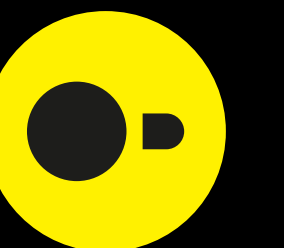


Why would you "import duckdb" in your Python project?

GÁBOR SZÁRNYAS
(DUCKDB LABS)





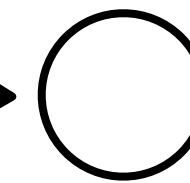
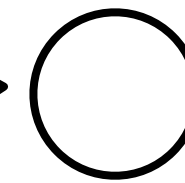
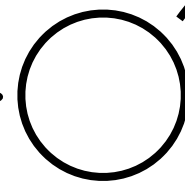
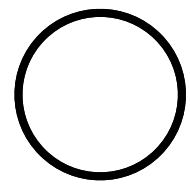
MSc

PhD

CWI



DuckDB Labs

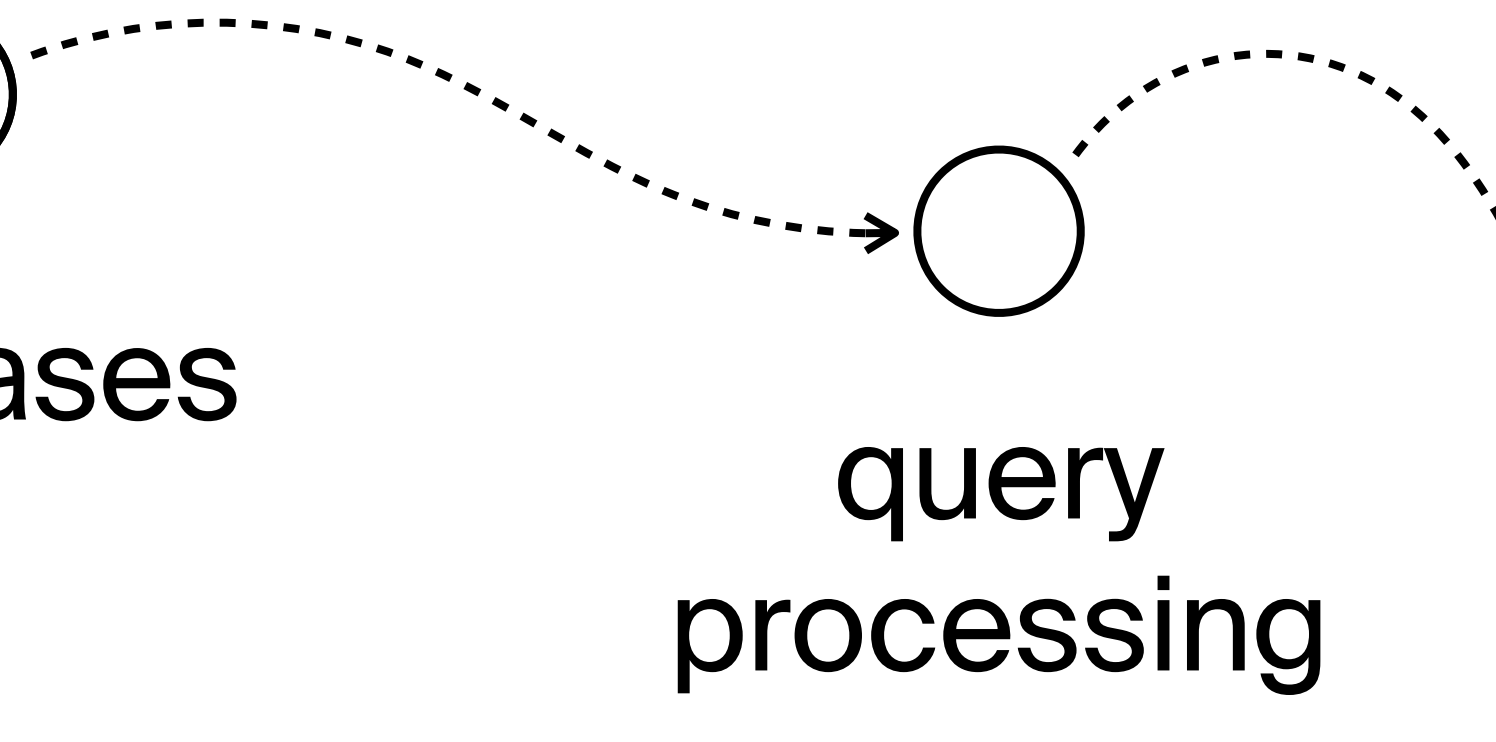


databases

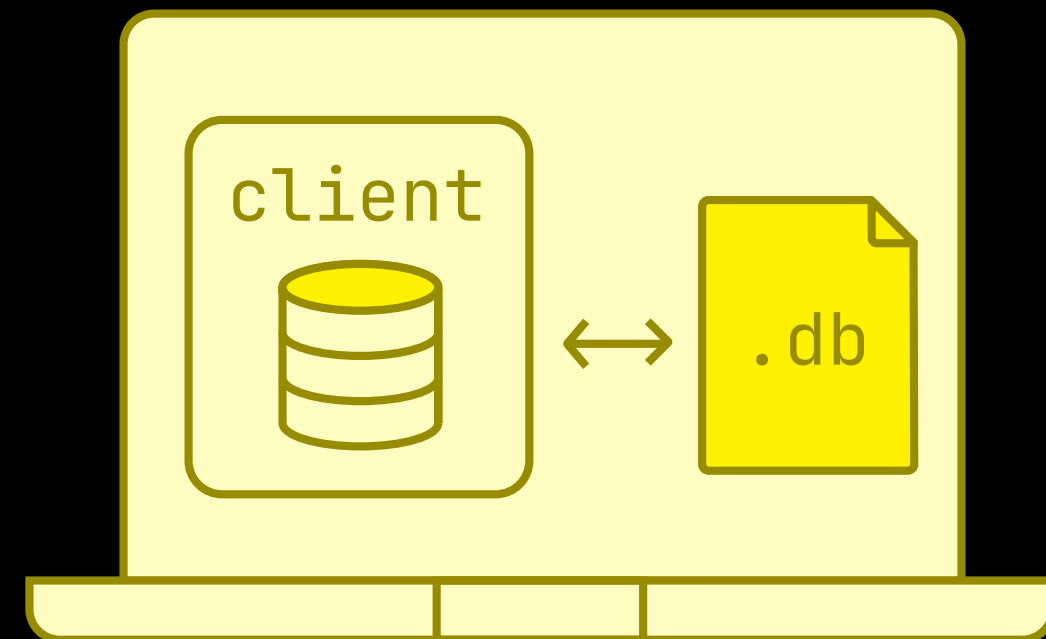
query
processing

postdoc

graph databases,
benchmarks



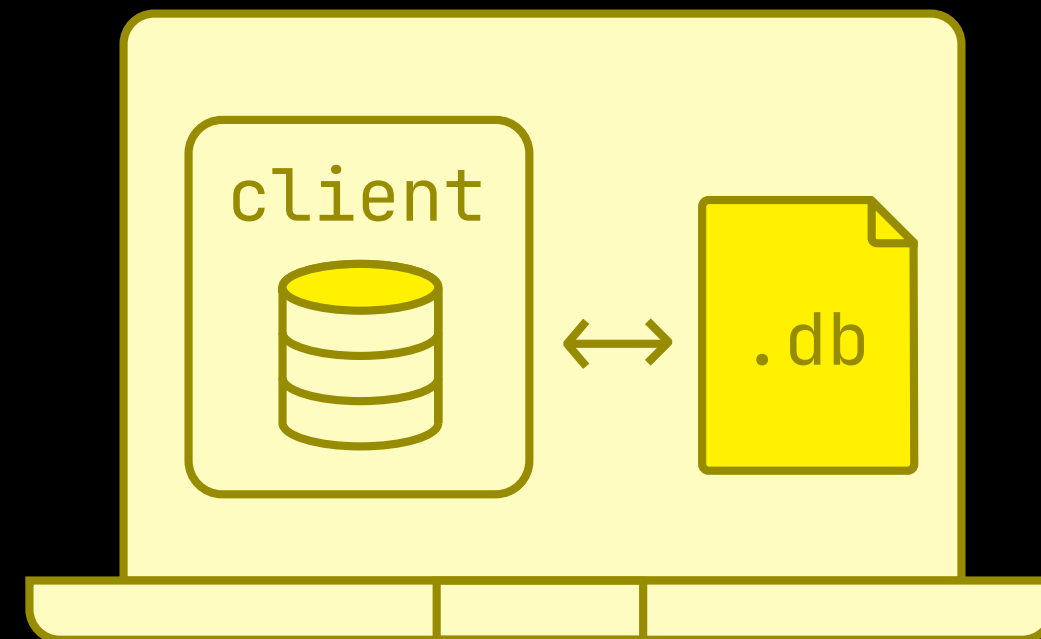
DuckDB = a SQL database for analytics



```
$ curl https://install.duckdb.org | sh  
$ duckdb my.db
```

```
D SELECT *  
FROM '/your/favorite/file.csv';
```

DuckDB = a SQL database for analytics



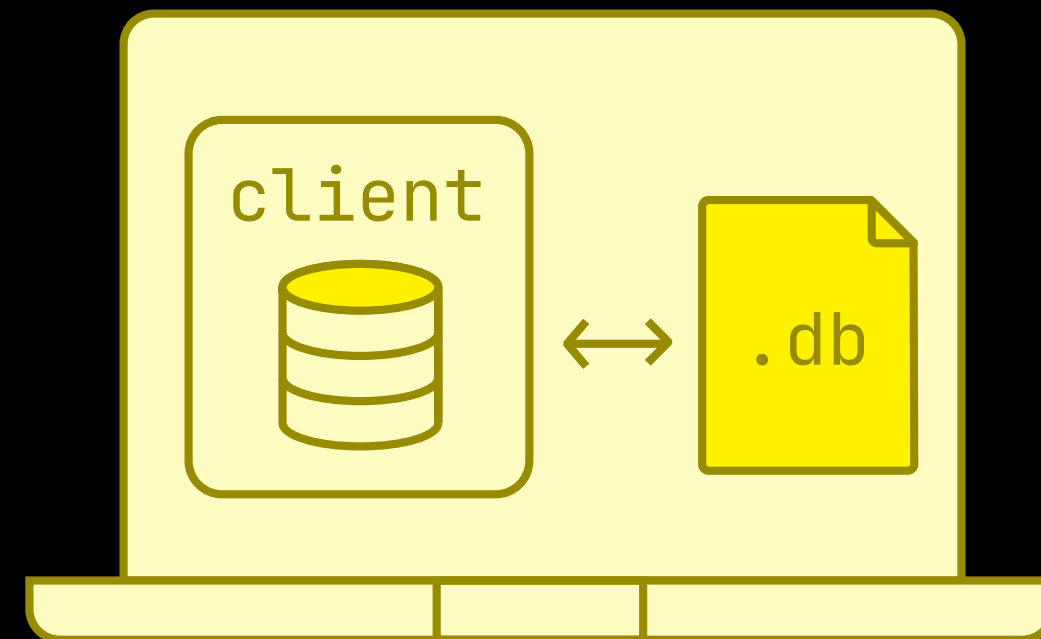
```
$ curl https://install.duckdb.org | sh  
$ duckdb my.db
```

```
D SELECT *  
FROM '/your/favorite/file.csv';
```

```
$ pip install duckdb  
$ python
```

```
>>> import duckdb  
>>> duckdb.sql("FROM 'my_data.csv'")
```

DuckDB = a SQL database for analytics



```
$ curl https://install.duckdb.org | sh  
$ duckdb my.db
```

```
D SELECT *  
FROM '/your/favorite/file.csv';
```

```
$ pip install duckdb  
$ python
```

```
>>> import duckdb  
>>> duckdb.sql("FROM 'my_data.csv'")
```

Zero-dependency C++ code: 15+ clients

PyPI Stats

[Search](#)

[All packages](#)

[Top packages](#)

[Track packages](#)

duckdb

[PyPI page](#)

[Home page](#)

Author: DuckDB Foundation

Summary: DuckDB in-process database

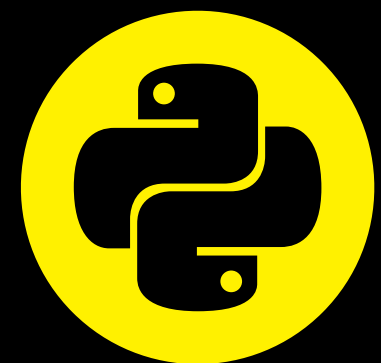
Latest version: 1.5.2

Optional dependencies: [adbc-driver-manager](#) | [fsspec](#) | [ipython](#) | [numpy](#) | [pandas](#) | [pyarrow](#)

