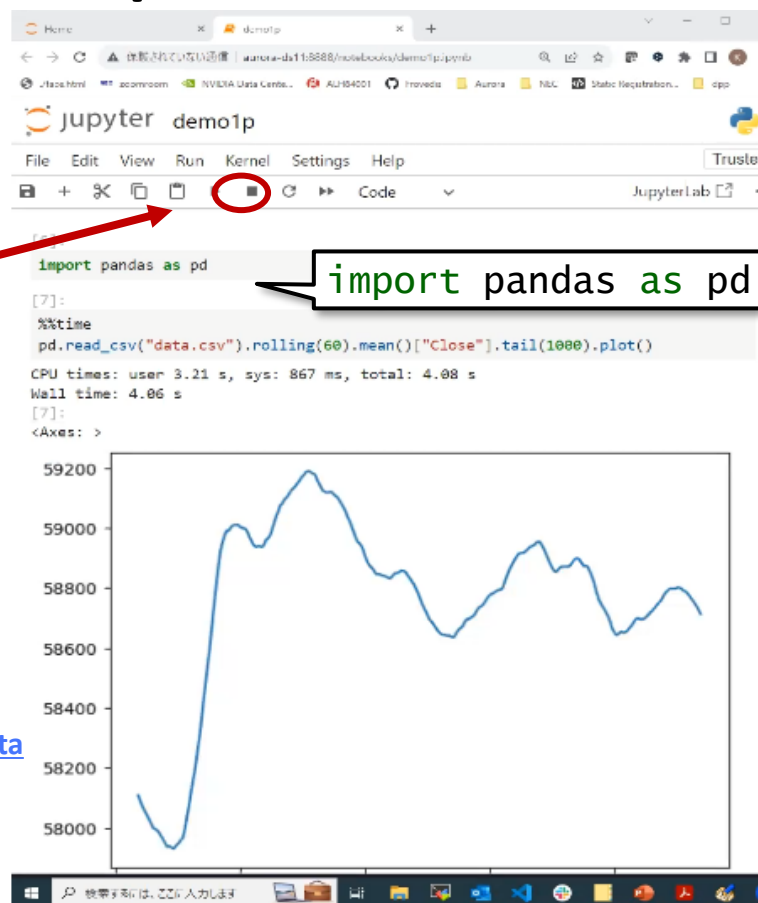
```
pd.read_csv("data.csv").rolling(60).mean()["Close"].tail(1000).plot()
```

**pandas** the difference is only in the import **FireDucks**

Program to calculate moving average



button to start execution

> import pandas as pd

> import fireducks.pandas as pd

data.csv:
**Bitcoin Historical Data**

pandas: 4.06s

**~15x**

FireDucks: 275ms

Orchestrating a brighter world **NEC**

# Usage of FireDucks

## 1. Explicit Import

easy to import

```
# import pandas as pd
import fireducks.pandas as pd
```

simply change the import statement

## 2. Import Hook

FireDucks provides command line option to automatically replace "**pandas**" with "**fireducks.pandas**"

```
$ python -m fireducks.pandas program.py
```