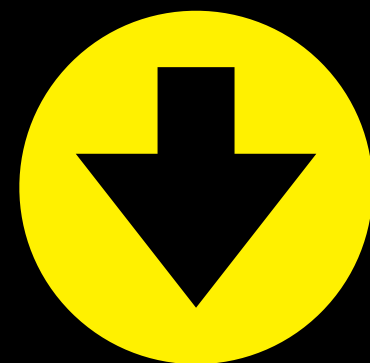
Downloads last day: 762,722

Downloads last week: 10,178,060

Downloads last month: 39,030,034



Top 1,000 in PyPI



0.5B Py downloads



38k GitHub stars



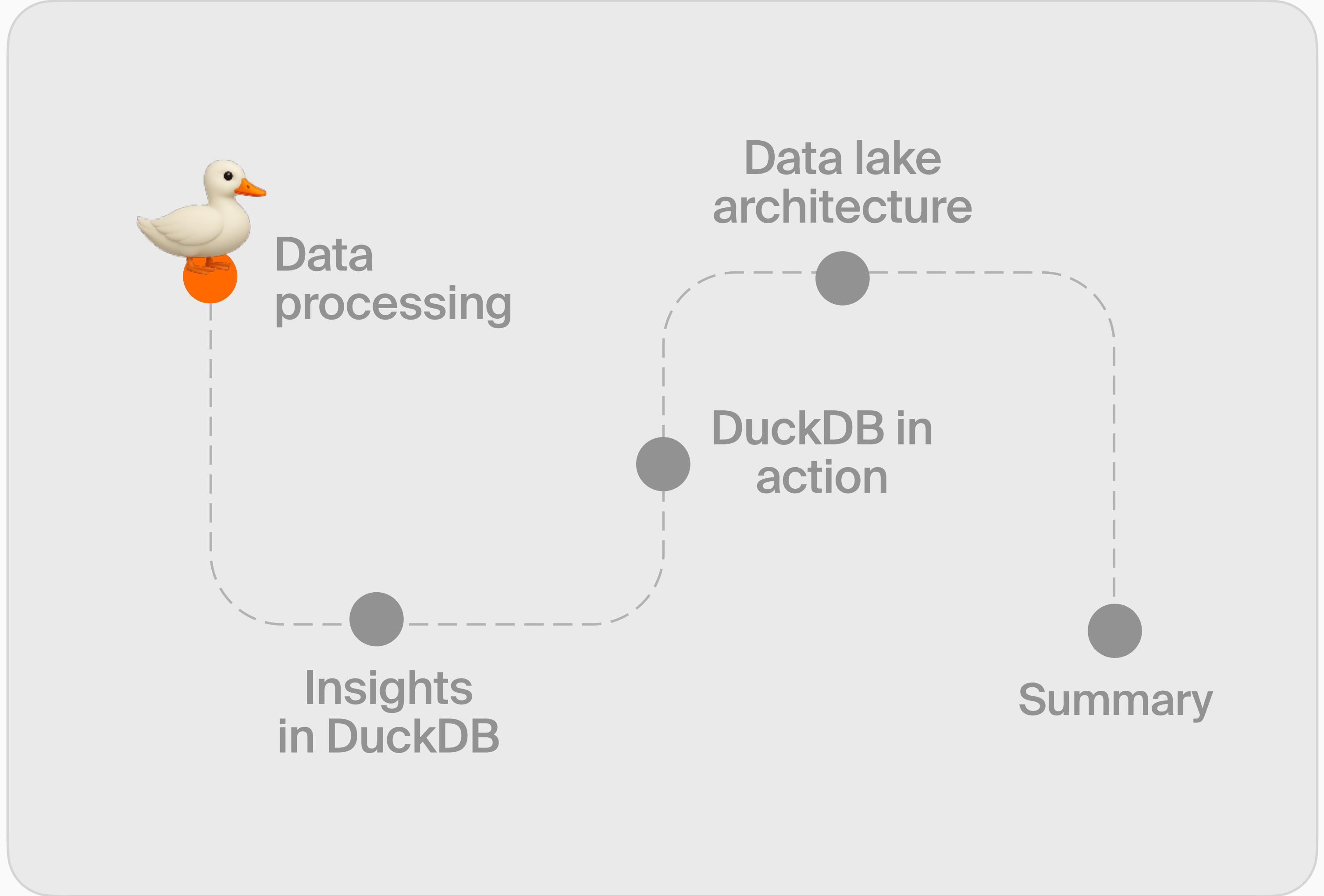
Data processing

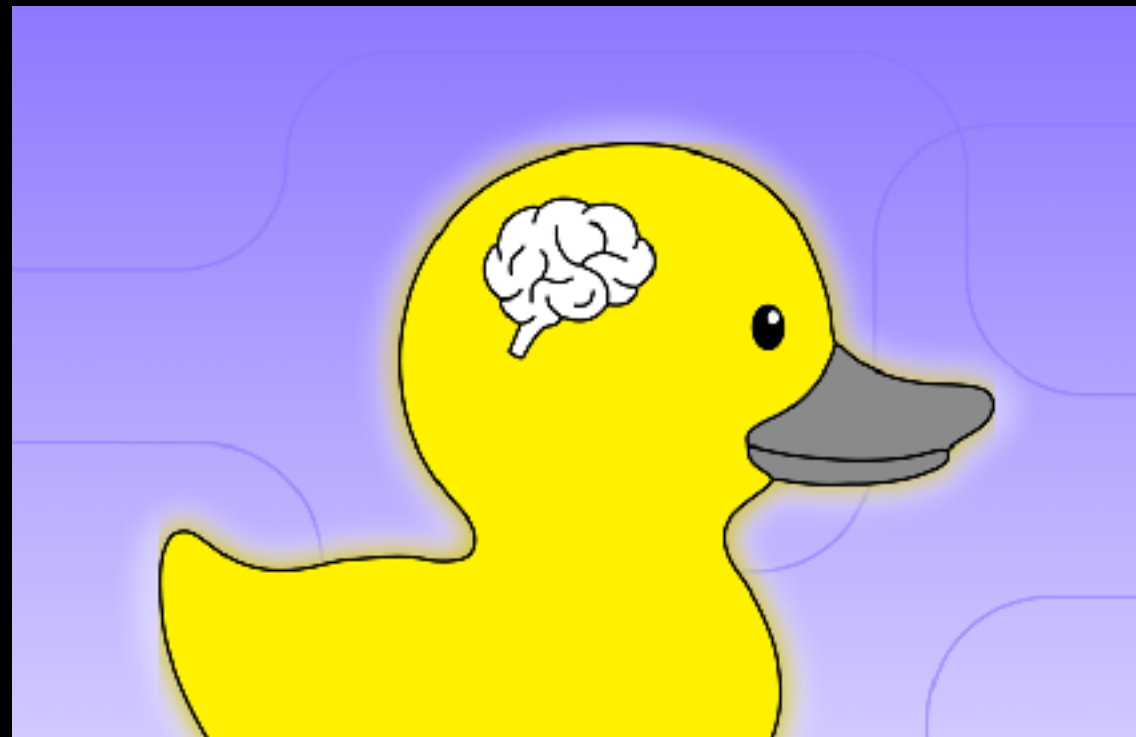
Data lake architecture

DuckDB in action

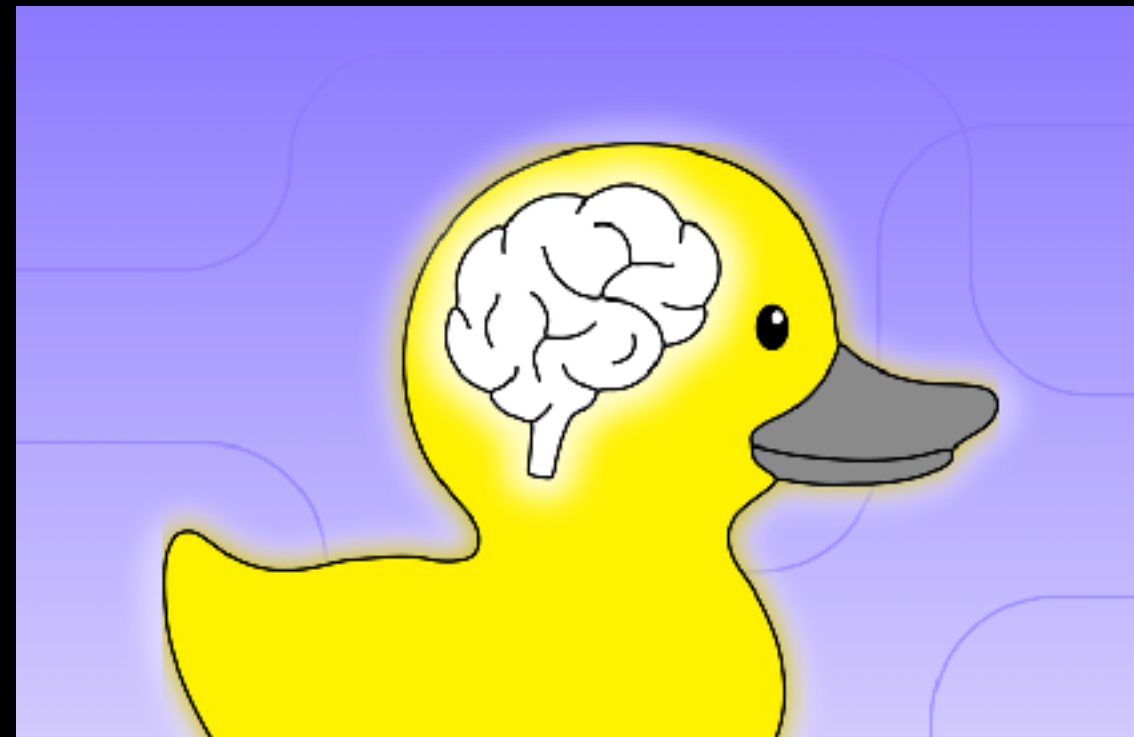
Insights in DuckDB

Summary

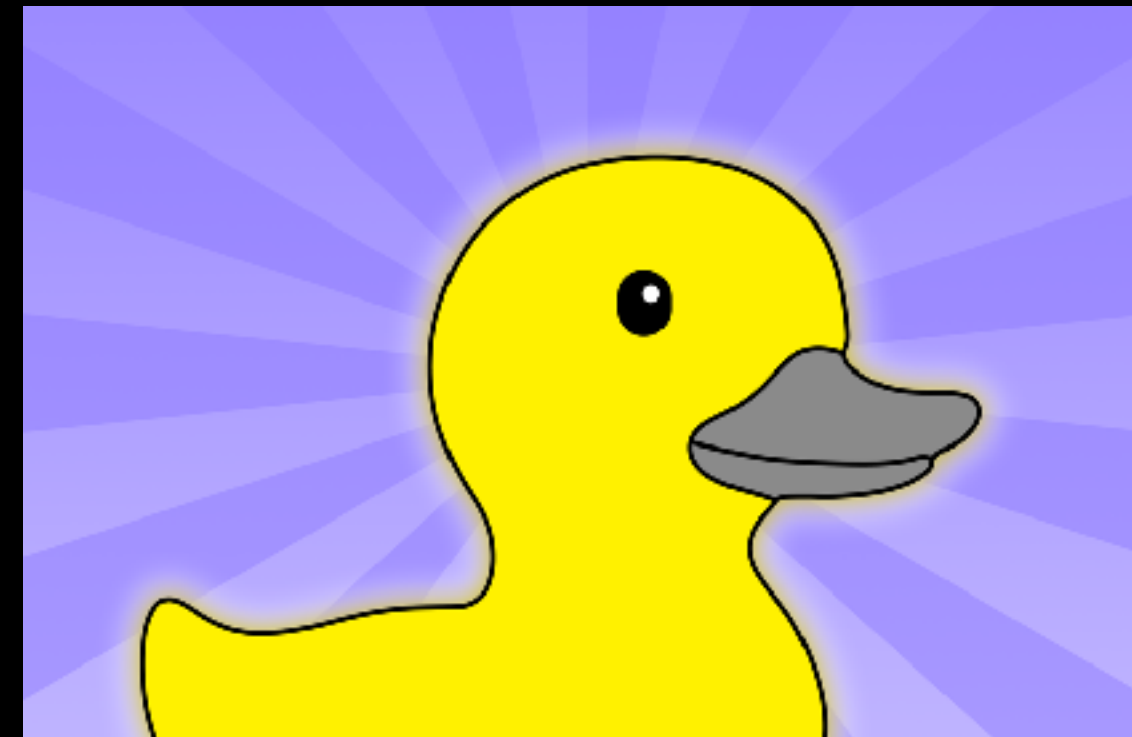




Text files



Binary files

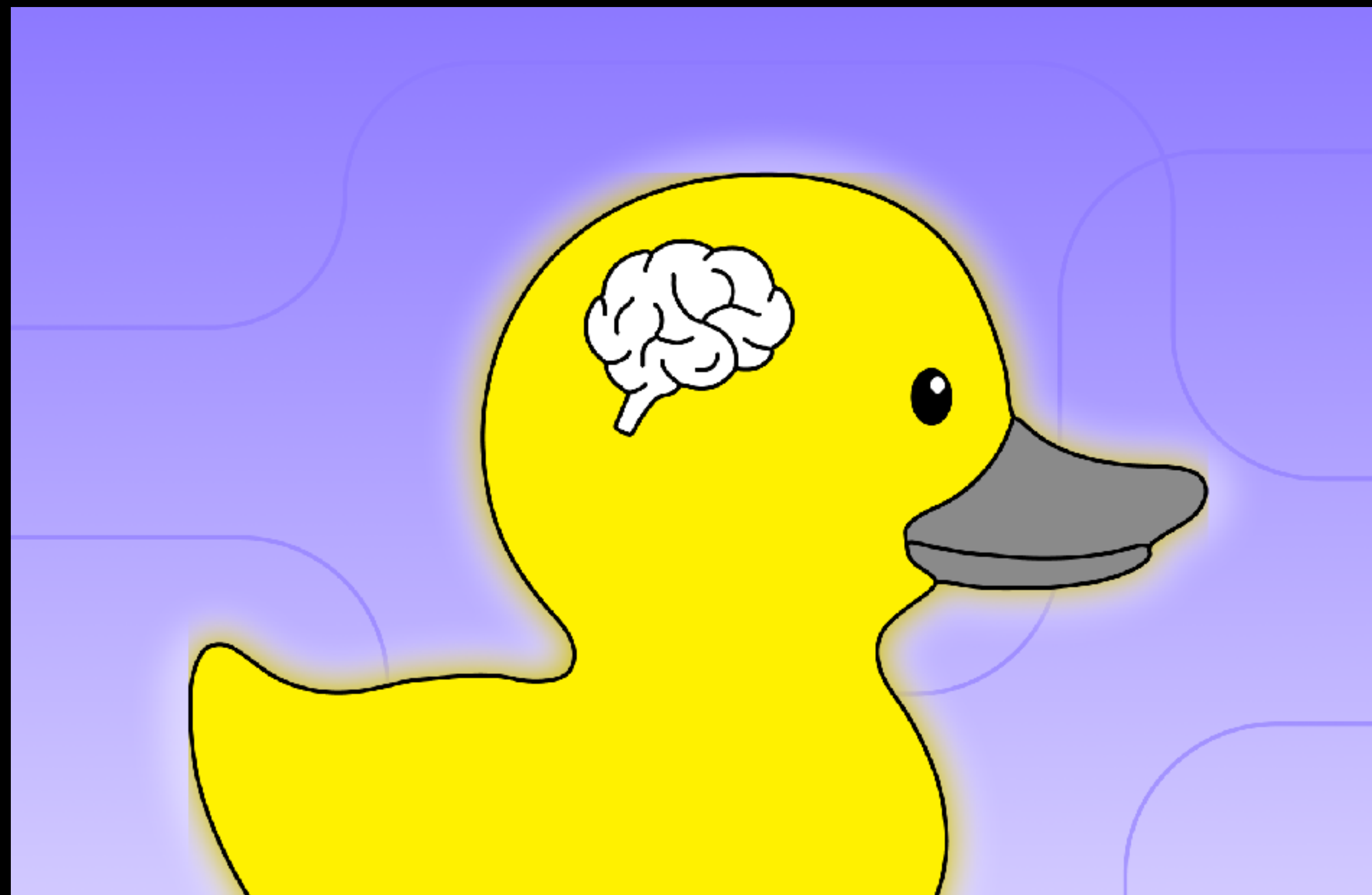


Database



Data lake

Text files




Open dataset: Railway services



Rijden de Treinen



 Download the free Rijden de Treinen app for iPhone and Android

Are the trains running?

Yes!

There are no large disruptions. There are engineering works on 9 lines.

Services between
2019 and 2025

150 million record
as CSV files

2.3 GB compressed,
20.0 GB uncompressed

Open dataset: Railway services

id int64	date date	station varchar	delay int64
738804	2019-01-01	Rotterdam Centraal	NULL
738804	2019-01-01	Delft	0
738804	2019-01-01	Den Haag HS	1
738804	2019-01-01	Leiden Centraal	0
738804	2019-01-01	Schiphol Airport	0

Open dataset: Railway services

id int64	date date	station varchar	delay int64
738804	2019-01-01	Rotterdam Centraal	NULL
738804	2019-01-01	Delft	0
738804	2019-01-01	Den Haag HS	1
738804	2019-01-01	Leiden Centraal	0
738804	2019-01-01	Schiphol Airport	0

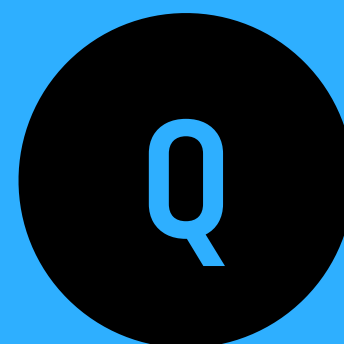


arrival delay

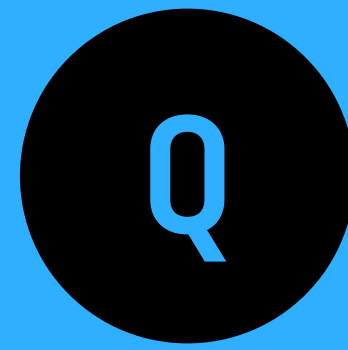
Open dataset: Railway services

id int64	date date	station varchar	delay int64
738804	2019-01-01	Rotterdam Centraal	NULL
738804	2019-01-01	Delft	0
738804	2019-01-01	Den Haag HS	1
738804	2019-01-01	Leiden Centraal	0
738804	2019-01-01	Schiphol Airport	0

arrival delay



Which train stations had the worst delays on the weekends of Feb 2025?

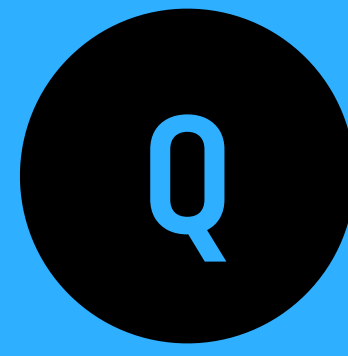


Which train stations had the worst delays on the weekends of Feb 2025?

```
import pandas as pd

df = pd.read_csv("services.csv")
df["date"] = pd.to_datetime(df["date"])

result = (
    df[
        df["date"].between("2025-02-01", "2025-02-28") &
        df["date"].dt.isocalendar().day.isin([6, 7])
    ]
    .groupby("station", as_index=False)["delay"]
    .mean()
    .rename(columns={"delay": "avg_delay"})
    .sort_values("avg_delay", ascending=False)
    .head(3)
)
```



Which train stations had the worst delays on the weekends of Feb 2025?

```
import pandas as pd
.
df = pd.read_csv("services.csv")
df["date"] = pd.to_datetime(df["date"])

result = (
    df[
        df["date"].between("2025-02-01", "2025-02-28") &
        df["date"].dt.isocalendar().day.isin([6, 7])
    ]
    .groupby("station", as_index=False)["delay"]
    .mean()
    .rename(columns={"delay": "avg_delay"})
    .sort_values("avg_delay", ascending=False)
    .head(3)
)
```



Which train stations had the worst delays on the weekends of Feb 2025?

```
import pandas as pd

df = pd.read_csv("services.csv")
df["date"] = pd.to_datetime(df["date"])

result = (
    df[
        df["date"].between("2025-02-01", "2025-02-28") &
        df["date"].dt.isocalendar().day.isin([6, 7])
    ]
    .groupby("station", as_index=False)["delay"]
    .mean()
    .rename(columns={"delay": "avg_delay"})
    .sort_values("avg_delay", ascending=False)
    .head(3)
)
```



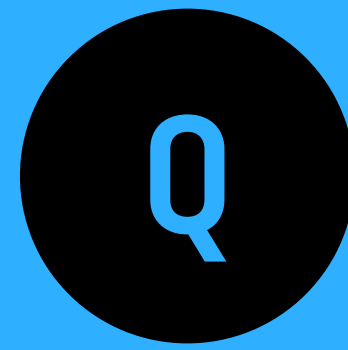
Which train stations had the worst delays on the weekends of Feb 2025?

```
import pandas as pd
```

```
df = pd.read_csv("services.csv")
```

```
df["date"] = pd.to_datetime(df["date"])
```

```
result = (  
    df[  
        df["date"].between("2025-02-01", "2025-02-28") &  
        df["date"].dt.isocalendar().day.isin([6, 7])  
    ]  
    .groupby("station", as_index=False)["delay"]  
    .mean()  
    .rename(columns={"delay": "avg_delay"})  
    .sort_values("avg_delay", ascending=False)  
    .head(3)  
)
```

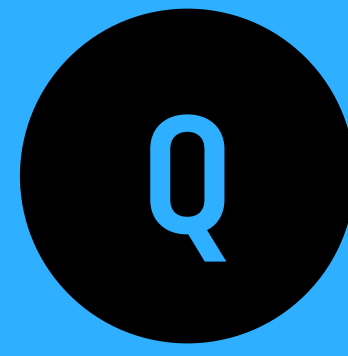


Which train stations had the worst delays on the weekends of Feb 2025?

```
import pandas as pd

df = pd.read_csv("services.csv")
df["date"] = pd.to_datetime(df["date"])

result = (
    df[
        df["date"].between("2025-02-01", "2025-02-28") &
        df["date"].dt.isocalendar().day.isin([6, 7])
    ]
    .groupby("station", as_index=False)["delay"]
    .mean()
    .rename(columns={"delay": "avg_delay"})
    .sort_values("avg_delay", ascending=False)
    .head(3)
)
```



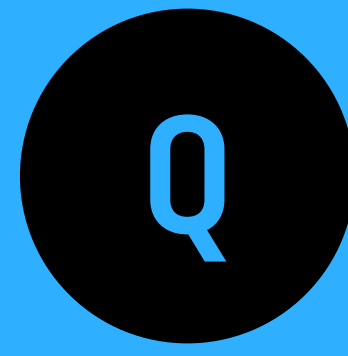
Which train stations had the worst delays on the weekends of Feb 2025?

```
import pandas as pd

df = pd.read_csv("services.csv")
df["date"] = pd.to_datetime(df["date"])

result = (
    df[
        df["date"].between("2025-02-01", "2025-02-28") &
        df["date"].dt.isocalendar().day.isin([6, 7])
    ]
    .groupby("station", as_index=False)["delay"]
    .mean()
    .rename(columns={"delay": "avg_delay"})
    .sort_values("avg_delay", ascending=False)
    .head(3)
)
```

	station	avg_delay
415	Siegburg/Bonn	5.583333
177	Emmerich-Elten	3.965398
103	Brussel-Zuid	3.864253



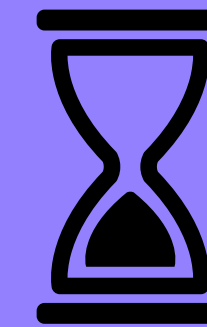
Which train stations had the worst delays on the weekends of Feb 2025?

```
import pandas as pd

df = pd.read_csv("services.csv")
df["date"] = pd.to_datetime(df["date"])

result = (
    df[
        df["date"].between("2025-02-01", "2025-02-28") &
        df["date"].dt.isocalendar().day.isin([6, 7])
    ]
    .groupby("station", as_index=False)["delay"]
    .mean()
    .rename(columns={"delay": "avg_delay"})
    .sort_values("avg_delay", ascending=False)
    .head(3)
)
```

	station	avg_delay
415	Siegburg/Bonn	5.583333
177	Emmerich-Elten	3.965398
103	Brussel-Zuid	3.864253



184 s





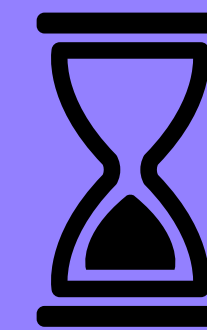
Which train stations had the worst delays on the weekends of Feb 2025?

```
import pandas as pd

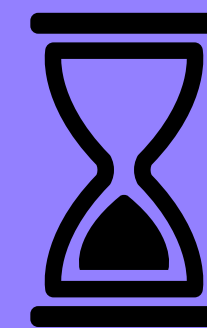
df = pd.read_csv("services.csv")
df["date"] = pd.to_datetime(df["date"])

result = (
    df[
        df["date"].between("2025-02-01", "2025-02-28") &
        df["date"].dt.isocalendar().day.isin([6, 7])
    ]
    .groupby("station", as_index=False)["delay"]
    .mean()
    .rename(columns={"delay": "avg_delay"})
    .sort_values("avg_delay", ascending=False)
    .head(3)
)
```

	station	avg_delay
415	Siegburg/Bonn	5.583333
177	Emmerich-Elten	3.965398
103	Brussel-Zuid	3.864253



184 s



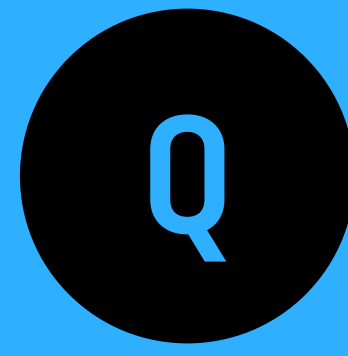
65 s



\$ duckdb



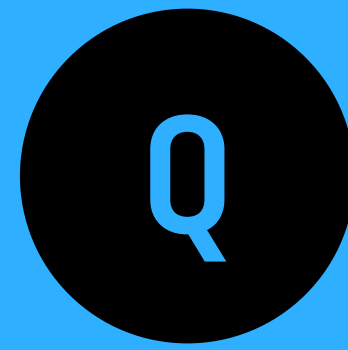
Which train stations had the worst delays on the weekends of Feb 2025?



Which train stations had the worst delays on the weekends of Feb 2025?

```
$ duckdb
```

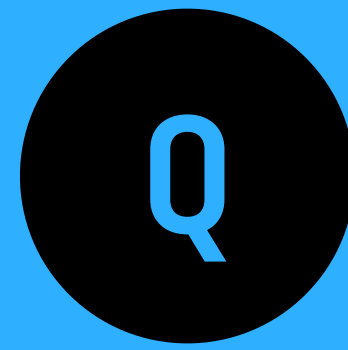
```
D SELECT
  station,
  avg(delay) AS avg_delay
FROM 'services.csv'
WHERE date BETWEEN '2025-02-01' AND '2025-02-28'
  AND extract(isodow FROM date) IN (6, 7)
GROUP BY station
ORDER BY avg_delay DESC
LIMIT 3;
```



Which train stations had the worst delays on the weekends of Feb 2025?

```
$ duckdb
```

```
D SELECT
  station,
  avg(delay) AS avg_delay
FROM 'services.csv'
WHERE date BETWEEN '2025-02-01' AND '2025-02-28'
  AND extract(isodow FROM date) IN (6, 7)
GROUP BY station
ORDER BY avg_delay DESC
LIMIT 3;
```



Which train stations had the worst delays on the weekends of Feb 2025?

```
$ duckdb
```

```
D SELECT
  station,
  avg(delay) AS avg_delay
FROM 'services.csv'
WHERE date BETWEEN '2025-02-01' AND '2025-02-28'
  AND extract(isodow FROM date) IN (6, 7)
GROUP BY station
ORDER BY avg_delay DESC
LIMIT 3;
```



Which train stations had the worst delays on the weekends of Feb 2025?

```
$ duckdb
```

```
D SELECT
```

```
    station,
```

```
    avg(delay) AS avg_delay
```

```
FROM 'services.csv'
```

```
WHERE date BETWEEN '2025-02-01' AND '2025-02-28'
```

```
    AND extract(isodow FROM date) IN (6, 7)
```

```
GROUP BY station
```

```
ORDER BY avg_delay DESC
```

```
LIMIT 3;
```



Which train stations had the worst delays on the weekends of Feb 2025?

```
$ duckdb
```

```
D SELECT
  station,
  avg(delay) AS avg_delay
FROM 'services.csv'
WHERE date BETWEEN '2025-02-01' AND '2025-02-28'
  AND extract(isodow FROM date) IN (6, 7)
GROUP BY station
ORDER BY avg_delay DESC
LIMIT 3;
```



Which train stations had the worst delays on the weekends of Feb 2025?

```
$ duckdb
```

```
D SELECT
  station,
  avg(delay) AS avg_delay
FROM 'services.csv'
WHERE date BETWEEN '2025-02-01' AND '2025-02-28'
  AND extract(isodow FROM date) IN (6, 7)
GROUP BY station
ORDER BY avg_delay DESC
LIMIT 3;
```

station	avg_delay
Siegburg/Bonn	5.58
Emmerich-Elten	3.97
Brussel-Zuid	3.86



Which train stations had the worst delays on the weekends of Feb 2025?

```
$ duckdb
```

```
D SELECT
  station,
  avg(delay) AS avg_delay
FROM 'services.csv'
WHERE date BETWEEN '2025-02-01' AND '2025-02-
  AND extract(isodow FROM date) IN (6, 7)
GROUP BY station
ORDER BY avg_delay DESC
LIMIT 3;
```



184 s

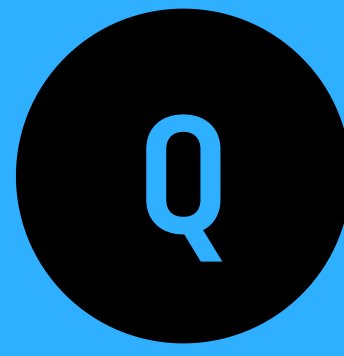
 pandas



65 s

 Polars

station	avg_delay
Siegburg/Bonn	5.58
Emmerich-Elten	3.97
Brussel-Zuid	3.86



Which train stations had the worst delays on the weekends of Feb 2025?

```
$ duckdb
```

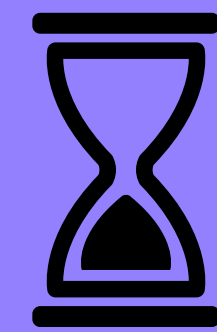
```
D SELECT
  station,
  avg(delay) AS avg_delay
FROM 'services.csv'
WHERE date BETWEEN '2025-02-01' AND '2025-02-
  AND extract(isodow FROM date) IN (6, 7)
GROUP BY station
ORDER BY avg_delay DESC
LIMIT 3;
```

station	avg_delay
Siegburg/Bonn	5.58
Emmerich-Elten	3.97
Brussel-Zuid	3.86



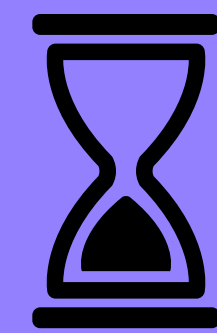
184 s

 pandas



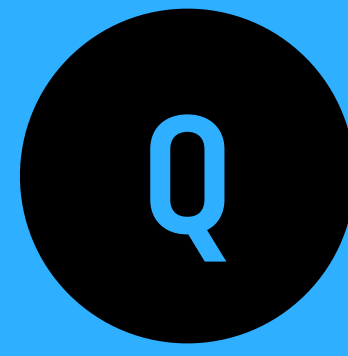
65 s

 Polars



7 s

 DuckDB



Which train stations had the worst delays on the weekends of Feb 2025?

```
$ duckdb
```

```
D SELECT
  station,
  avg(delay) AS avg_delay
FROM 'services.csv'
WHERE date BETWEEN '2025-02-01' AND '2025-02-
  AND extract(isodow FROM date) IN (6, 7)
GROUP BY station
ORDER BY avg_delay DESC
LIMIT 3;
```

station	avg_delay
---------	-----------



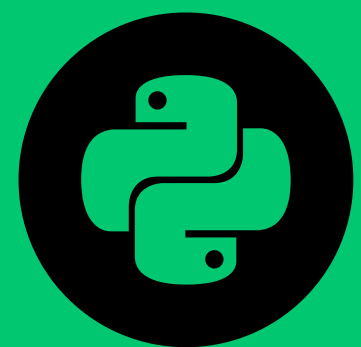
184 s



65 s



7 s



How about using Python?



Which train stations had the worst delays on the weekends of Feb 2025?

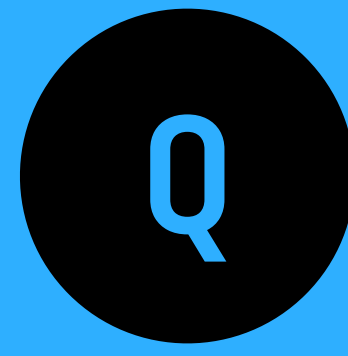
```
$ python
```

```
import duckdb
```

```
con = duckdb.connect()
```

```
rel = con.sql("""  
    SELECT station, avg(delay) AS avg_delay  
    FROM 'services.csv'  
    WHERE date BETWEEN '2025-02-01' AND '2025-02-28'  
        AND extract(isodow FROM date) IN (6, 7)  
    GROUP BY station  
    ORDER BY avg_delay DESC  
    LIMIT 3;""")
```

```
rel.show()
```



Which train stations had the worst delays on the weekends of Feb 2025?

```
$ python
```

```
import duckdb
```

```
con = duckdb.connect()
```

```
rel = con.sql("""  
    SELECT station, avg(delay) AS avg_delay  
    FROM 'services.csv'  
    WHERE date BETWEEN '2025-02-01' AND '2025-02-28'  
    AND extract(isodow FROM date) IN (6, 7)  
    GROUP BY station  
    ORDER BY avg_delay DESC  
    LIMIT 3;""")
```

```
rel.show()
```



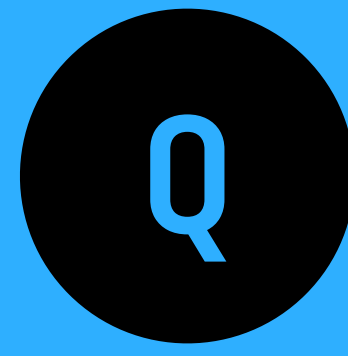
7 s

```
$ python
```

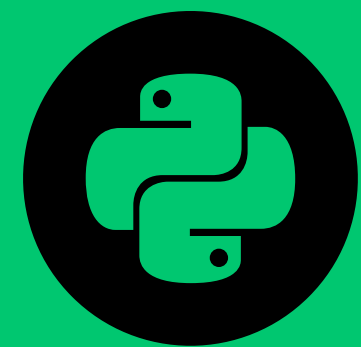
```
import duckdb
```

```
con = duckdb.connect()
```

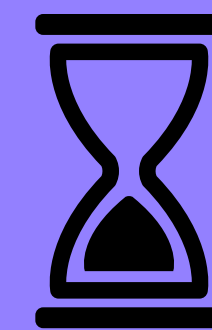
```
rel = con.sql("""  
    SELECT station, avg(delay) AS avg_delay  
    FROM 'services.csv'  
    WHERE date BETWEEN '2025-02-01' AND '2025-02-28'  
        AND extract(isodow FROM date) IN (6, 7)  
    GROUP BY station  
    ORDER BY avg_delay DESC  
    LIMIT 3;""")
```



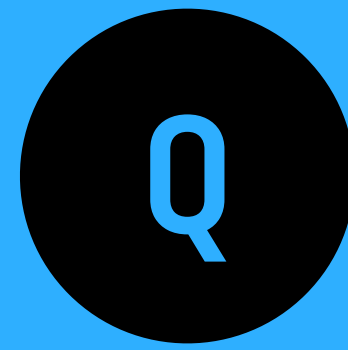
Which train stations had the worst delays on the weekends of Feb 2025?



But how about *being Pythonic*?



7 s



Which train stations had the worst delays on the weekends of Feb 2025?

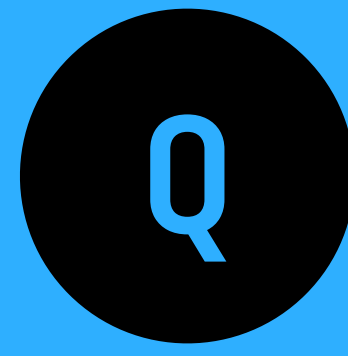
```
$ python
```

```
import duckdb
```

```
con = duckdb.connect()
```

```
rel = (  
    con.table("services.csv")  
    .filter("""date BETWEEN '2025-02-01' AND '2025-02-28'  
            AND extract(isodow FROM date) IN (6, 7)""")  
    .aggregate("station, avg(delay) AS avg_delay", "station")  
    .order("avg_delay DESC")  
    .limit(3)  
)
```

```
rel.show()
```



Which train stations had the worst delays on the weekends of Feb 2025?

```
$ python

import duckdb

con = duckdb.connect()

rel = (
    con.table("services.csv")
    .filter("""date BETWEEN '2025-02-01' AND '2025-02-28'
            AND extract(isodow FROM date) IN (6, 7)""")
    .aggregate("station, avg(delay) AS avg_delay", "station")
    .order("avg_delay DESC")
    .limit(3)
)

rel.show()
```



7s

```
$ python
```

```
import duckdb
```

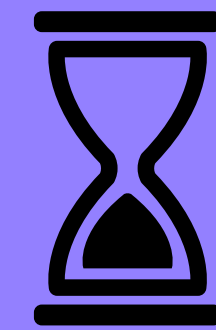
```
con = duckdb.connect()
```

```
rel = (  
    con.table("services.csv")  
    .filter("""date BETWEEN '2025-02-01' AND '2025-02-28'  
            AND extract(isodow FROM date) IN (6, 7)""")  
    .aggregate("station, avg(delay) AS avg_delay", "station")  
    .order("avg_delay DESC")  
    .limit(3)  
)
```

```
rel.show()
```



Which train stations had the worst delays on the weekends of Feb 2025?

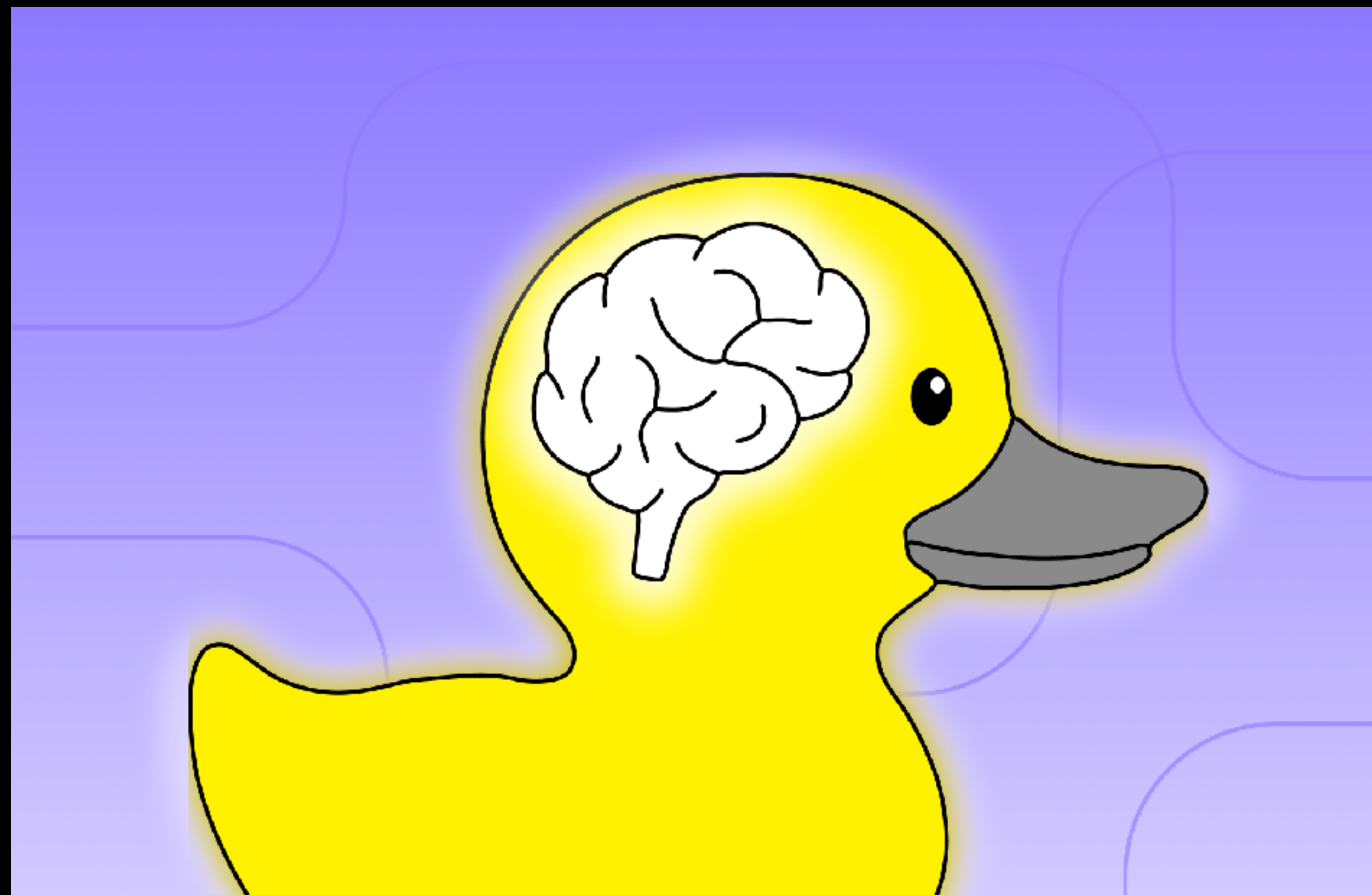


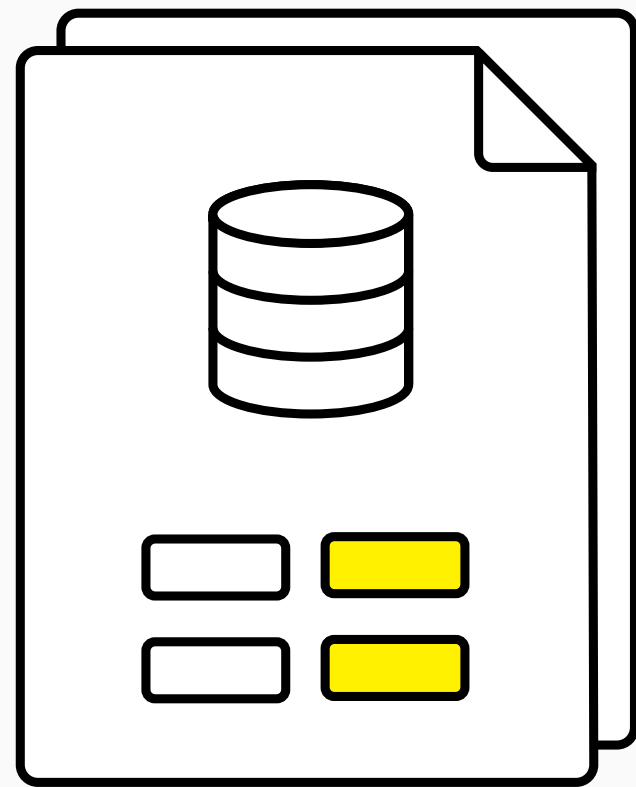
7s



Which one do you prefer?

Binary files





Tables on disk:

Row- or column-based storage?

Row-based storage

date	station	delay
Feb 28	Delft	0
Feb 28	Edam	0
Feb 28	Breda	0
Feb 28	Delft	2

Row-based storage

date	station	delay
Feb 28	Delft	0
Feb 28	Edam	0
Feb 28	Breda	0
Feb 28	Delft	2

Row-based storage

Feb 28

Delft

0

Feb 28

Edam

0

Feb 28

Breda

0

Feb 28

Delft

2

Row-based storage

Delft 0

Edam 0

Breda 0

Delft 2

Row-based storage

Delft 0

Edam 0

Breda 0

Delft 2

transactional databases, e.g., PostgreSQL

Avro format

Column-based storage

date	station	delay
Feb 28	Delft	0
Feb 28	Edam	0
Feb 28	Breda	0
Feb 28	Delft	2

Column-based storage

Feb 28

Feb 28

Feb 28

Feb 28

Delft

Edam

Breda

Delft

0

0

0

2

Column-based storage

Feb 28

Delft

Edam

Breda

Delft

0

2

Column-based storage

Feb 28

Delft

Edam

Breda

Delft

0 (x3)

2

Column-based storage

	Delft	Edam	Breda	Delft	0 (x3)	2
--	-------	------	-------	-------	--------	---

Column-based storage

	Delft	Edam	Breda	Delft	0 (x3)	2
--	-------	------	-------	-------	--------	---

analytical databases, e.g., DuckDB

Parquet format

Column-based storage

Delft Edam Breda Delft 0 (x3) 2

analytical databases, e.g., DuckDB

Parquet format



Who has heard of Apache Parquet?

```
$ python pandas-analysis-parquet.py
```



30 s

```
$ python polars-analysis-parquet.py
```



3 s

```
DuckDB: SELECT ... FROM '...parquet';
```



0.06 s

Parquet format

row group 1

row group 2

Metadata

schema,
statistics

date	station	delay
DATE	VARCHAR	INT

date	station	delay
DATE	VARCHAR	INT

min

Feb 28

Breda

0

min

March 1

Delft

0

max

Feb 28

Edam

2

max

March 2

Utrecht

5

Data

column-based
storage

Feb 28

Delft

0

Feb 28

Edam

0

Feb 28

Breda

0

Feb 28

Delft

2

March 1

Utrecht

3

March 1

Gouda

0

March 1

Edam

1

March 2

Delft

5

Parquet format

row group 1

row group 2

Metadata

schema,
statistics

date	station	delay
DATE	VARCHAR	INT

date	station	delay
DATE	VARCHAR	INT

min

Feb 28

Breda

0

min

March 1

Delft

0

max

Feb 28

Edam

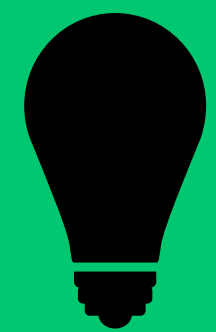
2

max

March 2

Utrecht

5



Built-in
min-max
indexes

Data

column-based
storage

Feb 28

Delft

0

Feb 28

Edam

0

Feb 28

Breda

0

Feb 28

Delft

2

March 1

Utrecht

3

March 1

Gouda

0

March 1

Edam

1

March 2

Delft

5

Parquet format

row group 1

row group 2

Metadata

schema,
statistics

date	station	delay
DATE	VARCHAR	INT

date	station	delay
DATE	VARCHAR	INT

min	Feb 28	Breda	0
max	Feb 28	Edam	2

min	March 1	Delft	0
max	March 2	Utrecht	5

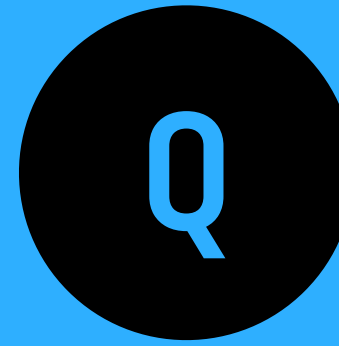
Data

column-based
storage

Feb 28	Delft	0
Feb 28	Edam	0
Feb 28	Breda	0
Feb 28	Delft	2

March 1	Utrecht	3
March 1	Gouda	0
March 1	Edam	1
March 2	Delft	5

Parquet format



Which train stations had the worst delays on the weekends of Feb (2025)?

row group 1

row group 2

Metadata

schema,
statistics

	date DATE	station VARCHAR	delay INT
min	Feb 28	Breda	0
max	Feb 28	Edam	2

	date DATE	station VARCHAR	delay INT
min	March 1	Delft	0
max	March 2	Utrecht	5

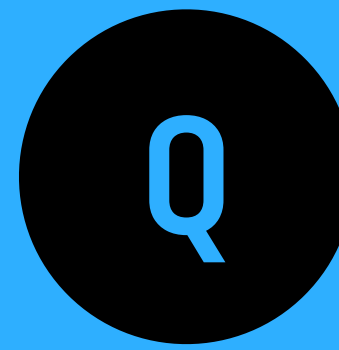
Data

column-based
storage

Feb 28	Delft	0
Feb 28	Edam	0
Feb 28	Breda	0
Feb 28	Delft	2

March 1	Utrecht	3
March 1	Gouda	0
March 1	Edam	1
March 2	Delft	5

Parquet format



Which train stations had the worst delays on the weekends of Feb (2025)?

row group 1

row group 2

Metadata

schema,
statistics

	date DATE	station VARCHAR	delay INT
min	Feb 28	Breda	0
max	Feb 28	Edam	2

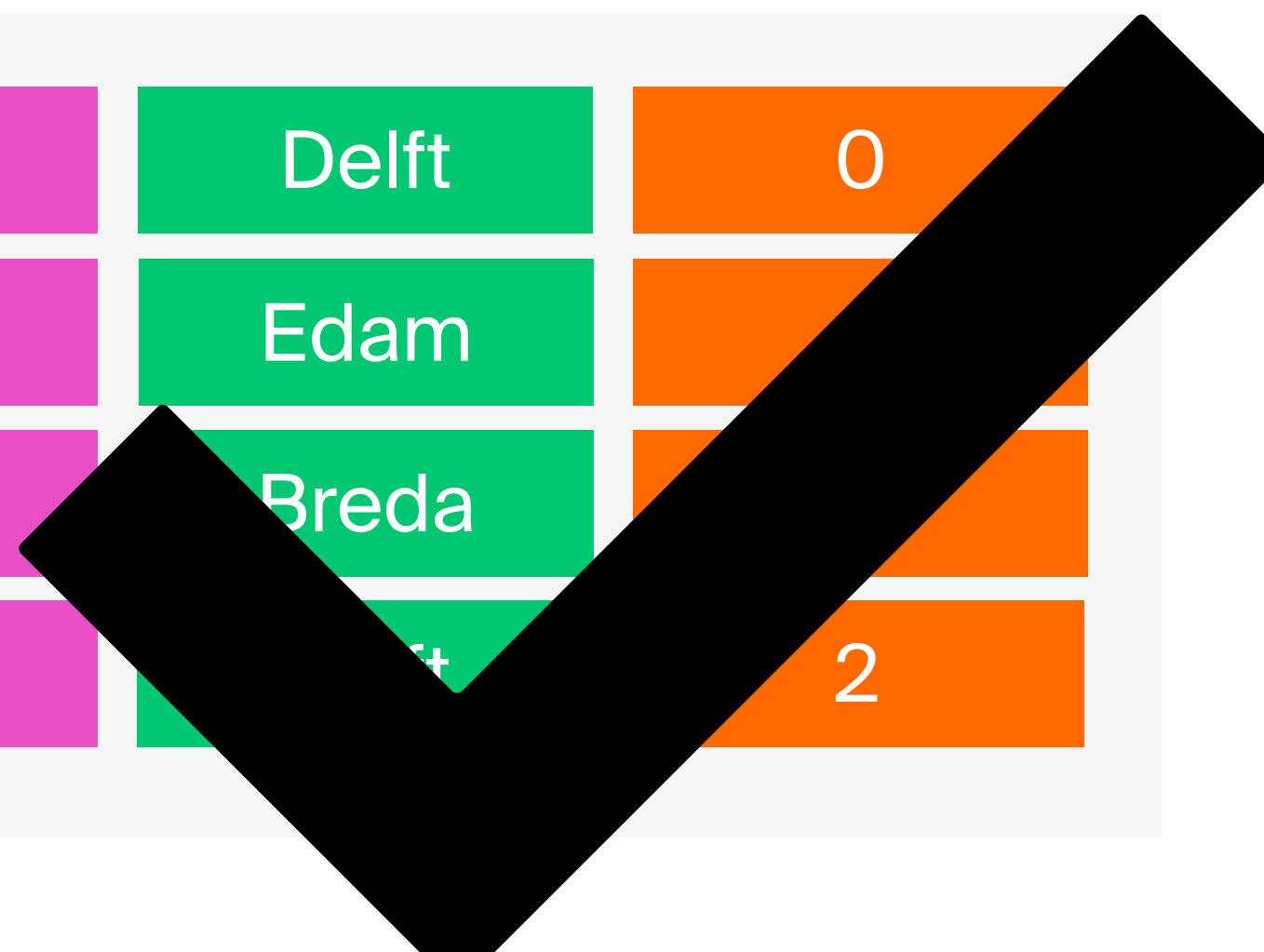
	date DATE	station VARCHAR	delay INT
min	March 1	Delft	0
max	March 2	Utrecht	5

Data

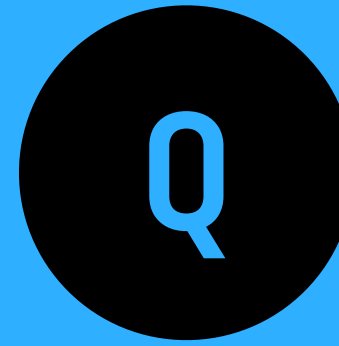
column-based
storage

Feb 28	Delft	0
Feb 28	Edam	0
Feb 28	Breda	0
Feb 28	Delft	2

March 1	Utrecht	3
March 1	Gouda	0
March 1	Edam	1
March 2	Delft	5



Parquet format



Which train stations had the worst delays on the weekends of Feb (2025)?

row group 1

row group 2

Metadata

schema,
statistics

	date DATE	station VARCHAR	delay INT
min	Feb 28	Breda	0
max	Feb 28	Edam	2

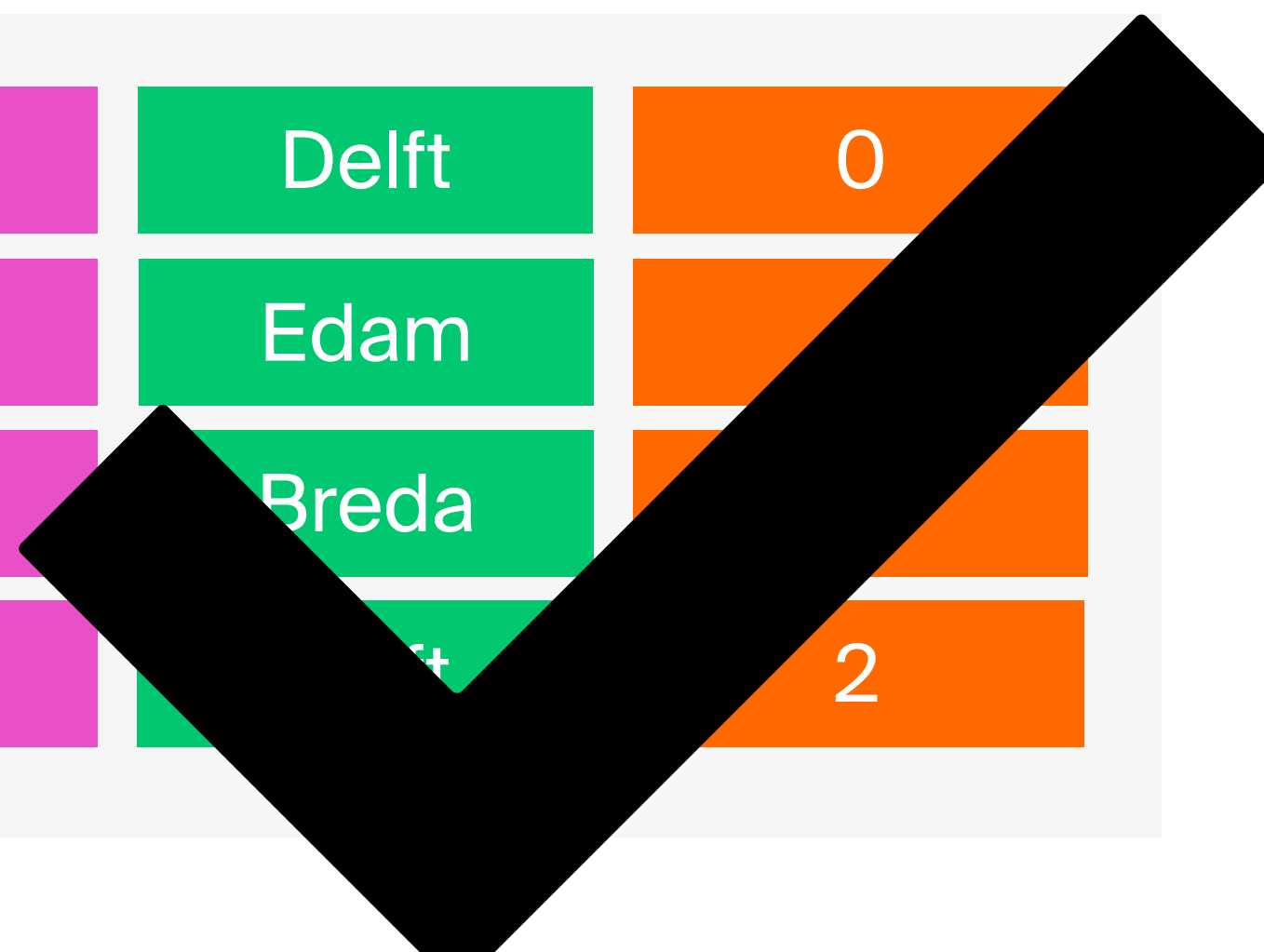
	date DATE	station VARCHAR	delay INT
min	March 1	Delft	0
max	March 2	Utrecht	5