zero code modification

```
import mod_A
import mod_B
import mod_C
import pandas as pd
:
```
program.py

```
import pandas as pd
:
```
mod_A.py

```
import pandas as pd
:
```
mod_B.py

```
import pandas as pd
:
```
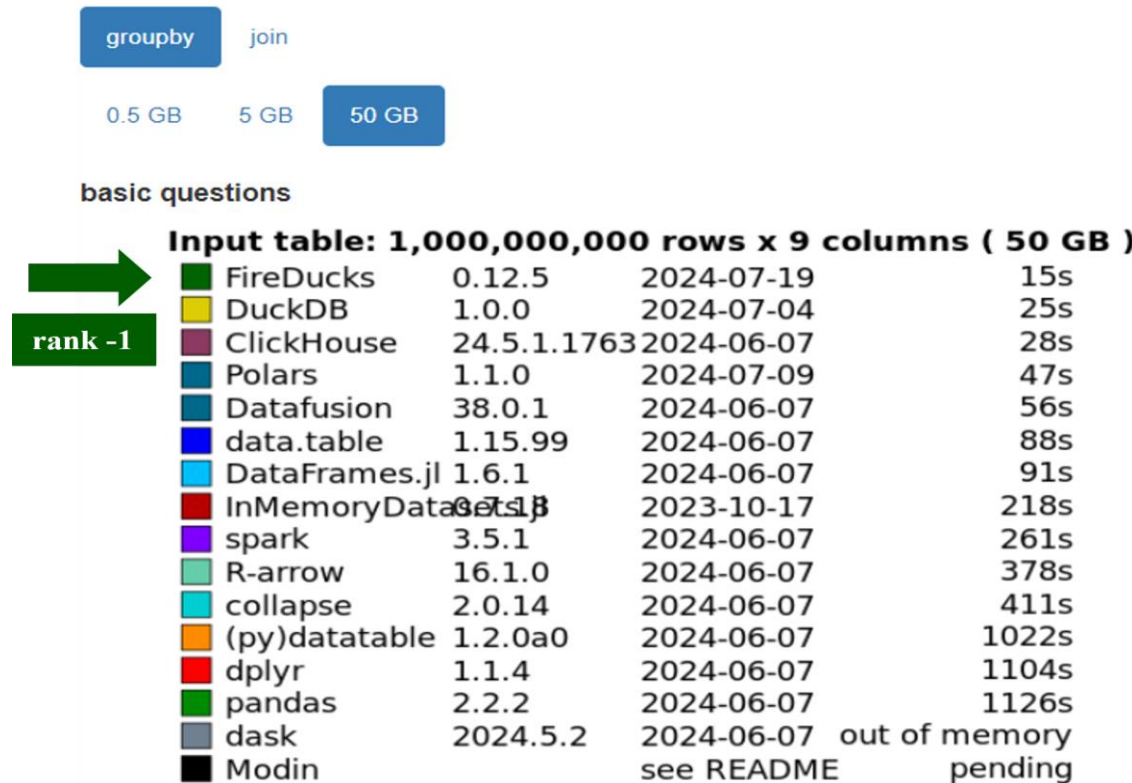mod_C.py

## 3. Notebook Extension

FireDucks provides simple import extension for interative notebooks.
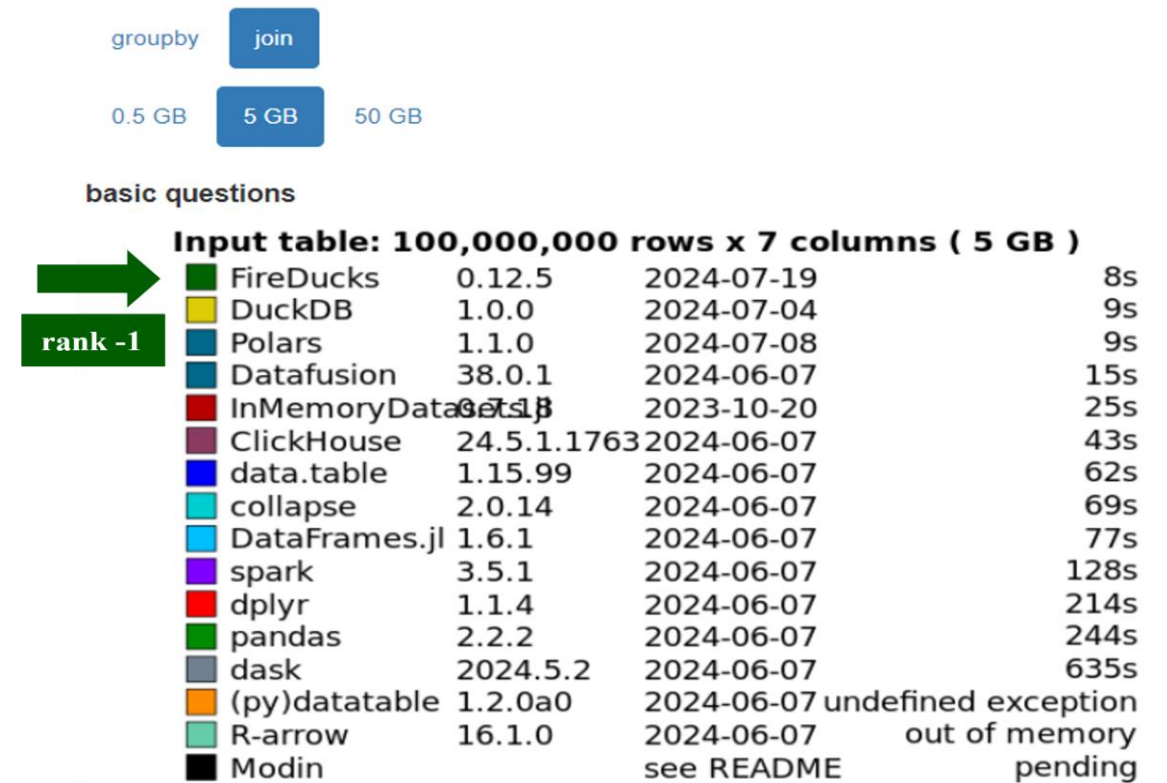
```
%load_ext fireducks.pandas
import pandas as pd
```