Data

column-based
storage

Feb 28	Delft	0
Feb 28	Edam	0
Feb 28	Breda	0
Feb 28	Delft	2

March 1	Utrecht	3
March 1	Gouda	0
March 1	Edam	1
March 2	Delft	5



Parquet format



Which train stations had the worst delays on the weekends of Feb (2025)?

row group 1

row group 2

Metadata

schema,
statistics

	date DATE	station VARCHAR	delay INT
min	Feb 28	Breda	0
max	Feb 28	Edam	2

	date DATE	station VARCHAR	delay INT
min	March 1	Delft	0
max	March 2	Utrecht	5

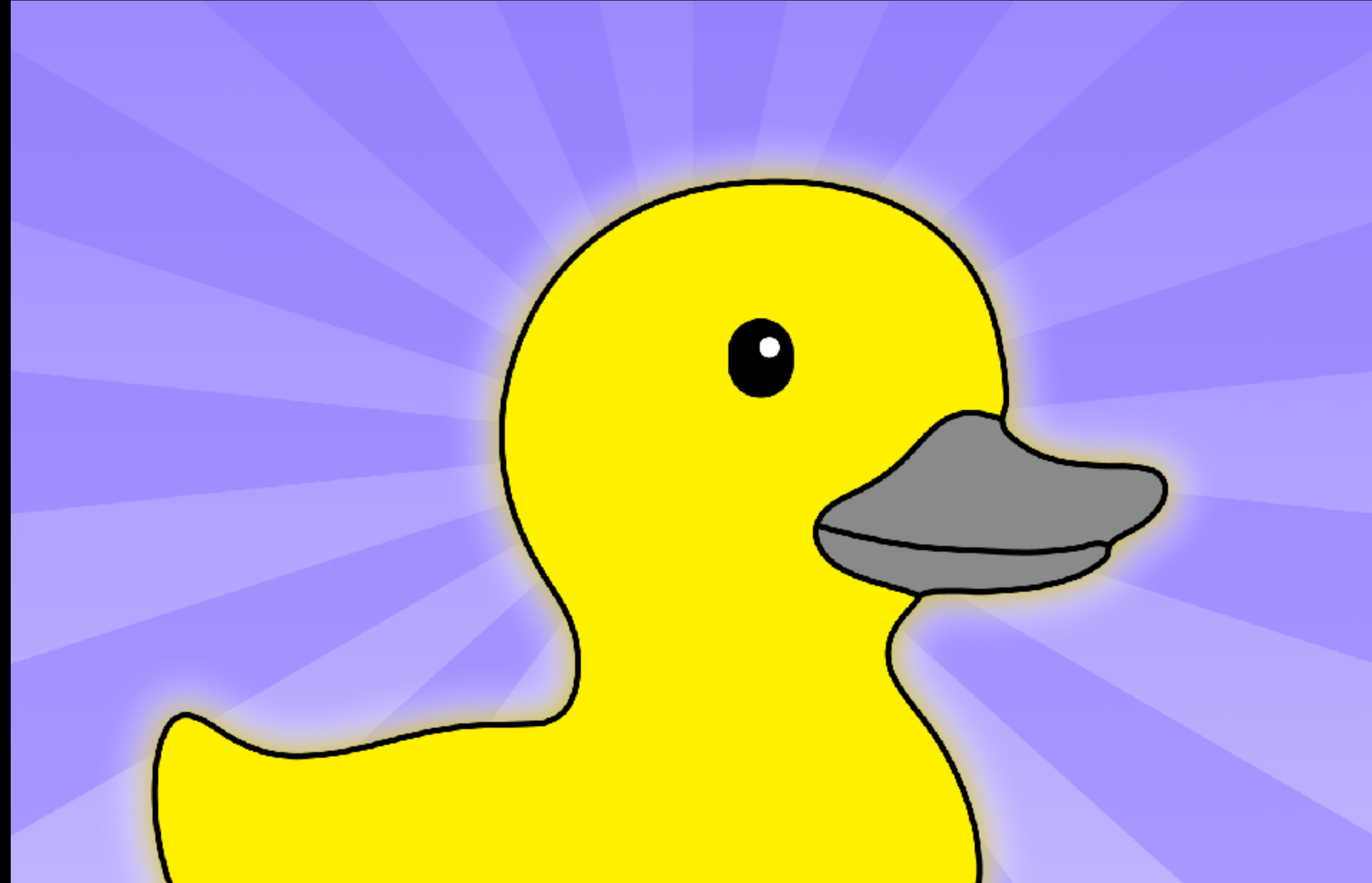
Data

column-based
storage

Feb 28	Delft	0
Feb 28	Edam	0
Feb 28	Breda	0
Feb 28	Delft	2

March 1	Utrecht	5
March 1	Delft	0
March 1	Breda	1
March 2	Edam	5

Database



```
$ duckdb services.duckdb
```

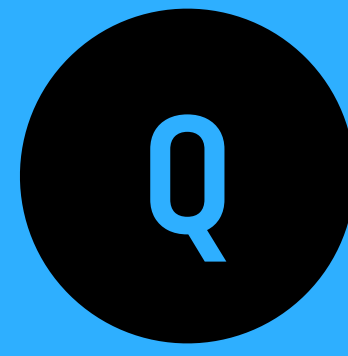
```
D CREATE TABLE services AS  
FROM 'services.csv';
```

```
$ duckdb services.duckdb
```

```
D CREATE TABLE services AS  
FROM 'services.csv';
```



17 s

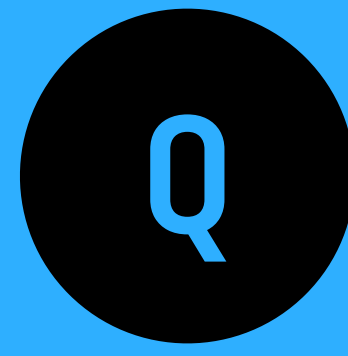


Which train stations had the worst delays on the weekends of Feb 2025?

```
$ duckdb services.duckdb
```

```
D SELECT
  station,
  avg(delay) AS avg_delay
FROM services
WHERE date BETWEEN '2025-02-01' AND '2025-02-28'
  AND extract(isodow FROM date) IN (6, 7)
GROUP BY station
ORDER BY avg_delay DESC
LIMIT 3;
```

station	avg_delay
Siegburg/Bonn	5.58
Emmerich-Elten	3.97
Brussel-Zuid	3.86

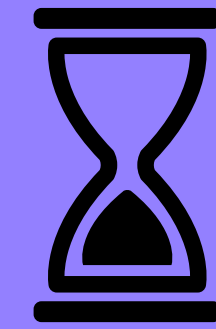


Which train stations had the worst delays on the weekends of Feb 2025?

```
$ duckdb services.duckdb
```

```
D SELECT
  station,
  avg(delay) AS avg_delay
FROM services
WHERE date BETWEEN '2025-02-01' AND '2025-02-28'
  AND extract(isodow FROM date) IN (6, 7)
GROUP BY station
ORDER BY avg_delay DESC
LIMIT 3;
```

station	avg_delay
Siegburg/Bonn	5.58
Emmerich-Elten	3.97
Brussel-Zuid	3.86



0.03s

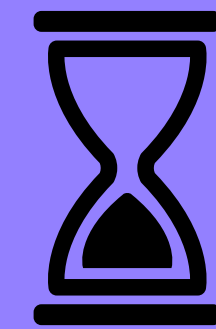


Which train stations had the worst delays on the weekends of Feb 2025?

```
$ duckdb services.duckdb
```

```
D SELECT
  station,
  avg(delay) AS avg_delay
FROM services
WHERE date BETWEEN '2025-02-01' AND '2025-02-28'
  AND extract(isodow FROM date) IN (6, 7)
GROUP BY station
ORDER BY avg_delay DESC
LIMIT 3;
```

station	avg_delay
Siegburg/Bonn	5.58
Emmerich-Elten	3.97
Brussel-Zuid	3.86

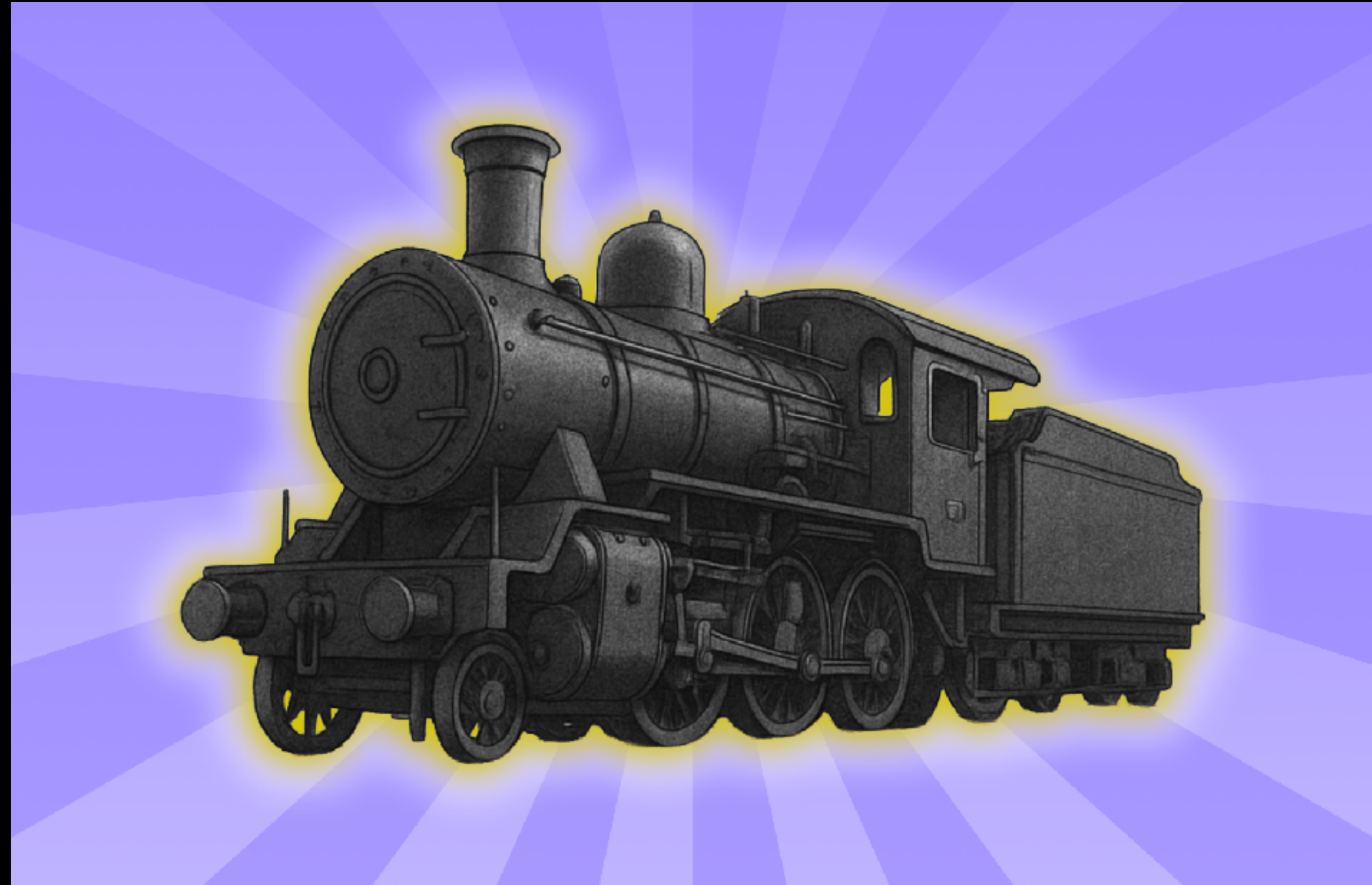


0.03s



The min-max indexes are created automatically for each row group

Railway operators



The biggest delay at Amsterdam Centraal is almost 3 hours!

```
SELECT date, train_number, delay
FROM services
WHERE date = '2025-02-27'
      AND station = 'Amsterdam Centraal'
      AND train_number = 420;
```

date	train_number	delay
2025-02-27	420	174

Let's renumber the services with large delays!

```
UPDATE services
SET train_number =
    train_number + (random() * 100 + 1)::INT
WHERE delay > 150;
```

(23 rows updated)

Let's renumber the services with large delays!

```
UPDATE services
SET train_number =
    train_number + (random() * 100 + 1)::INT
WHERE delay > 150;
```

(23 rows updated)

```
SELECT date, train_number, delay
FROM services
WHERE date = '2025-02-27'
    AND train_number = 420;
```

date	train_number	delay
0 rows		

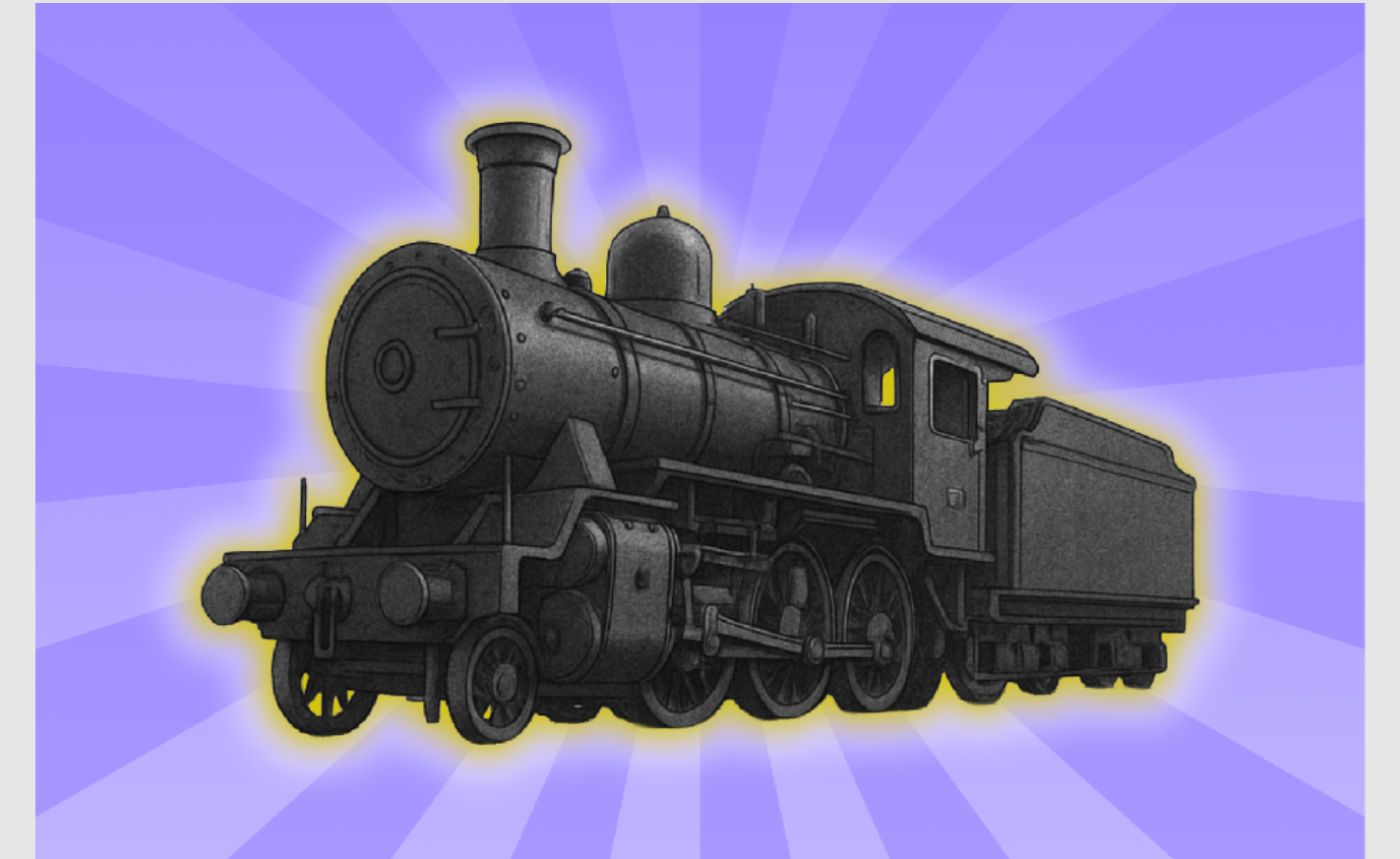
Let's renumber the services with large delays!

```
UPDATE services
SET train_number =
    train_number + (random() * 100 + 1)::INT
WHERE delay > 150;
```

(23 rows updated)

```
SELECT date, train_number, delay
FROM services
WHERE date = '2025-02-27'
    AND train_number = 420;
```

date	train_number	delay
0 rows		



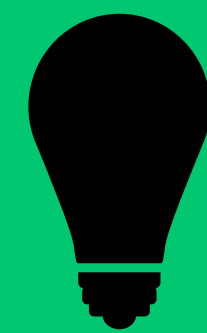
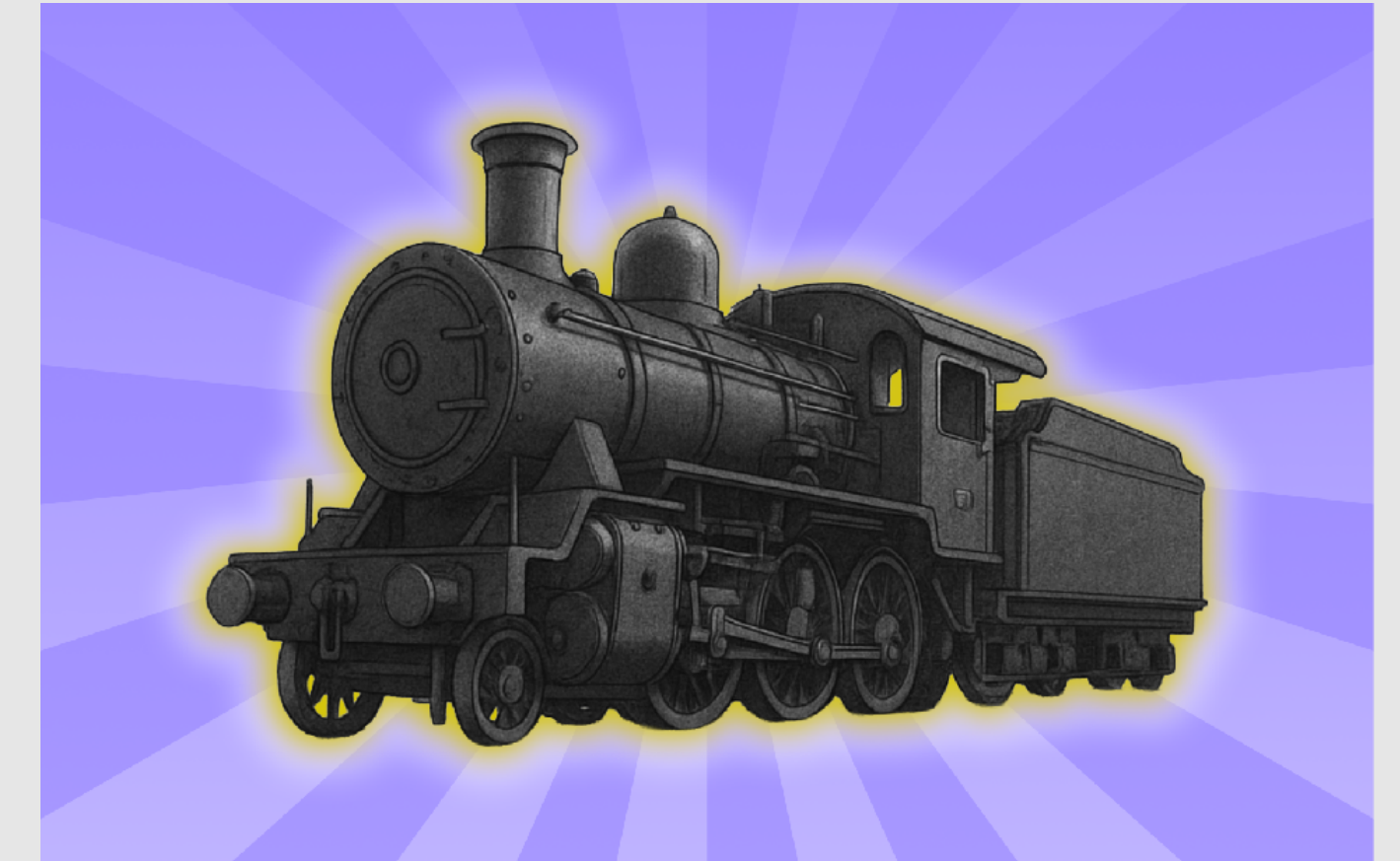
Let's renumber the services with large delays!

```
UPDATE services
SET train_number =
    train_number + (random() * 100 + 1)::INT
WHERE delay > 150;
```

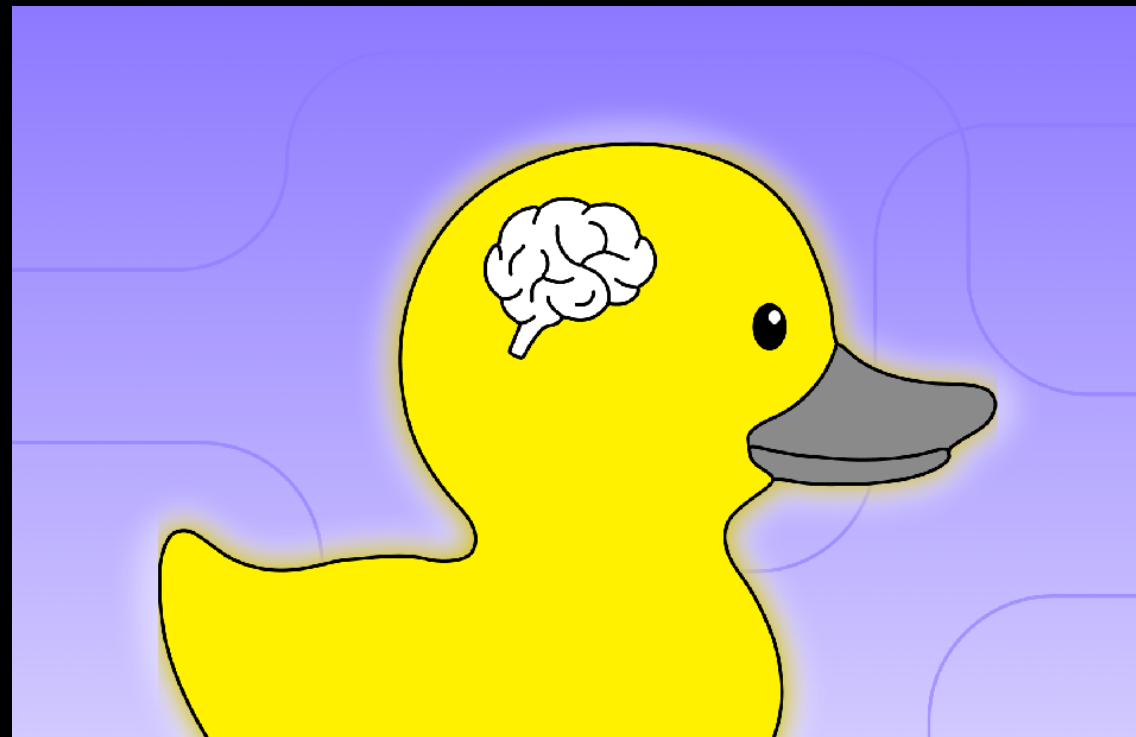
(23 rows updated)

```
SELECT date, train_number, delay
FROM services
WHERE date = '2025-02-27'
    AND train_number = 420;
```

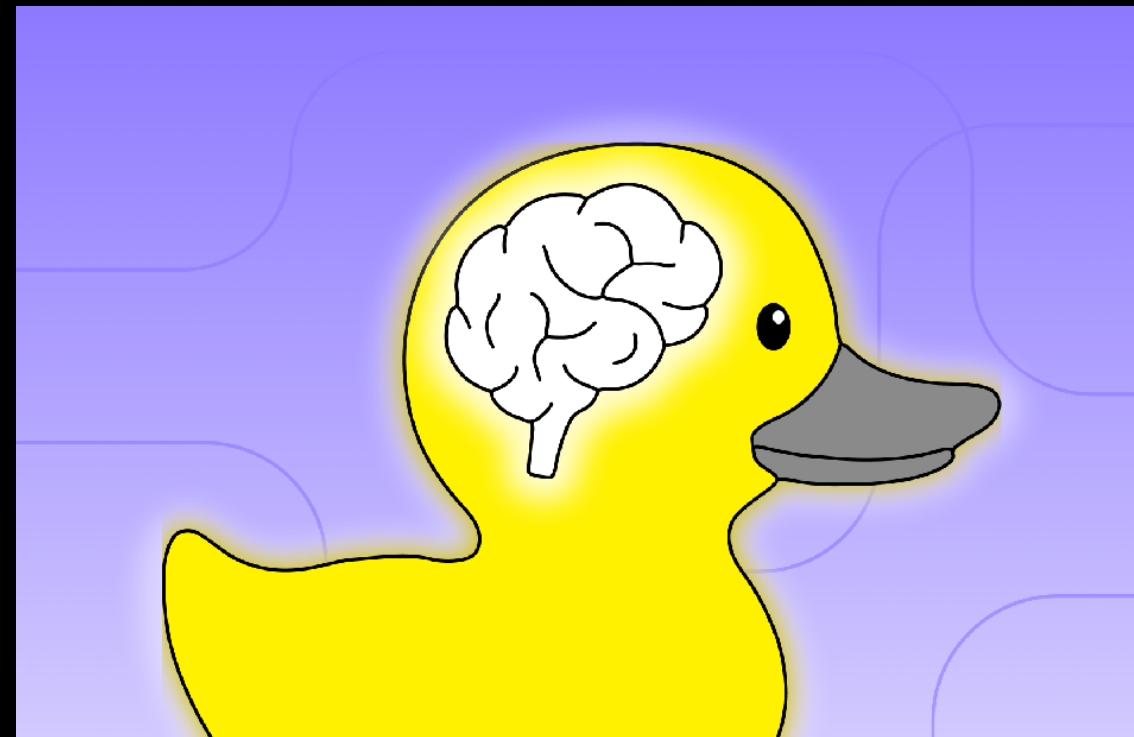
date	train_number	delay
0 rows		



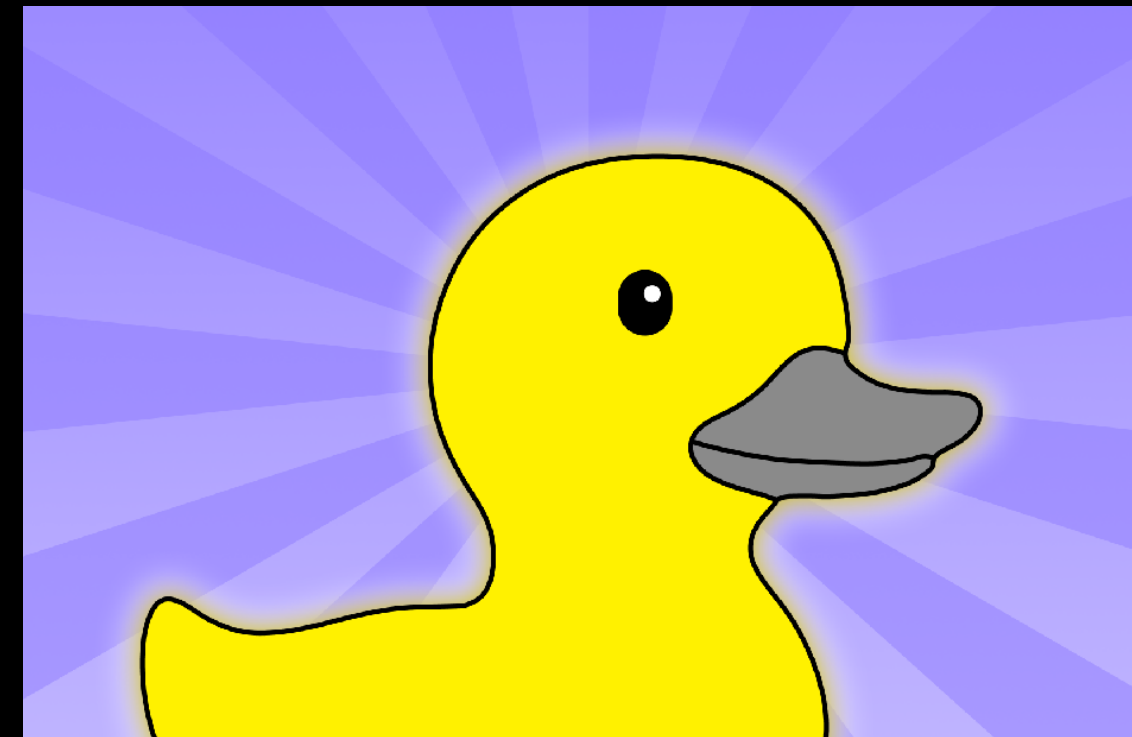
Databases are updatable!



Text files



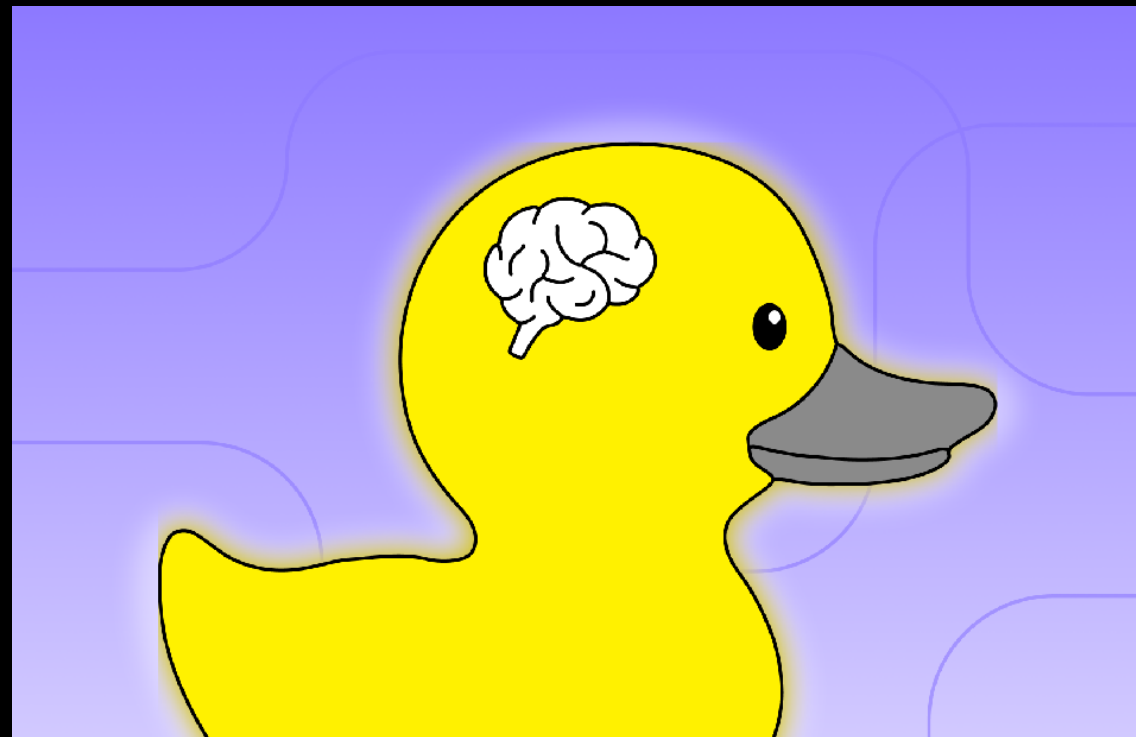
Binary files



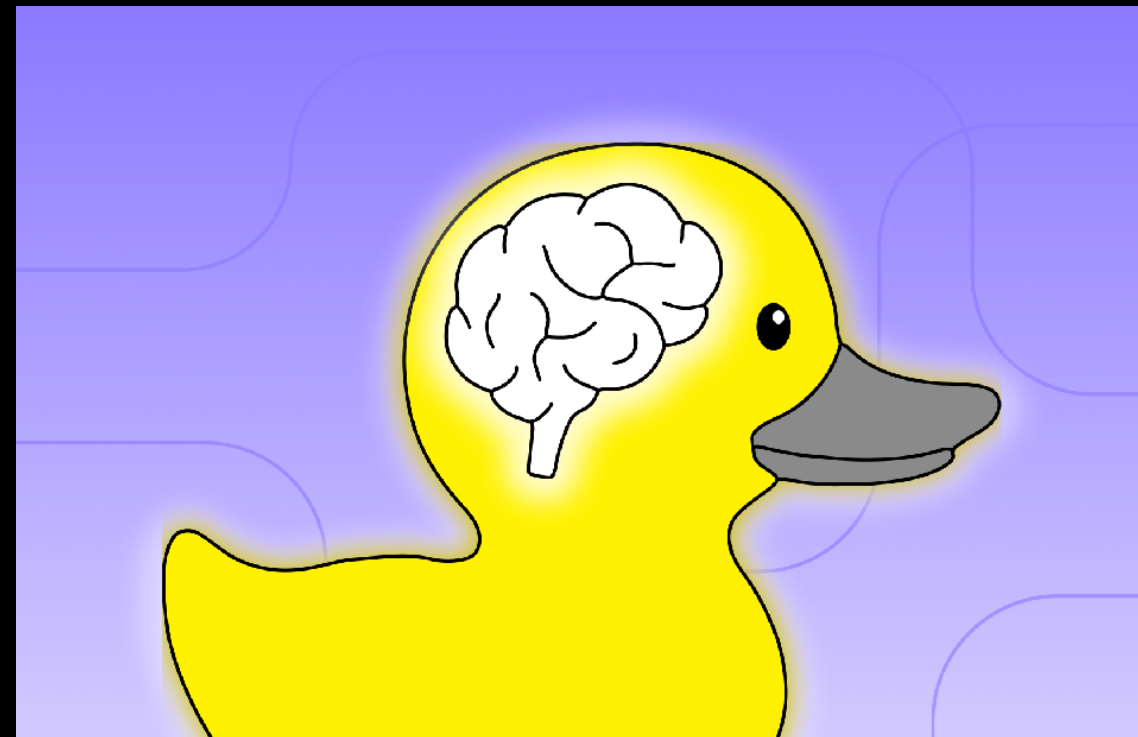
Database



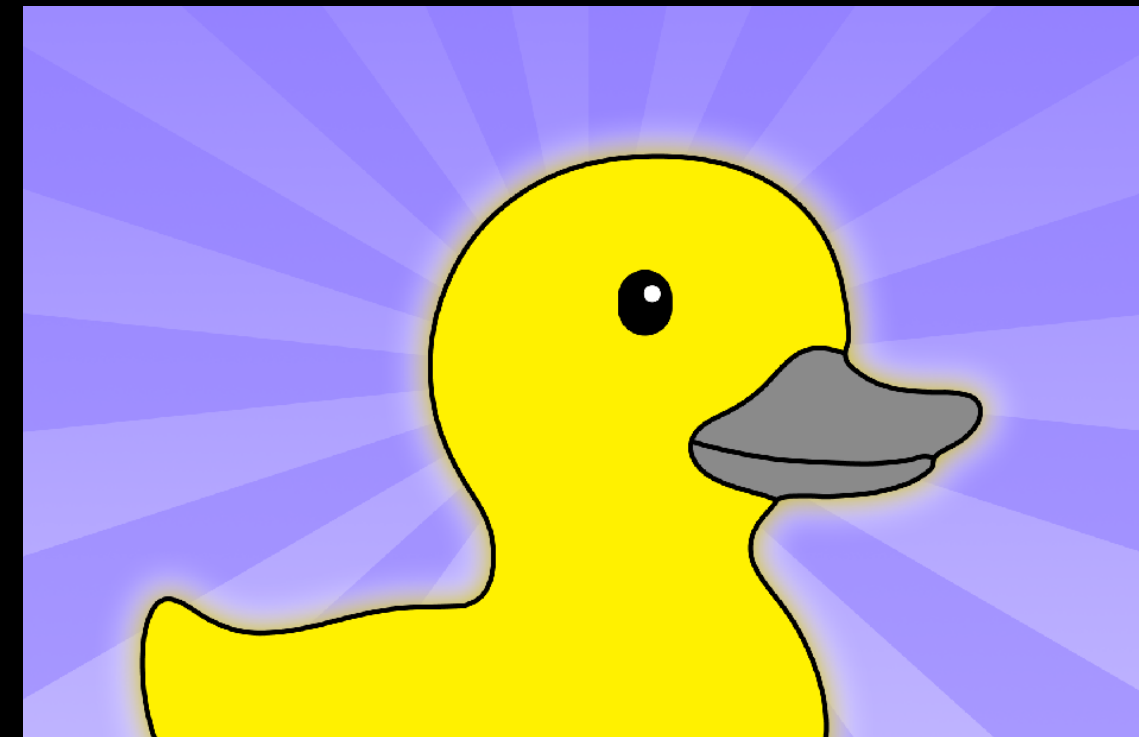
Data lake



Text files

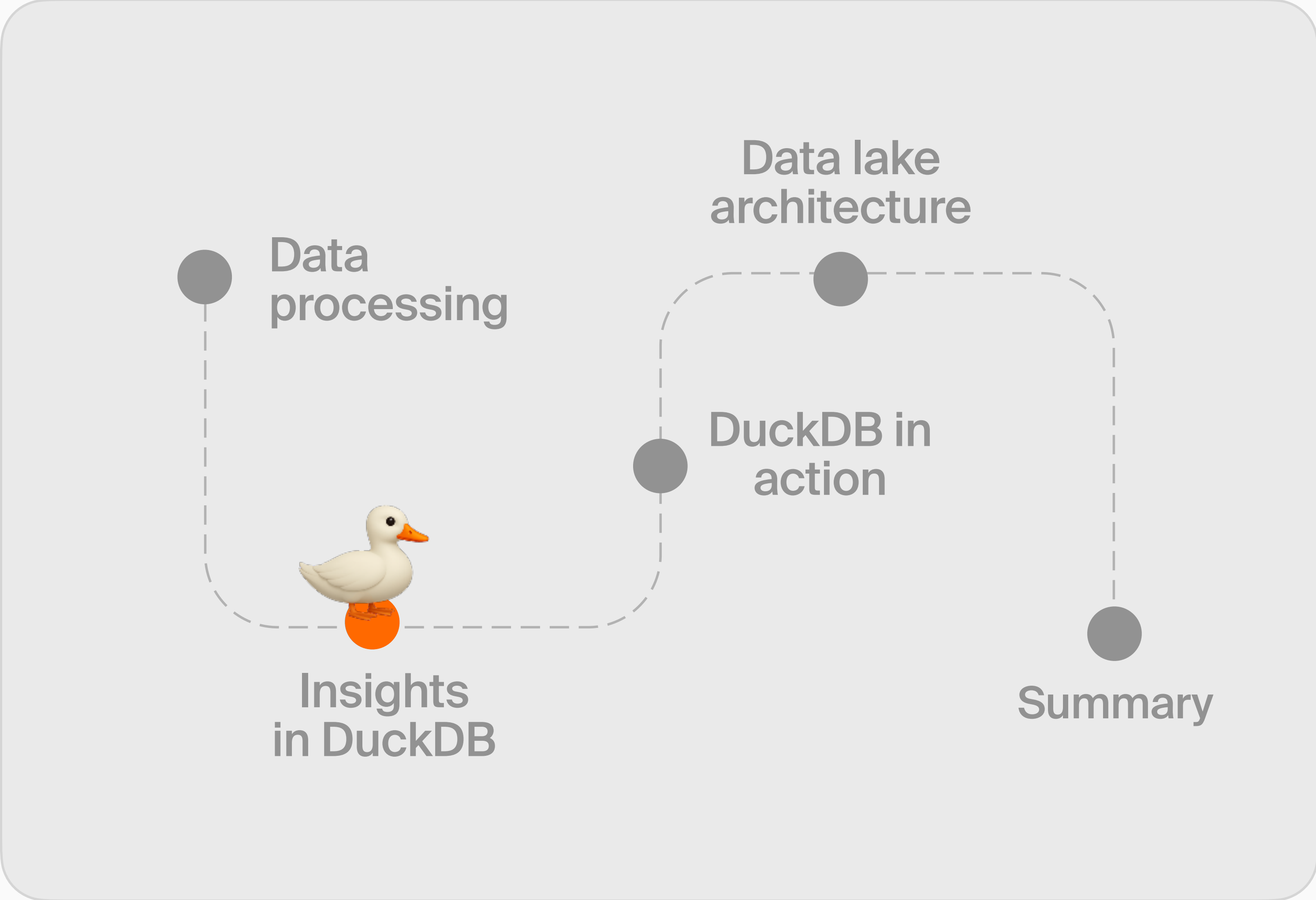


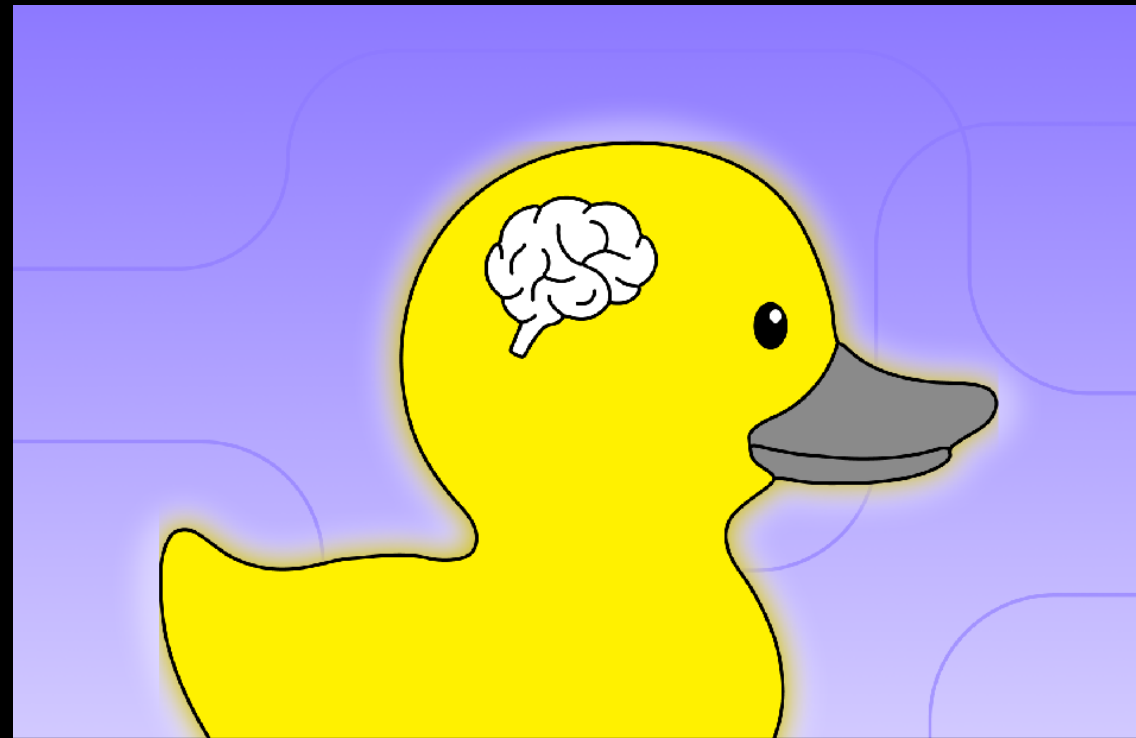
Binary files



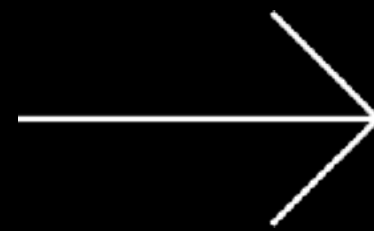
Database







Text files



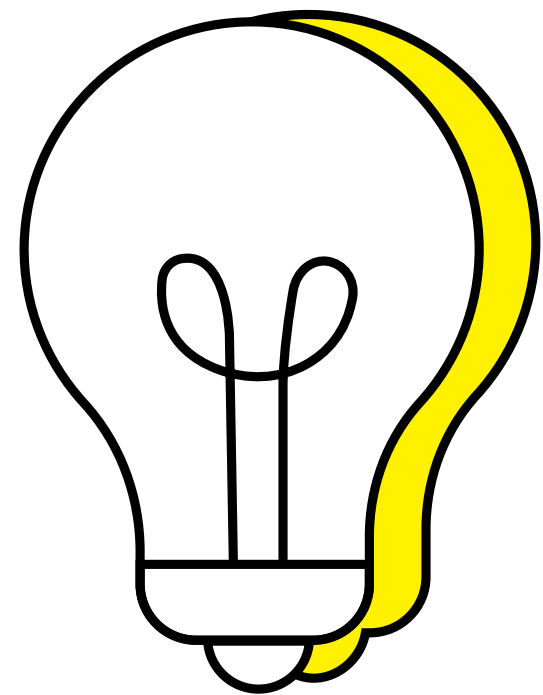
Database

Concise queries

Safe

Fast

Updatable



INSIGHT #1

**Databases can target
data science workloads**

Data science landscape, 2016

 pandas

 *hadoop*

?

data.tables

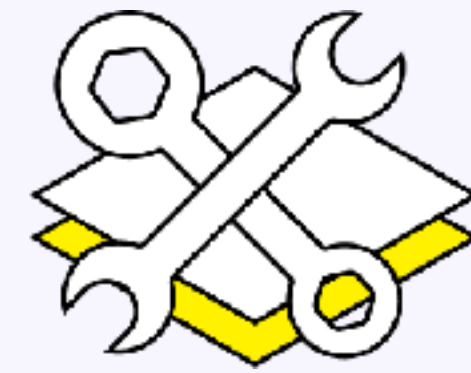
APACHE
*Spark*TM

data size 

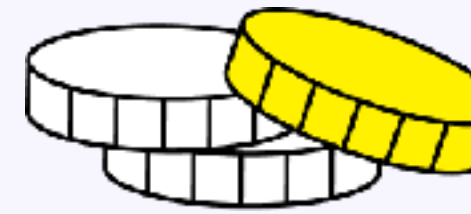
Resistance: Why would I use a database?

```
$ sudo apt install your_favorite_database
$ ... configure for 60 minutes ...
$ sudo service start your_favorite_database
$ wget https://blobs.duckdb.org/data/services.csv.gz
$ gunzip services.csv.gz
$ db_client -u admin -p admin

# CREATE TABLE ...; COPY ...; SELECT ...;
```