simple integration in a notebook

\Orchestrating a brighter world    NEC

# Benchmark (1): DB-Benchmark

Database-like ops benchmark (https://duckdblabs.github.io/db-benchmark)



<table>
<tr><td>groupby</td><td>join</td></tr>
<tr><td>0.5 GB</td><td>5 GB</td><td>50 GB</td></tr>
</table>

**basic questions**

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

| | | | | |
|---|---|---|---|---|
| FireDucks | 0.12.5 | 2024-07-19 | | 15s |
| DuckDB | 1.0.0 | 2024-07-04 | | 25s |
| ClickHouse | 24.5.1.1763 | 2024-06-07 | | 28s |
| Polars | 1.1.0 | 2024-07-09 | | 47s |
| Datafusion | 38.0.1 | 2024-06-07 | | 56s |
| data.table | 1.15.99 | 2024-06-07 | | 88s |
| DataFrames.jl | 1.6.1 | 2024-06-07 | | 91s |
| InMemoryDatasets | 0.7.13 | 2023-10-17 | | 218s |
| spark | 3.5.1 | 2024-06-07 | | 261s |
| R-arrow | 16.1.0 | 2024-06-07 | | 378s |
| collapse | 2.0.14 | 2024-06-07 | | 411s |
| (py)datatable | 1.2.0a0 | 2024-06-07 | | 1022s |
| dplyr | 1.1.4 | 2024-06-07 | | 1104s |
| pandas | 2.2.2 | 2024-06-07 | | 1126s |
| dask | 2024.5.2 | 2024-06-07 | out of memory | |
| Modin | | | see README | pending |

**rank -1**

groupby



<table>
<tr><td>groupby</td><td>join</td></tr>
<tr><td>0.5 GB</td><td>5 GB</td><td>50 GB</td></tr>
</table>

**basic questions**

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

| | | | | |
|---|---|---|---|---|
| FireDucks | 0.12.5 | 2024-07-19 | | 8s |
| DuckDB | 1.0.0 | 2024-07-04 | | 9s |
| Polars | 1.1.0 | 2024-07-08 | | 9s |
| Datafusion | 38.0.1 | 2024-06-07 | | 15s |
| InMemoryDatasets | 0.7.13 | 2023-10-20 | | 25s |
| ClickHouse | 24.5.1.1763 | 2024-06-07 | | 43s |
| data.table | 1.15.99 | 2024-06-07 | | 62s |
| collapse | 2.0.14 | 2024-06-07 | | 69s |
| DataFrames.jl | 1.6.1 | 2024-06-07 | | 77s |
| spark | 3.5.1 | 2024-06-07 | | 128s |
| dplyr | 1.1.4 | 2024-06-07 | | 214s |
| pandas | 2.2.2 | 2024-06-07 | | 244s |
| dask | 2024.5.2 | 2024-06-07 | | 635s |
| (py)datatable | 1.2.0a0 | 2024-06-07 | undefined exception | |
| R-arrow | 16.1.0 | 2024-06-07 | out of memory | |
| Modin | | | see README | pending |

**rank -1**

join

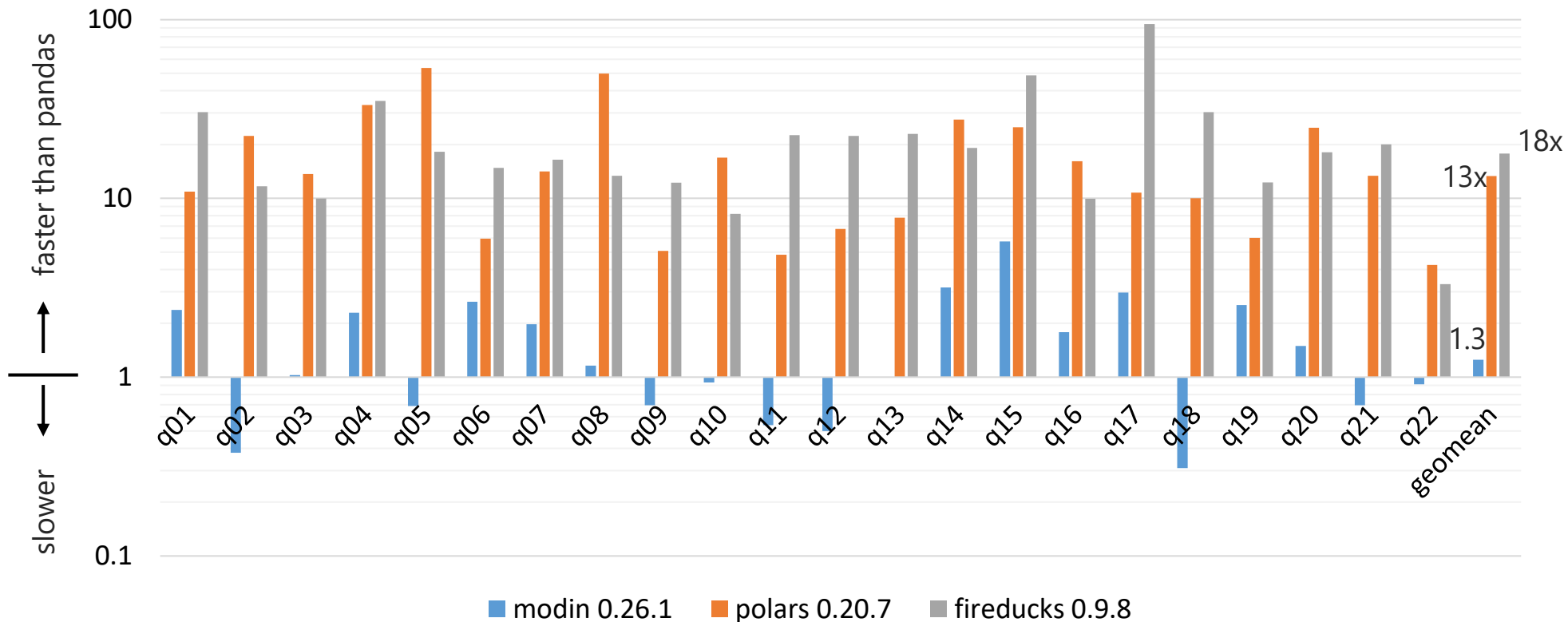\Orchestrating a brighter world   **NEC**