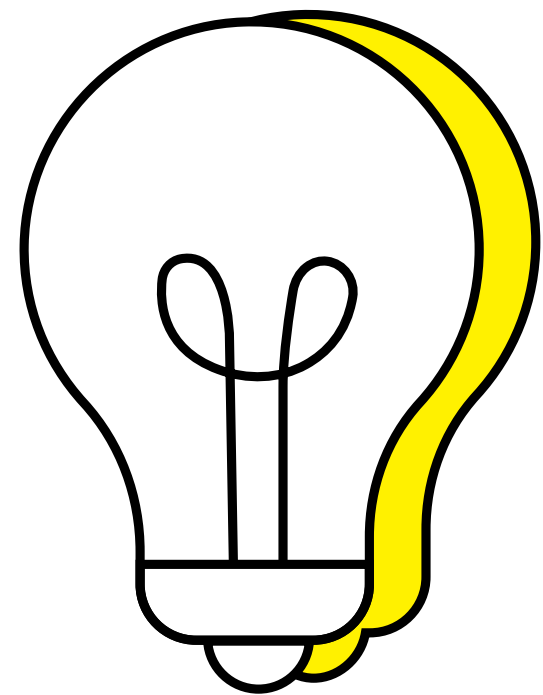
Difficult
operation



Expensive



Bad UX



INSIGHT #2

Focus on usability

End-to-end performance

setup

define schema

load

write queries

run queries

End-to-end performance

setup

define schema

write queries

load

run queries

End-to-end performance

setup

define schema

write queries

load

run queries

End-to-end performance

setup

define schema

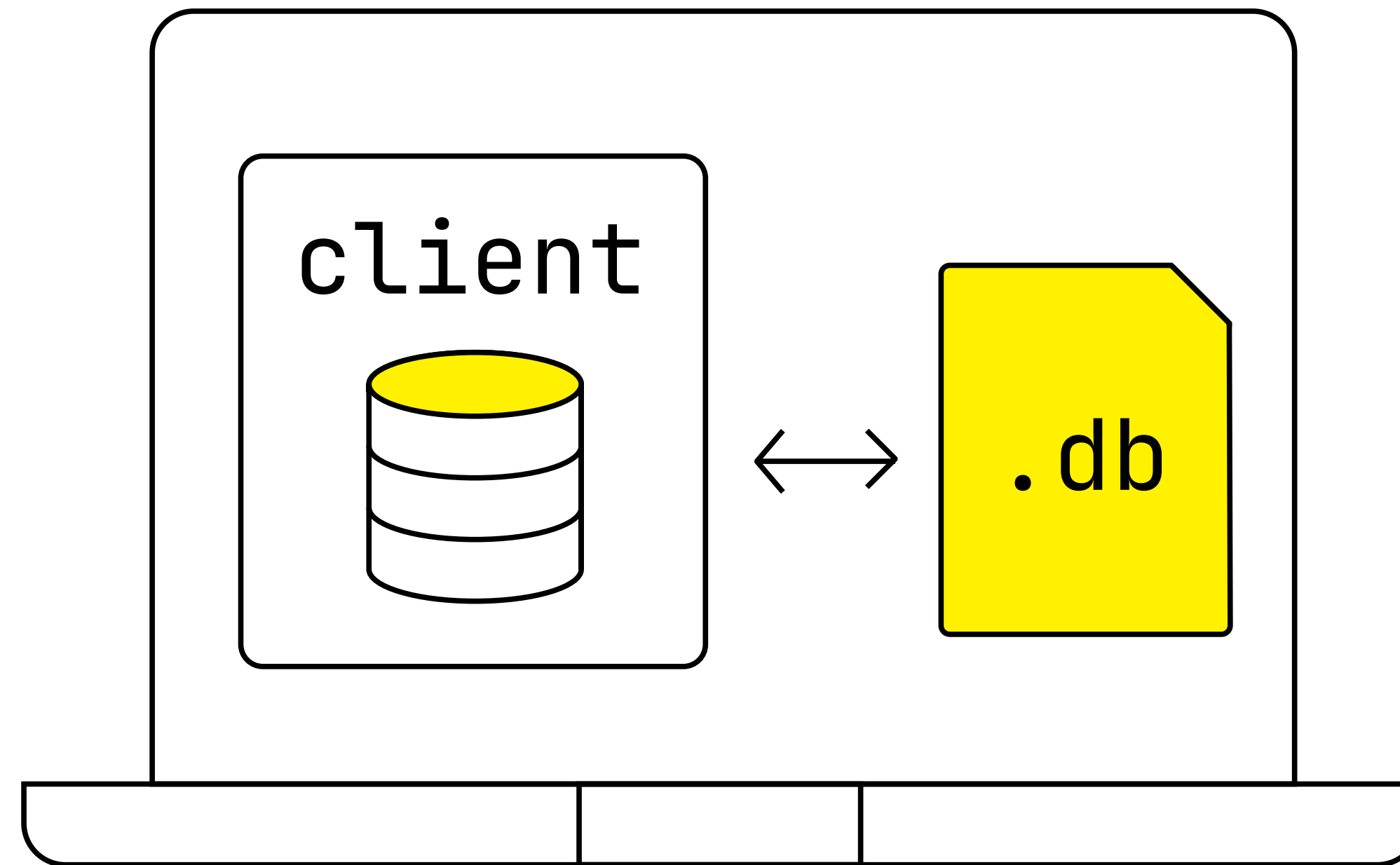
load

write queries

run queries

In-process architecture

setup



In-process architecture



```
pip install duckdb
```



```
curl https://install.duckdb.org | sh
```



```
go get github.com/duckdb/duckdb-go/v2
```



End-to-end performance

setup

define schema

load

write queries

run queries

“Friendly” SQL

define schema

write queries

```
SELECT
    date, station, avg(delay),
--    min(delay)
FROM 'https://blobs.duckdb.org/data/services.csv.gz'
GROUP BY ALL;
```

“Friendly” SQL

define schema

write queries

```
SELECT
    date, station, avg(delay),
--    min(delay)
FROM 'https://blobs.duckdb.org/data/services.csv.gz'
GROUP BY ALL;
```

“Friendly” SQL

define schema

write queries

```
SELECT
    date, station, avg(delay),
-- min(delay)
FROM 'https://blobs.duckdb.org/data/services.csv.gz'
GROUP BY ALL;
```

“Friendly” SQL

define schema

write queries

```
SELECT
    date, station, avg(delay),
-- min(delay)
FROM 'https://blobs.duckdb.org/data/services.csv.gz'
GROUP BY ALL;
```

“Friendly” SQL

define schema

write queries

```
SELECT
  date, station, avg(delay),
  -- min(delay)
FROM 'https://blobs.duckdb.org/data/services.csv.gz'
GROUP BY ALL;
```



2022



2022



2023



2023



2024



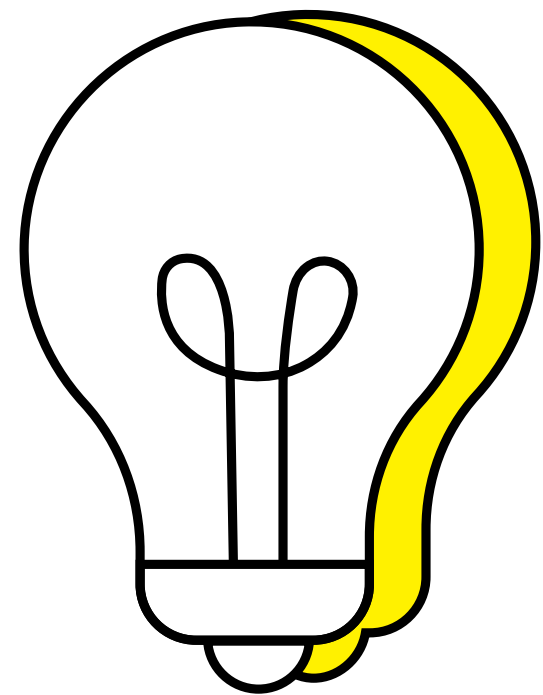
2025



2026



2028 (?)



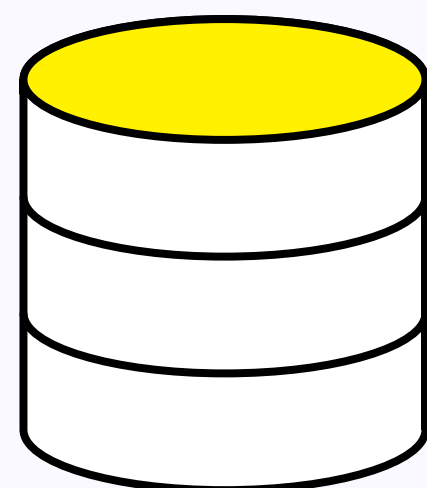
INSIGHT #3

“Big Data” is overrated

2011

Data processing landscape

Lots of data



Weak hardware



Inefficient software

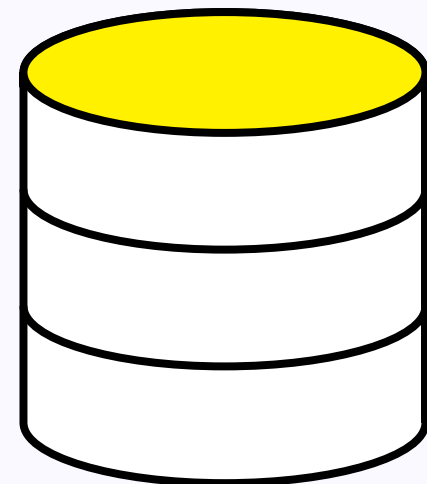
```
0[|||||100.0%] 1[      0.0%]  
2[      0.0%] 3[      0.0%]  
4[      0.0%] 5[      0.0%]  
6[      0.0%] 7[      0.0%]  
Mem[|||||      362M/27.4G]
```

2026

Data processing landscape

Lots of data

but most queries
only access a small amount



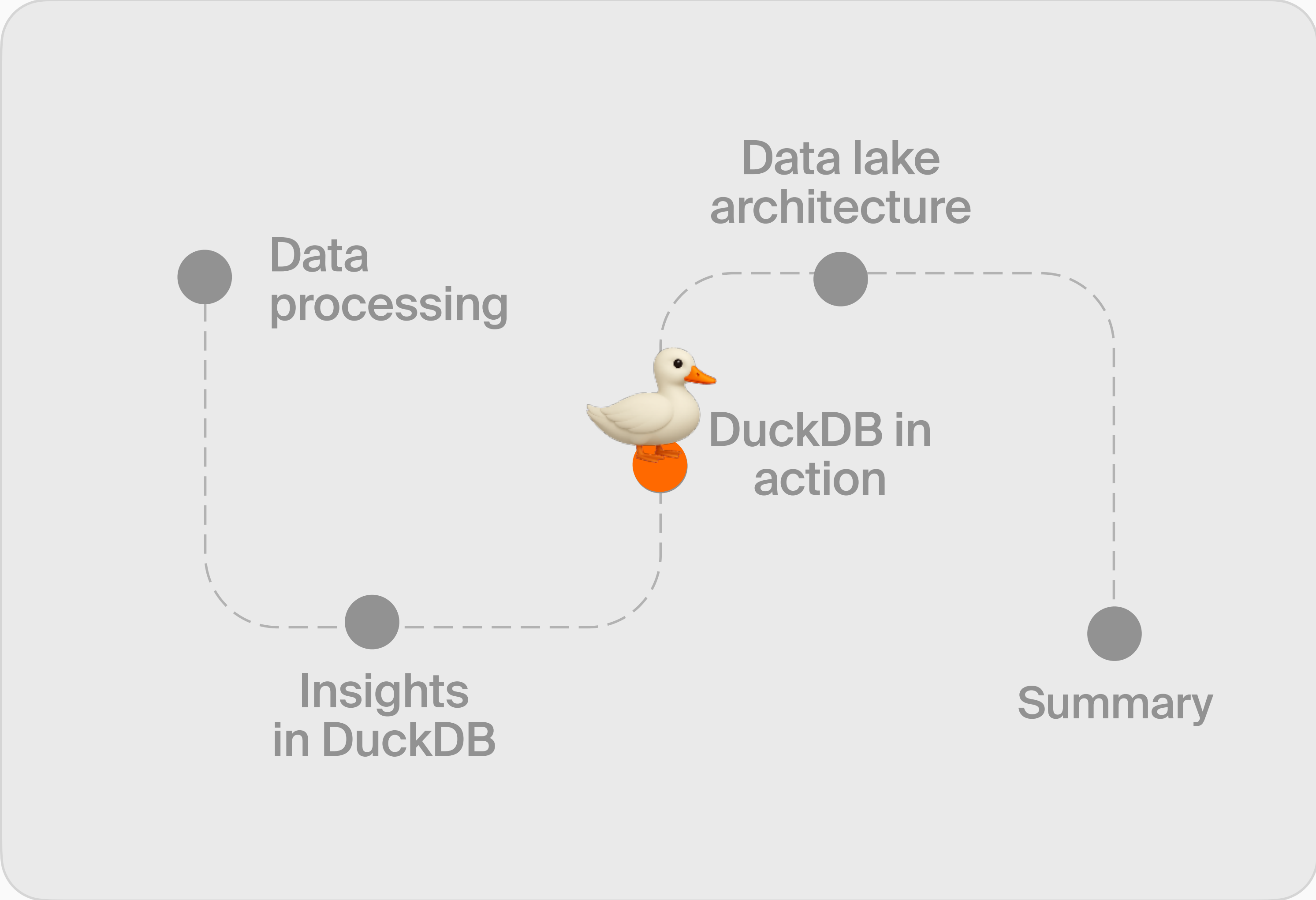
Strong hardware

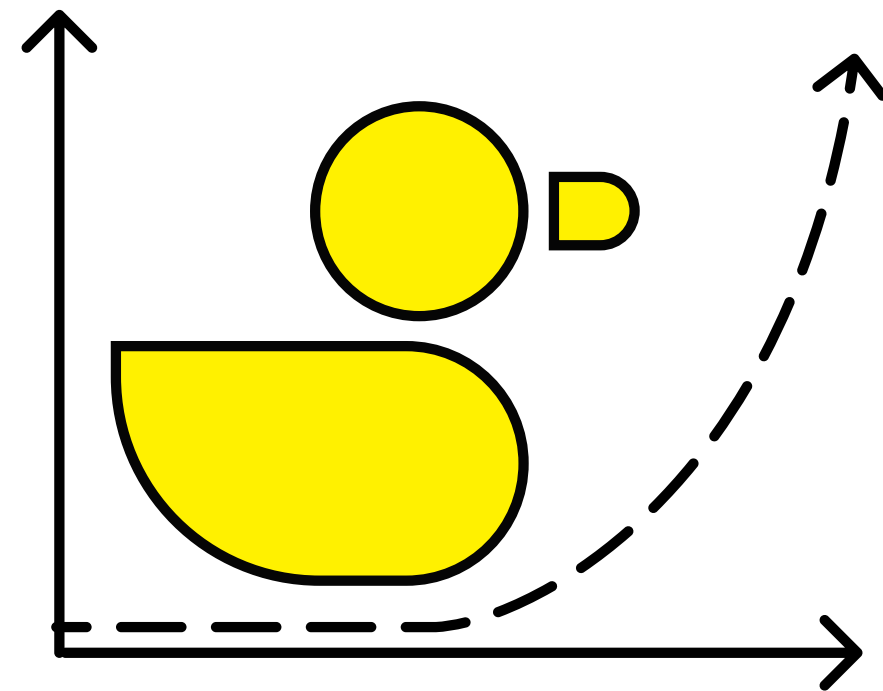


New techniques

columnar storage,
lightweight compression,
vectorized query execution







Performance & portability

Analytical queries

Analytical queries

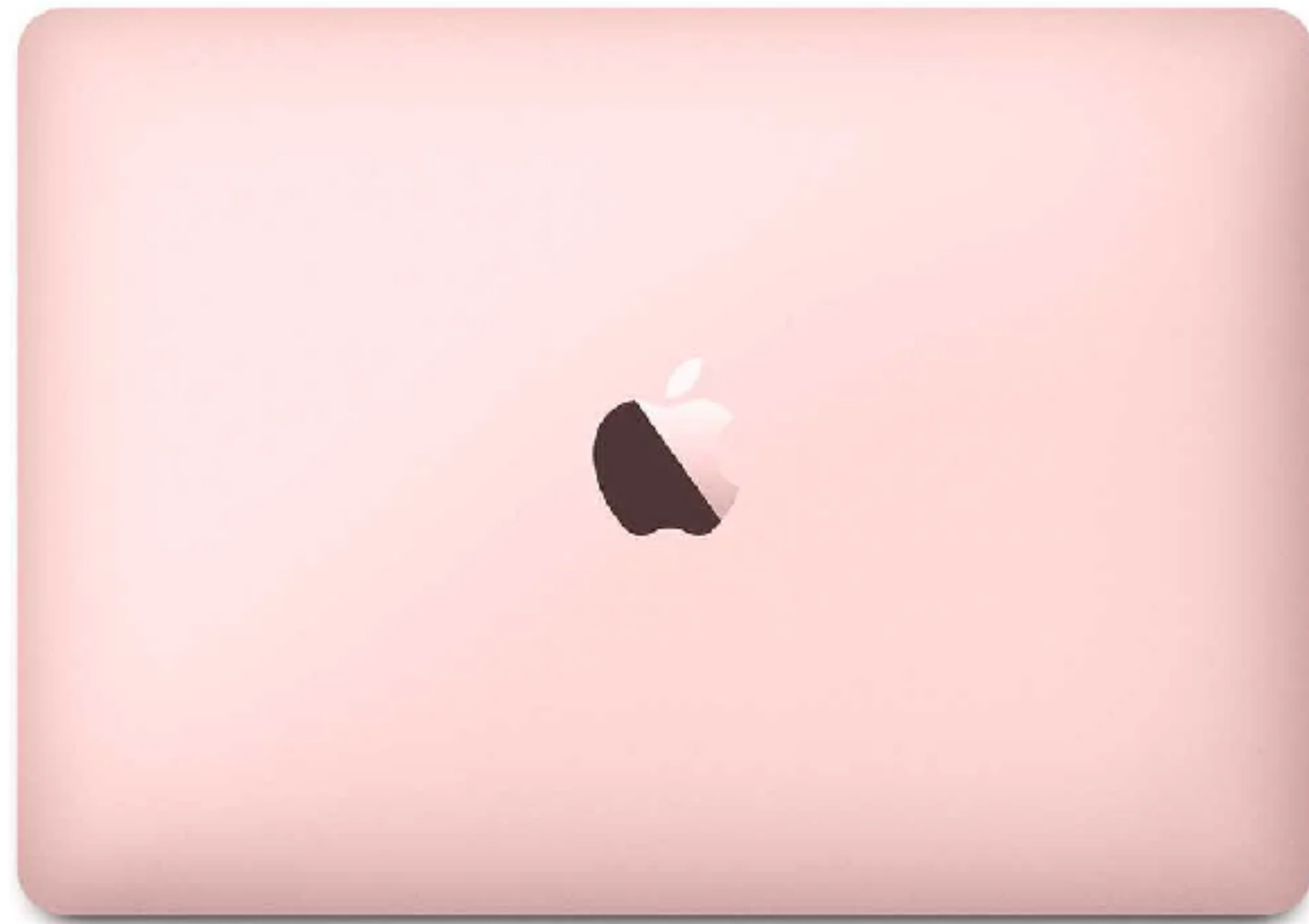


100 GB

Analytical queries



100 GB

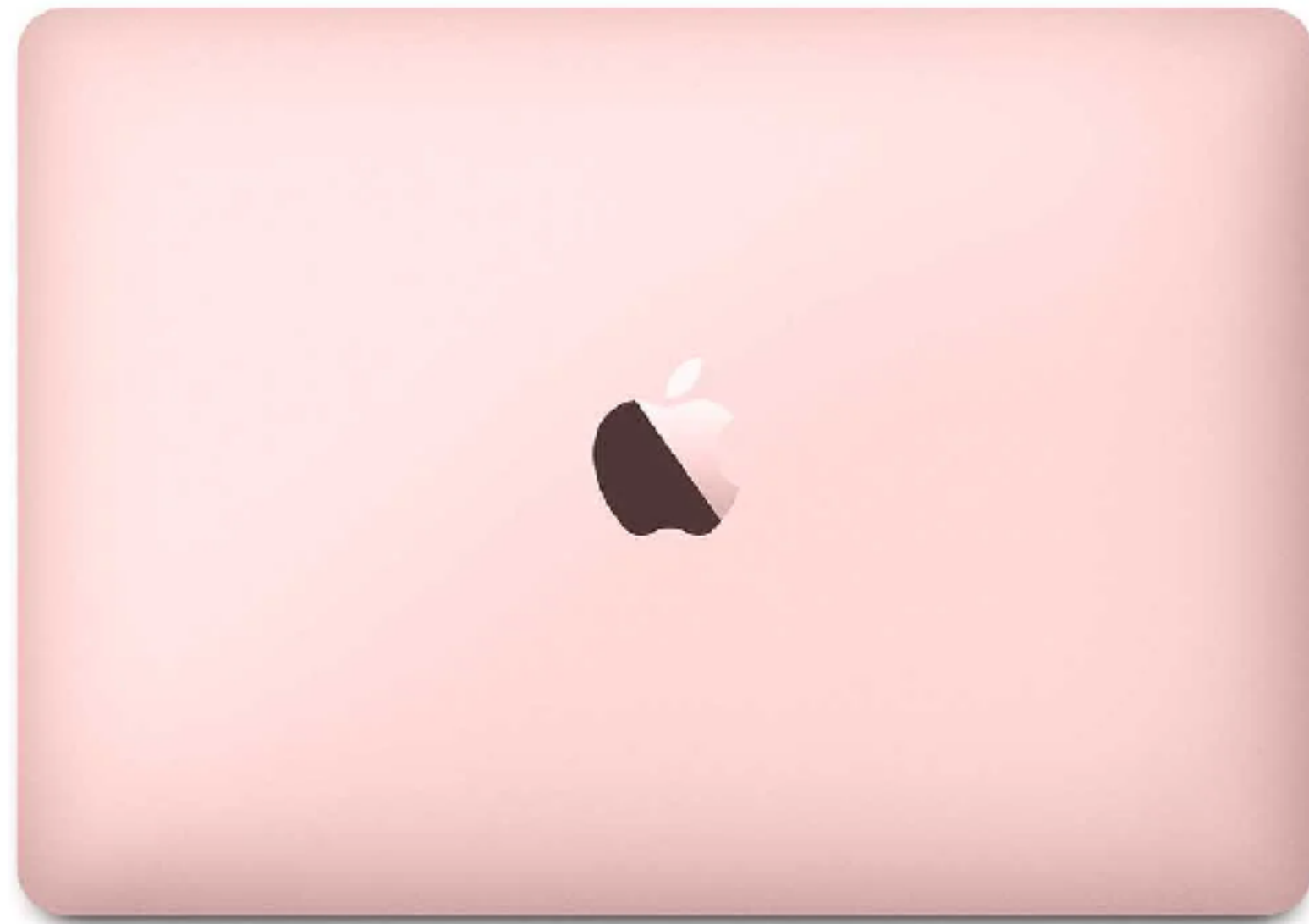


1 000 GB

Analytical queries



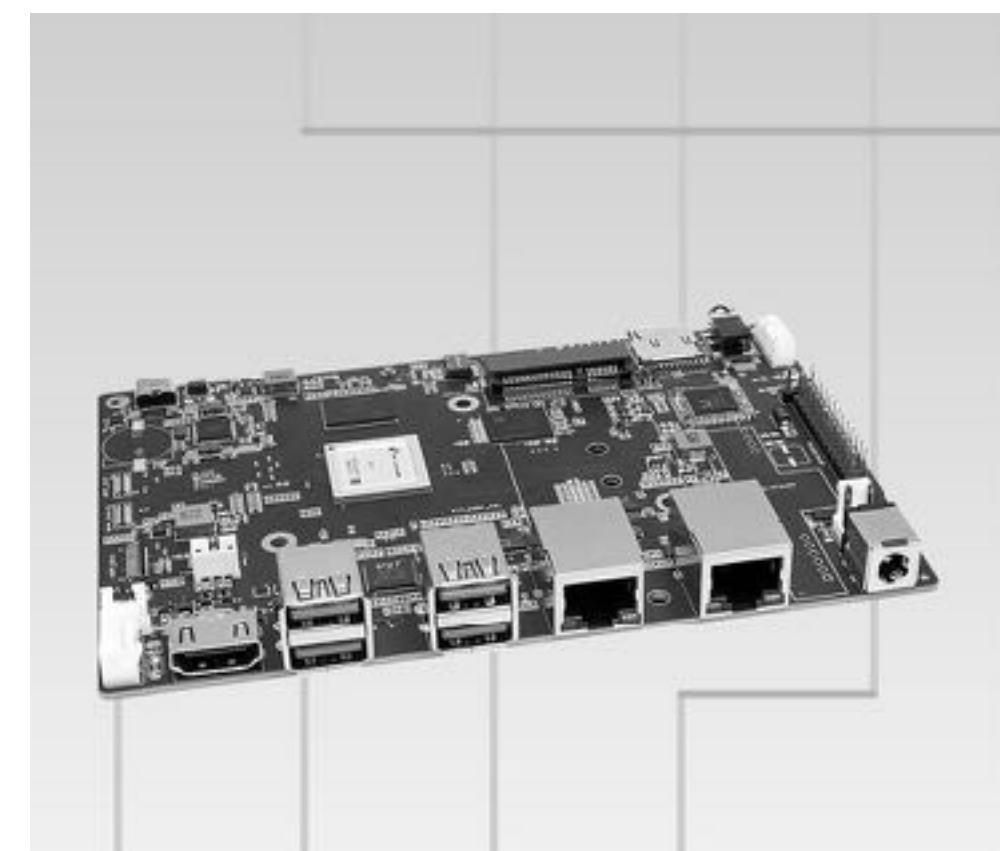
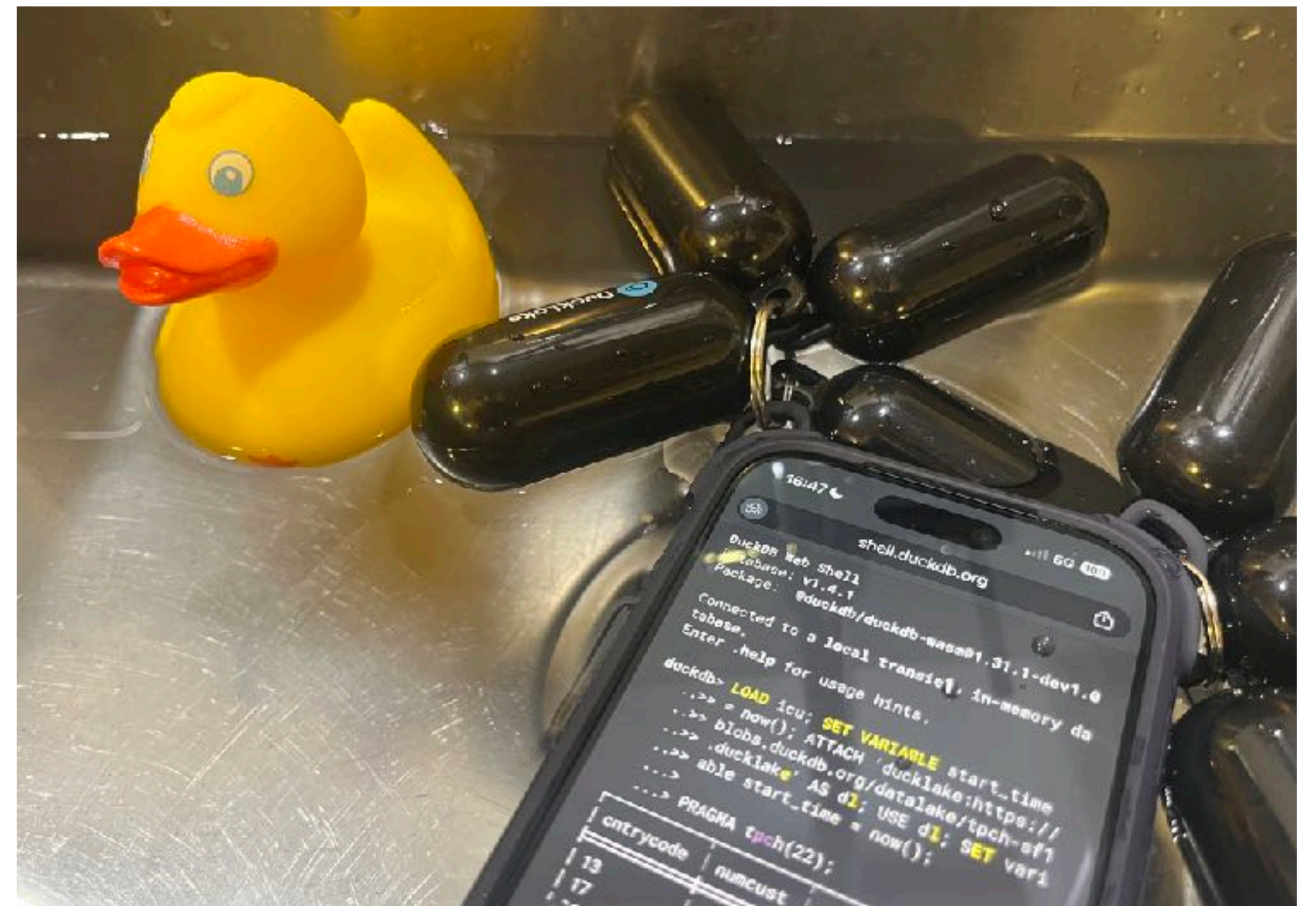
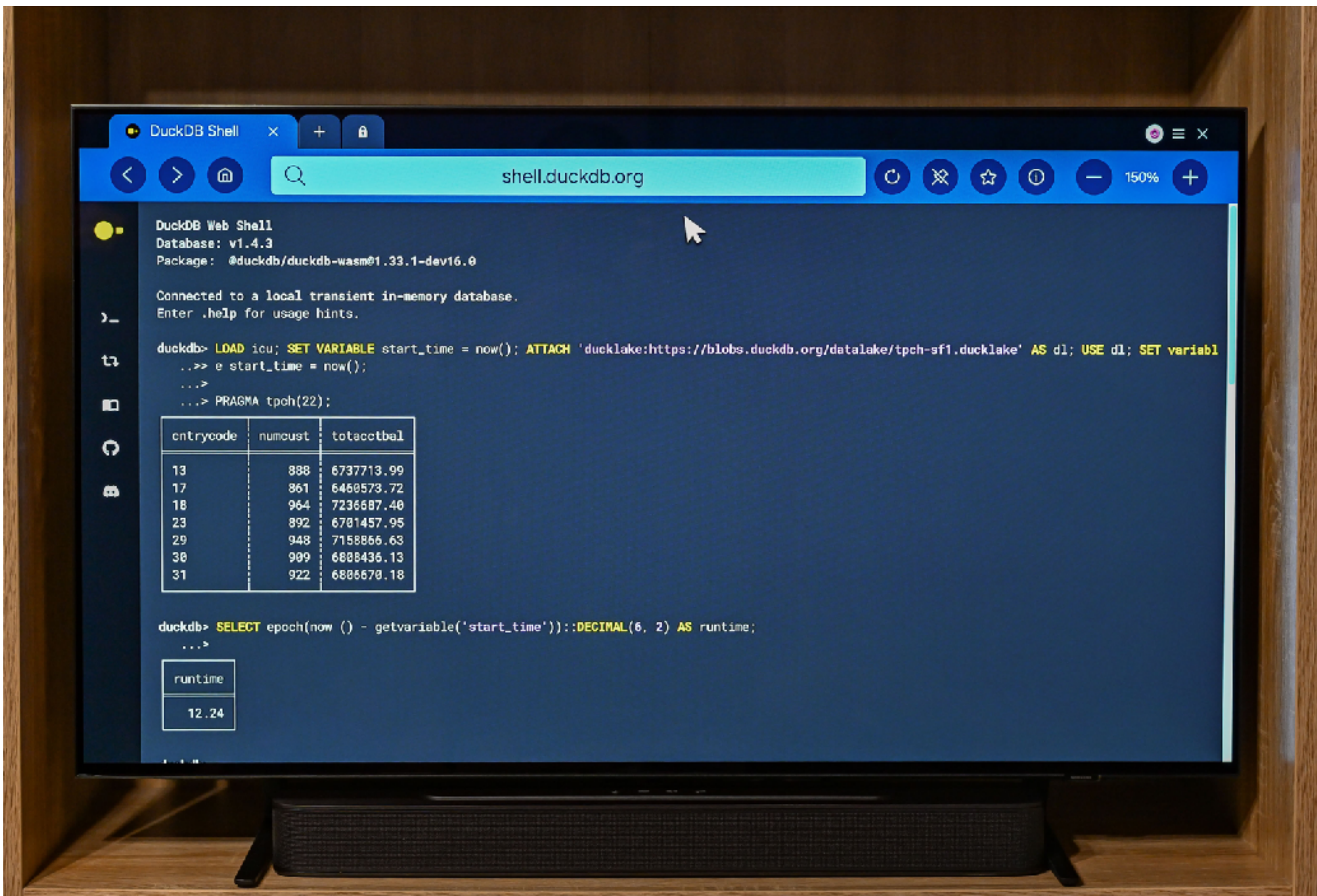
100 GB



1000 GB



10 000 GB





Python integrations

Pandas integration

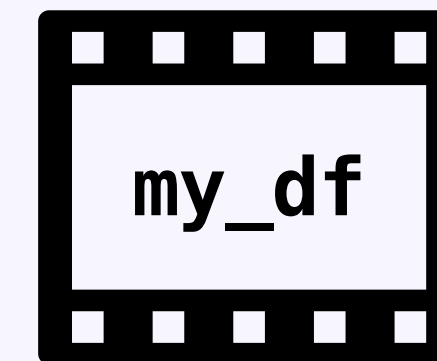
```
import duckdb
import pandas as pd

my_df = pd.DataFrame \
    .from_dict({'a': [42, 43]})

res = duckdb \
    .sql("SELECT avg(a) FROM my_df")

res.df()

# avg(a)
# 0 42.5
```



Pandas integration

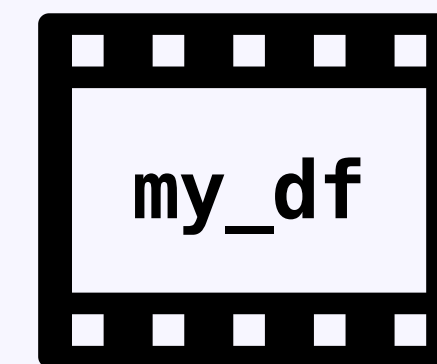
```
import duckdb
import pandas as pd

my_df = pd.DataFrame \
    .from_dict({'a': [42, 43]})

res = duckdb \
    .sql("SELECT avg(a) FROM my_df")

res.df()

#   avg(a)
# 0    42.5
```



 pandas

Pandas integration

```
import duckdb
import pandas as pd

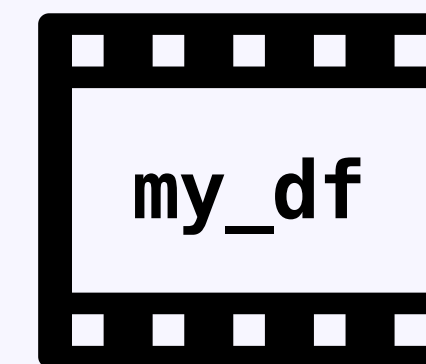
my_df = pd.DataFrame \
    .from_dict({'a': [42, 43]})

res = duckdb \
    .sql("SELECT avg(a) FROM my_df")

res.df()

# avg(a)
# 0 42.5
```

Replacement scan



 pandas

Pandas integration

```
import duckdb
import pandas as pd

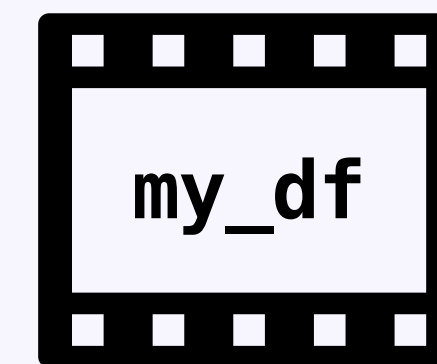
my_df = pd.DataFrame \
    .from_dict({'a': [42, 43]})

res = duckdb \
    .sql("SELECT avg(a) FROM my_df")

res.df()

#  avg(a)
#  0    42.5
```

Replacement scan



 pandas

Pandas integration

```
import duckdb
import pandas as pd

my_df = pd.DataFrame \
    .from_dict({'a': [42, 43]})

res = duckdb \
    .sql("SELECT avg(a) FROM my_df")

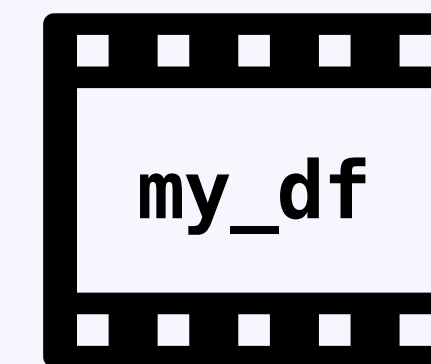
res.df()

# avg(a)
# 0 42.5
```

Replacement scan



Zero-copy access



 pandas

Polars integration

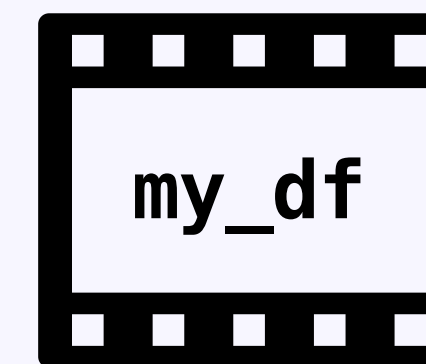
```
import duckdb
import polars as pl

my_df = pl.DataFrame({'a': [42, 43]})
res = duckdb.sql("SELECT avg(a) FROM my_df")
res.pl()
```

```
# shape: (1, 1)
```

```
#
# avg(a)
# ---
# f64
#
# 42.5
#
```

avg(a)
42.5



Polars

Polars integration

```
import duckdb
import polars as pl

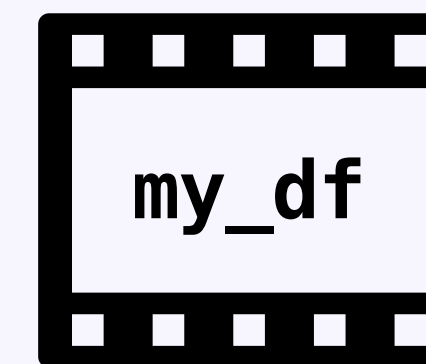
my_df = pl.DataFrame({'a': [42, 43]})
res = duckdb.sql("SELECT avg(a) FROM my_df")
res.pl()
```

```
# shape: (1, 1)
```

```
#
#  avg(a)
#  ---
#  f64
#
#  42.5
#
```



Using
Arrow



Polars

Polars integration

```
import duckdb
import polars as pl

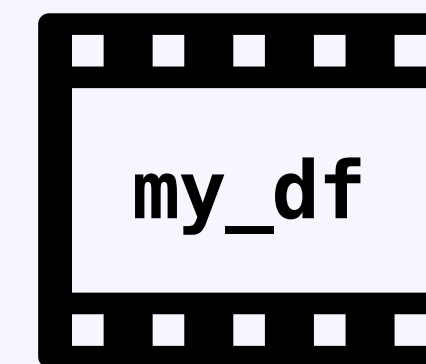
my_df = pl.DataFrame({'a': [42, 43]})
res = duckdb.sql("SELECT avg(a) FROM my_df")
res.pl()
```

```
# shape: (1, 1)
```

```
#
#   avg(a)
#   ---
#   f64
#
#   42.5
#
```



Using
Arrow



Polars

Numpy integration

```
import duckdb
import numpy as np

my_arr = duckdb \
    .sql("FROM (VALUES (23), (24)) t(a)") \
    .fetchnumpy()

np.sum(my_arr['a'])

# np.int64(47)
```



Numpy integration

```
import duckdb
import numpy as np

my_arr = duckdb \
    .sql("FROM (VALUES (23), (24)) t(a)") \
    .fetchnumpy()
```

```
np.sum(my_arr['a'])
```

```
# np.int64(47)
```



Support user-defined functions (UDFs)

```
import duckdb
from faker import Faker

def get_random_name():
    return Faker().name()

duckdb.create_function(
    "random_name", get_random_name, [],
    duckdb.typing.VARCHAR
)

duckdb.sql("SELECT random_name()") \
    .fetchall()
```



SELECT random_name()



get_random_name()

Support user-defined functions (UDFs)

```
import duckdb
from faker import Faker

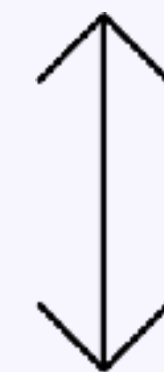
def get_random_name():
    return Faker().name()

duckdb.create_function(
    "random_name", get_random_name, [],
    duckdb.typing.VARCHAR
)

duckdb.sql("SELECT random_name()") \
    .fetchall()
```

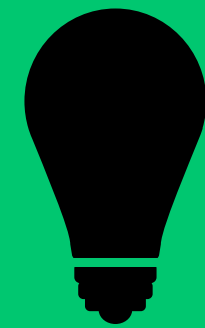


SELECT random_name()



get_random_name()

Support user-defined functions (UDFs)



Roadmap: harness no-GIL Python to parallelize

```
import duckdb
from faker import Faker

def get_random_name():
    return Faker().name()

duckdb.create_function(
    "random_name", get_random_name, [],
    duckdb.typing.VARCHAR
)

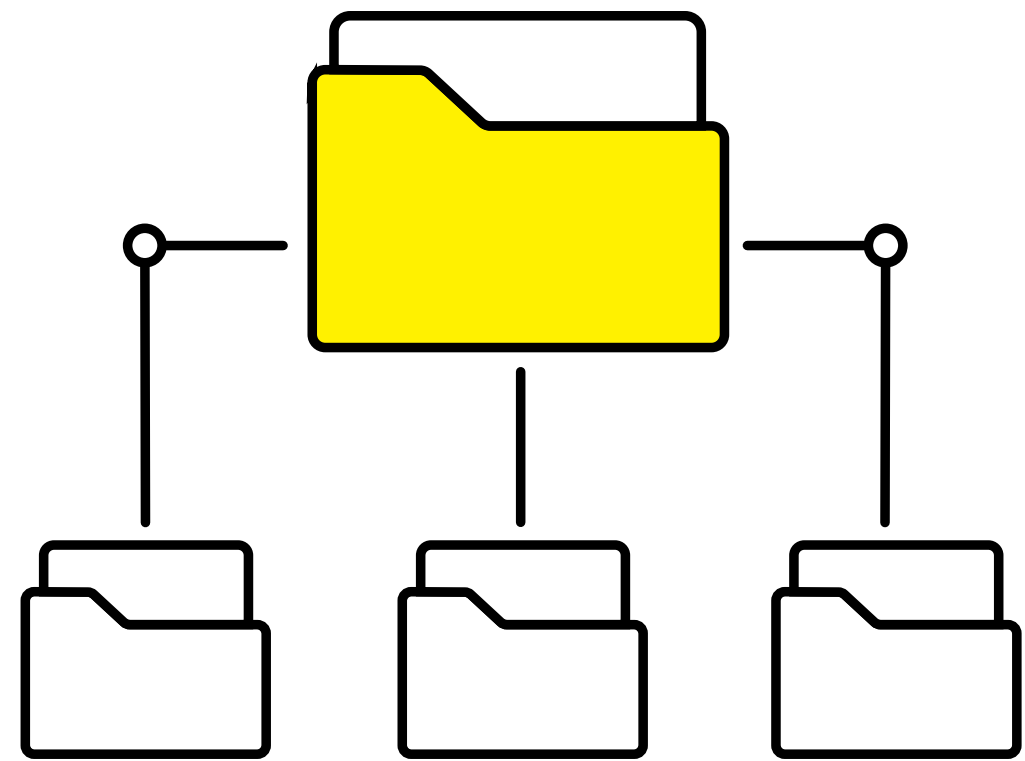
duckdb.sql("SELECT random_name()") \
    .fetchall()
```



SELECT random_name()

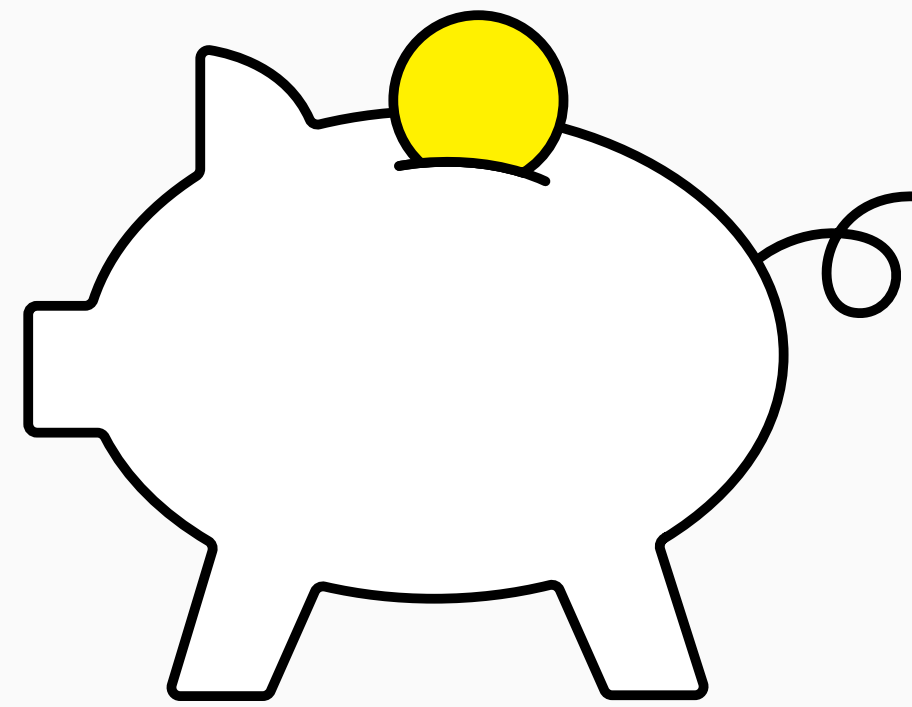


get_random_name()