# Benchmark: Speedup from pandas in TPC-H benchmark

## FireDucks is 95x faster than pandas at max

Server

Xeon Gold 5317 x2
(24 cores), 256GB

### Speedup from pandas 2.2.0 (Scale Factor=10)
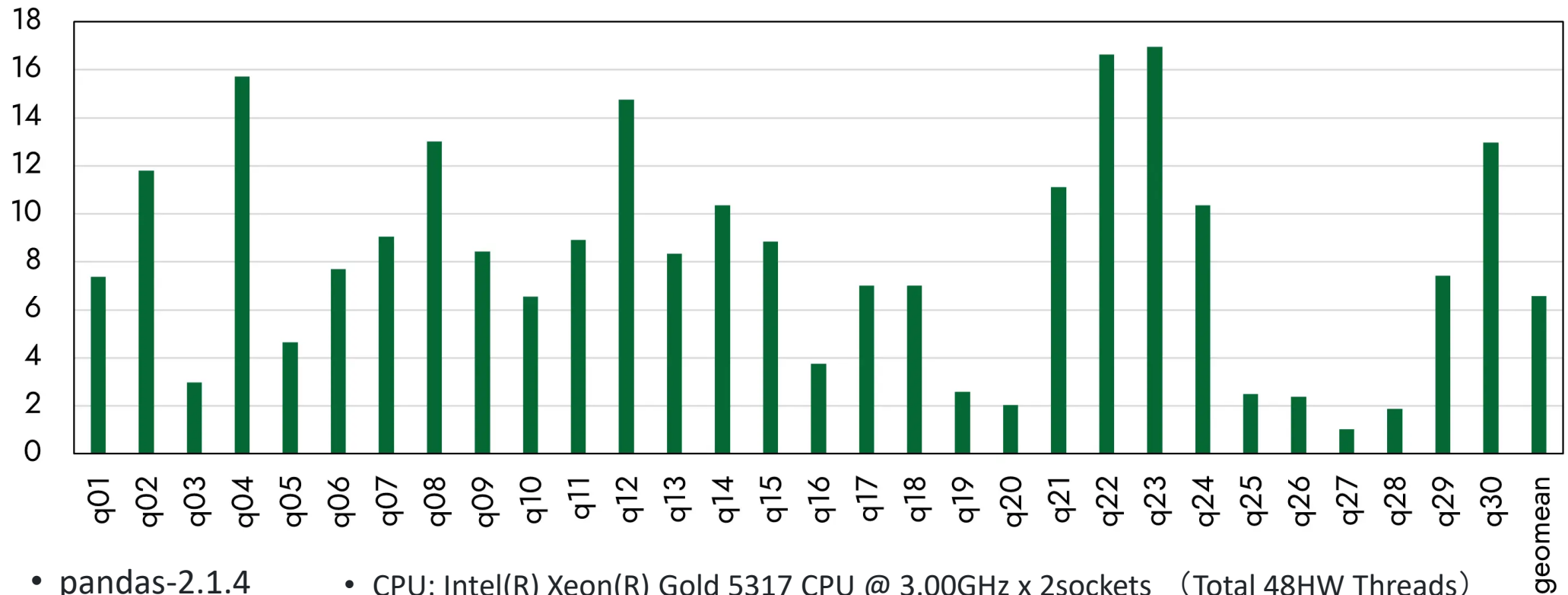


Comparison of
DataFrame libraries
(average speedup)

| | |
|---|---|
| **FireDucks** | **18x** |
| Polars | 13x |
| Modin | 1.3x |

■ modin 0.26.1   ■ polars 0.20.7   ■ fireducks 0.9.8

\Orchestrating a brighter world   NEC

# Benchmark (3): Speedup from pandas in TPCx-BB benchmark

## ETL(Extract, Transform, Load) and ML Workflow

FireDucks speedup from pandas



- pandas-2.1.4
- fireducks-0.9.3
- CPU: Intel(R) Xeon(R) Gold 5317 CPU @ 3.00GHz x 2sockets （Total 48HW Threads）
- Main memory: 256GB

\Orchestrating a brighter world    NEC

# Resource on FireDucks

**Web site (User guide, benchmark, blog)**

**https://fireducks-dev.github.io/**

**X(twitter) (Release information)**

**https://x.com/fireducksdev**

**Github (Issue report)**

**https://github.com/fireducks-dev/fireducks**

**Q/A, communication**

**https://join.slack.com/t/fireDucks/shared_invite/zt-2j4lucmtj-IGR7AWlXO62Lu605pnBJ2w**

## FireDucks

Compiler Accelerated DataFrame Library for Python with fully-compatible pandas API

Get Started

```
import fireducks.pandas as pd
```

News
Release fileducks-0.12.4 (Jul 09, 2024)
Have you ever thought of speeding up your data analysis in pandas with a compiler?(blog) (Jul 03, 2024)
Evaluation result of Database-like ops benchmark with FireDucks is now available. (Jun 18, 2024)

### Accelerate pandas without any manual code changes

Do you have a pandas-based program that is slow? FireDucks can speed-up your programs without any manual code changes. You can accelerate your data analysis without worrying about slow performance due to single-threaded execution in pandas.

\Orchestrating a brighter world    **NEC**

# Let's go for a test drive!

**https://colab.research.google.com/drive/1qpej-X7CZsIeOqKuhBg4kq-cbGuJf1Zp?usp=sharing**

\Orchestrating a brighter world   NEC

# Thank You!

◆ Focus more on in-depth data exploration using "**pandas**".

◆ Let the "**FireDucks**" take care of the optimization for you.

◆ Enjoy Green Computing!

\Orchestrating a brighter world NEC

# \Orchestrating a brighter world

NEC creates the social values of safety, security,

fairness and efficiency to promote a more sustainable world

where everyone has the chance to reach their full potential.

\Orchestrating a brighter world

NEC