Use cases

Cost saving

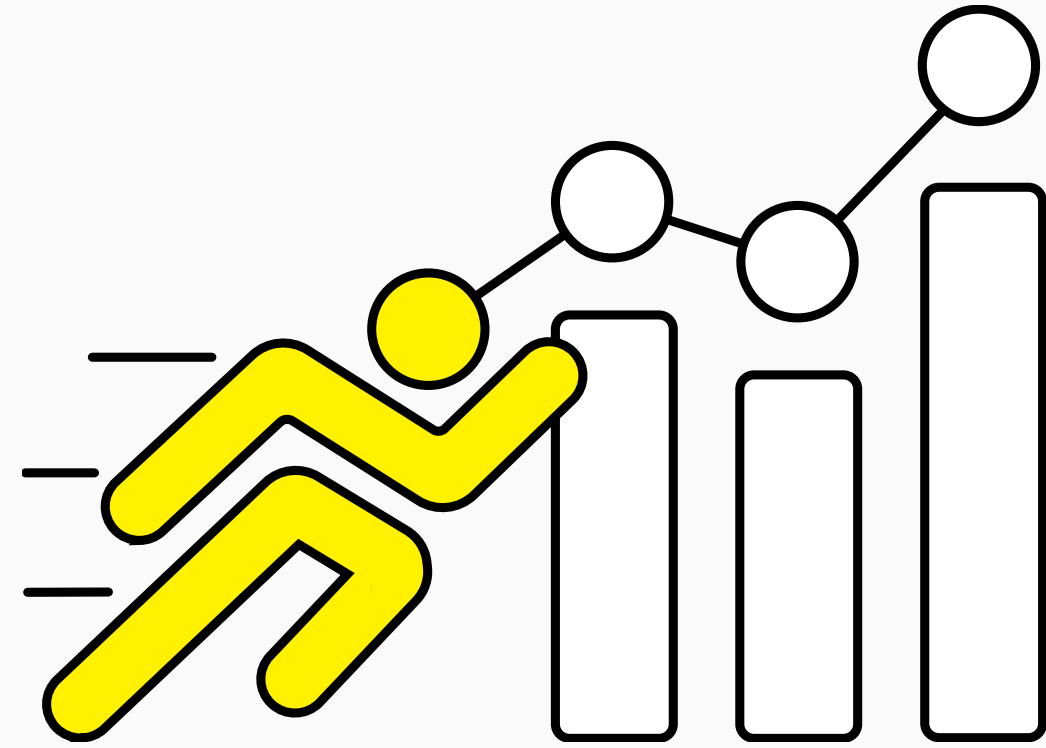


**Fast local
development:
time & egress fees**

**Replace
proprietary systems:
license cost**

**Replace distributed
systems:
hardware cost**

Last mile analytics



1 PB log →
pre-process with Spark

Process the 200 GB
aggregate DuckDB

Build fast dashboards

Education

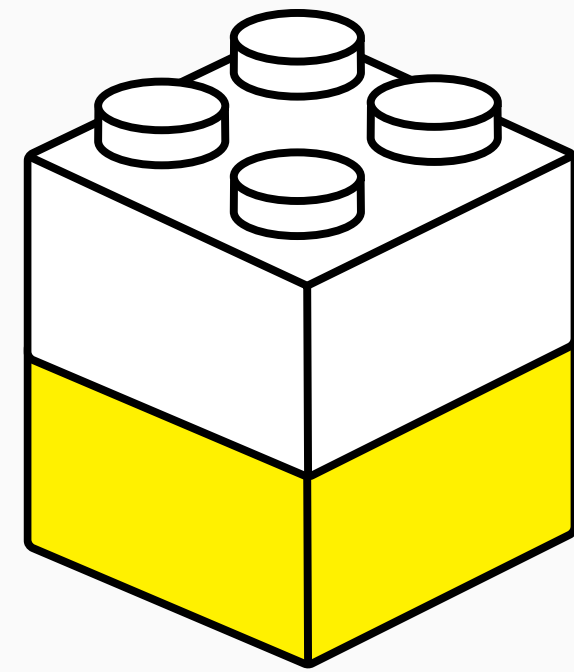


Easy to install

Open-source

No DBA needed

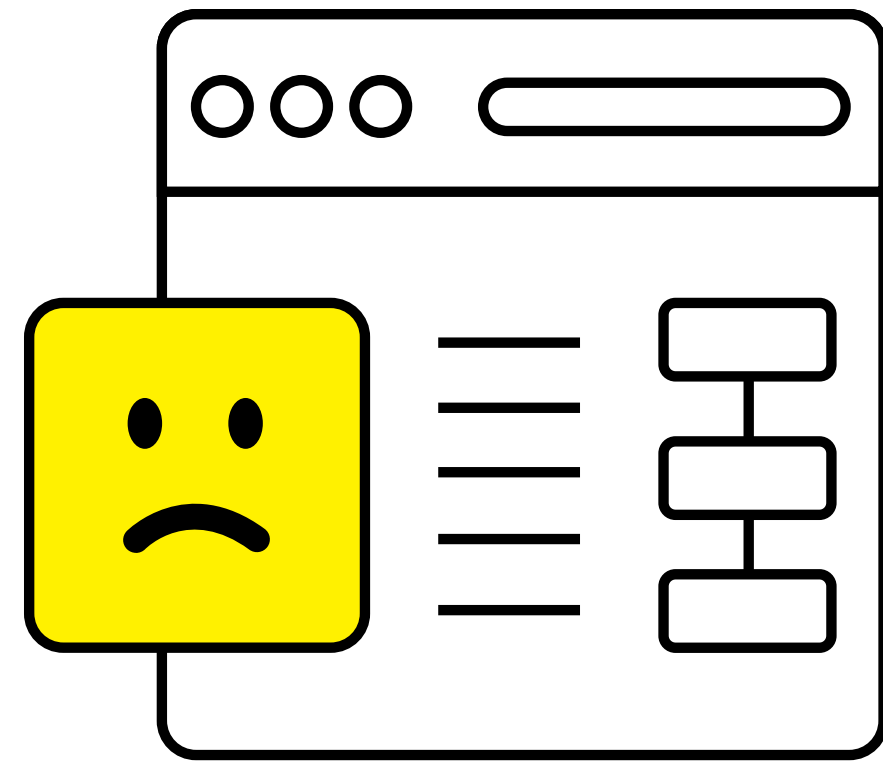
Modernizing codebases



Postgres-compatible
SQL dialect

Builds on open formats
(CSV, Parquet)

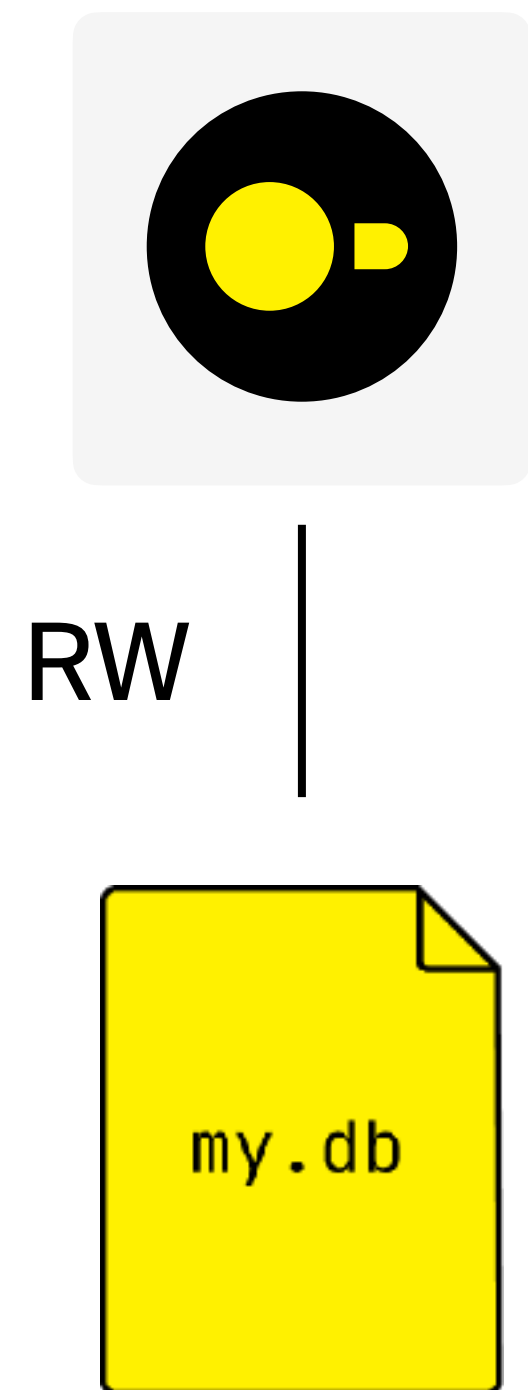
LLMs can refactor
code to use DuckDB



Limitations

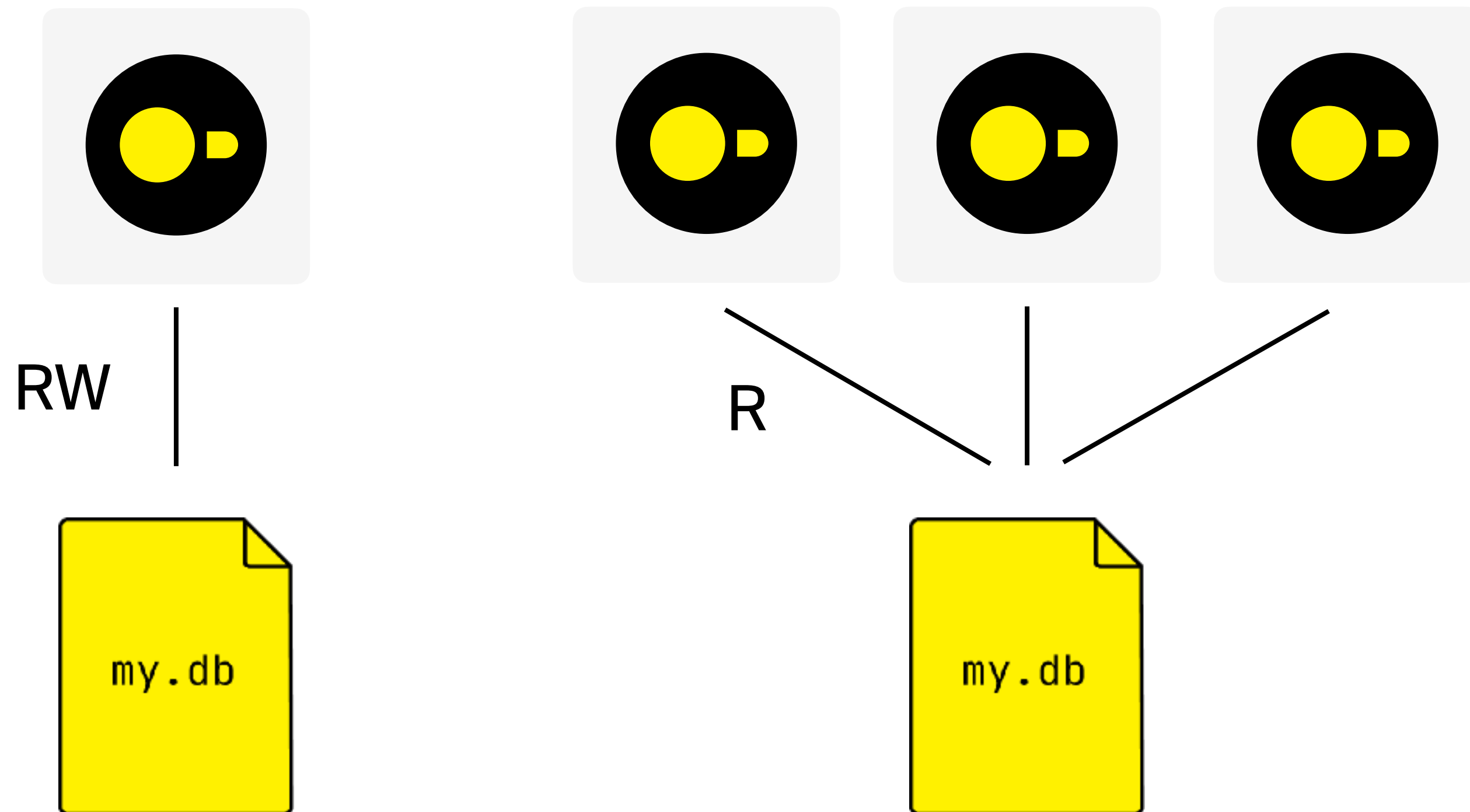
What are DuckDB's limitations?

Single-player mode: only a single process can attach in RW mode



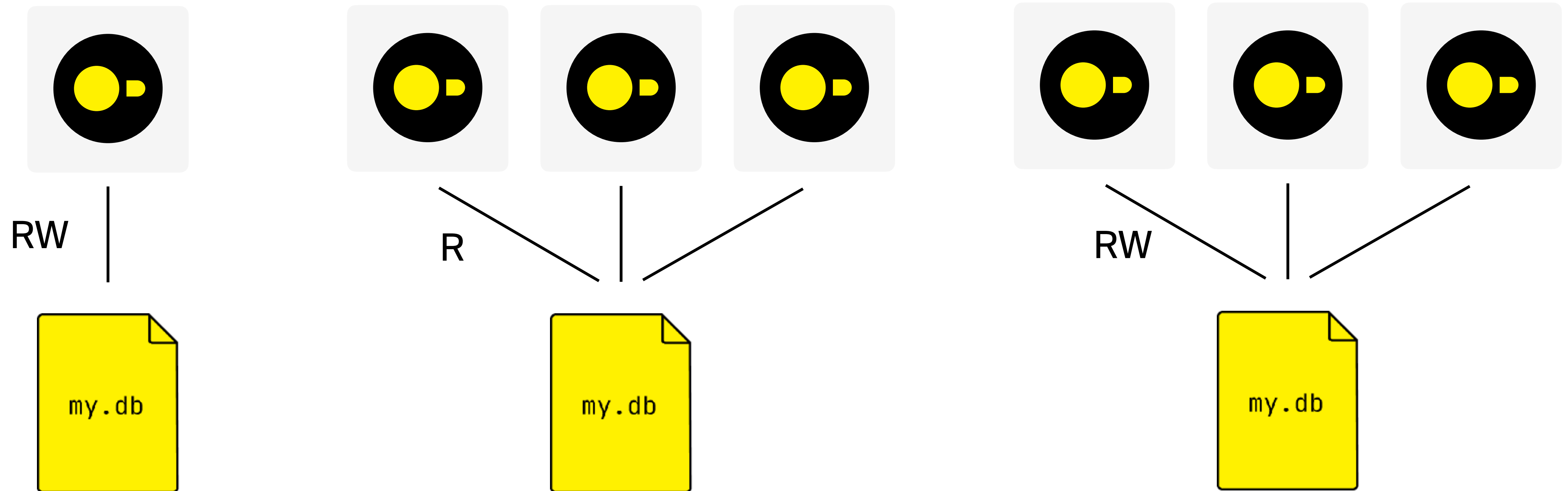
What are DuckDB's limitations?

Single-player mode: only a single process can attach in RW mode



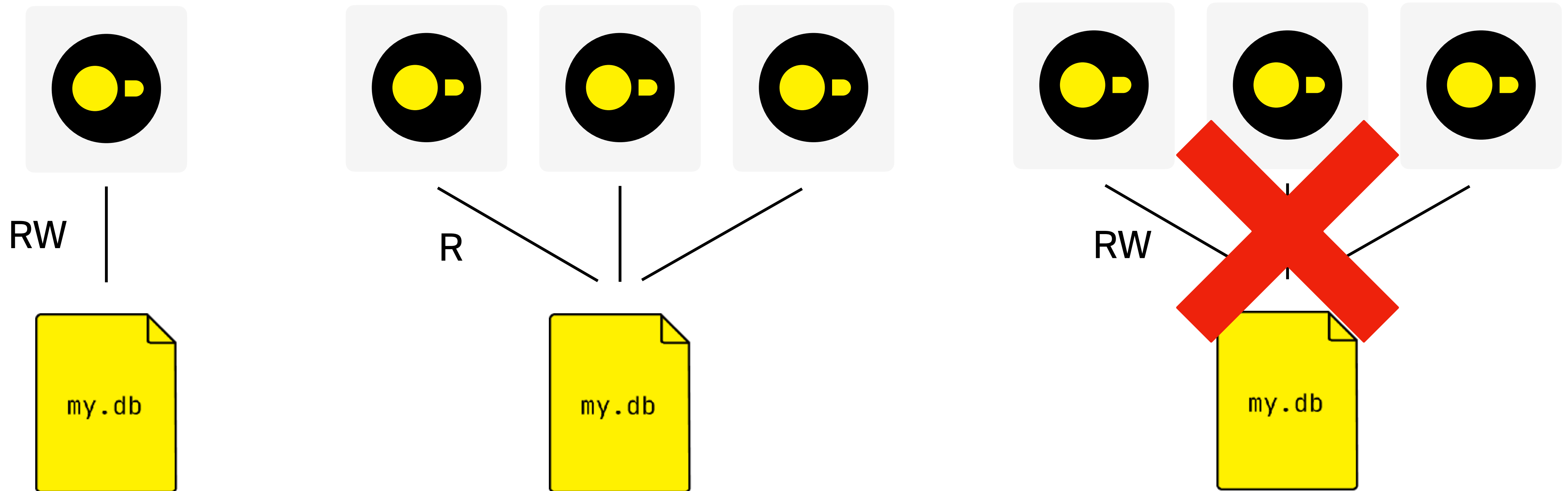
What are DuckDB's limitations?

Single-player mode: only a single process can attach in RW mode



What are DuckDB's limitations?

Single-player mode: only a single process can attach in RW mode



What are DuckDB's limitation?

10 TB CSV = 2.7 TB database file

**Optimized
for batch loading**

small writes are slow

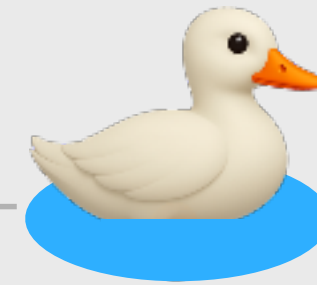


Huge database file

difficult to handle



Data lake architecture

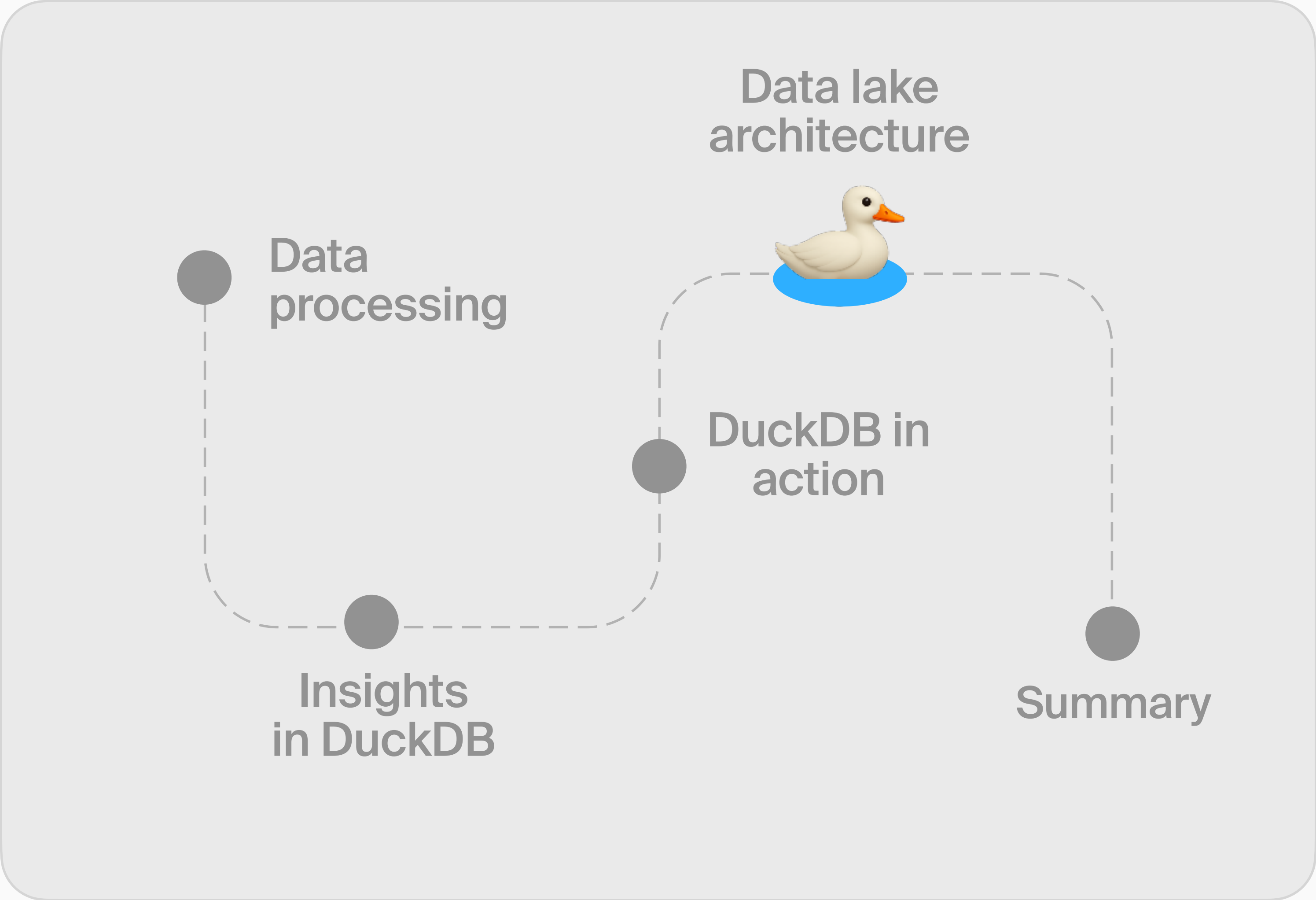


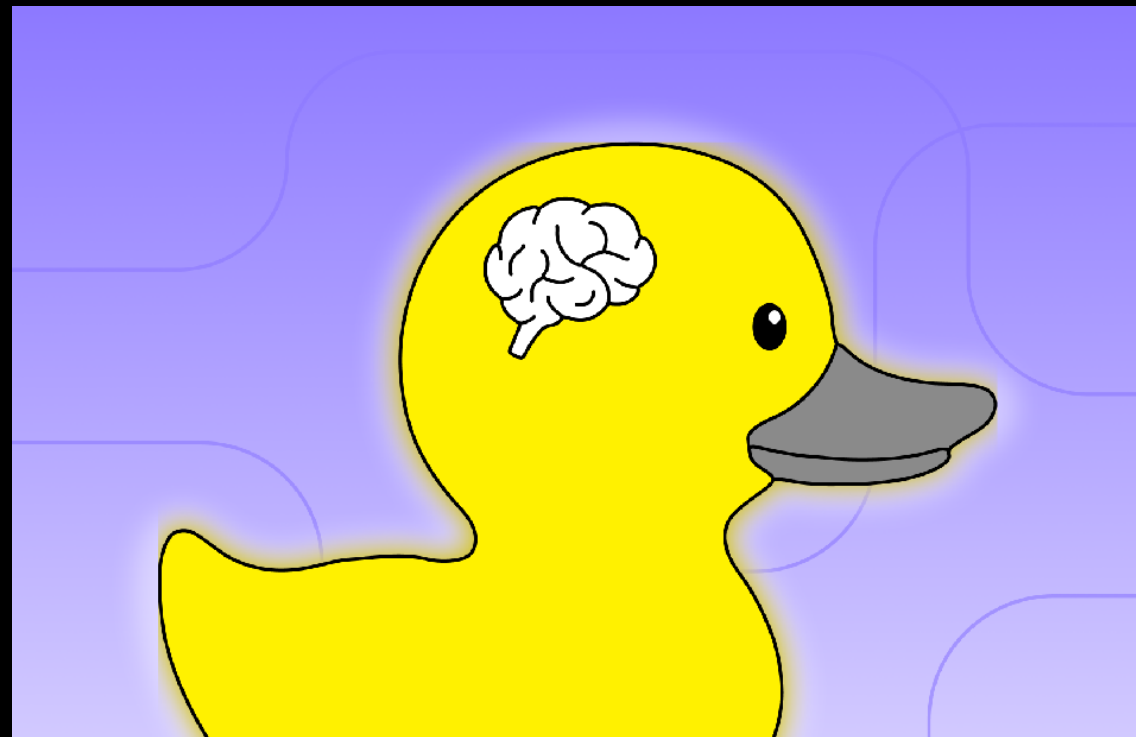
● Data processing

● DuckDB in action

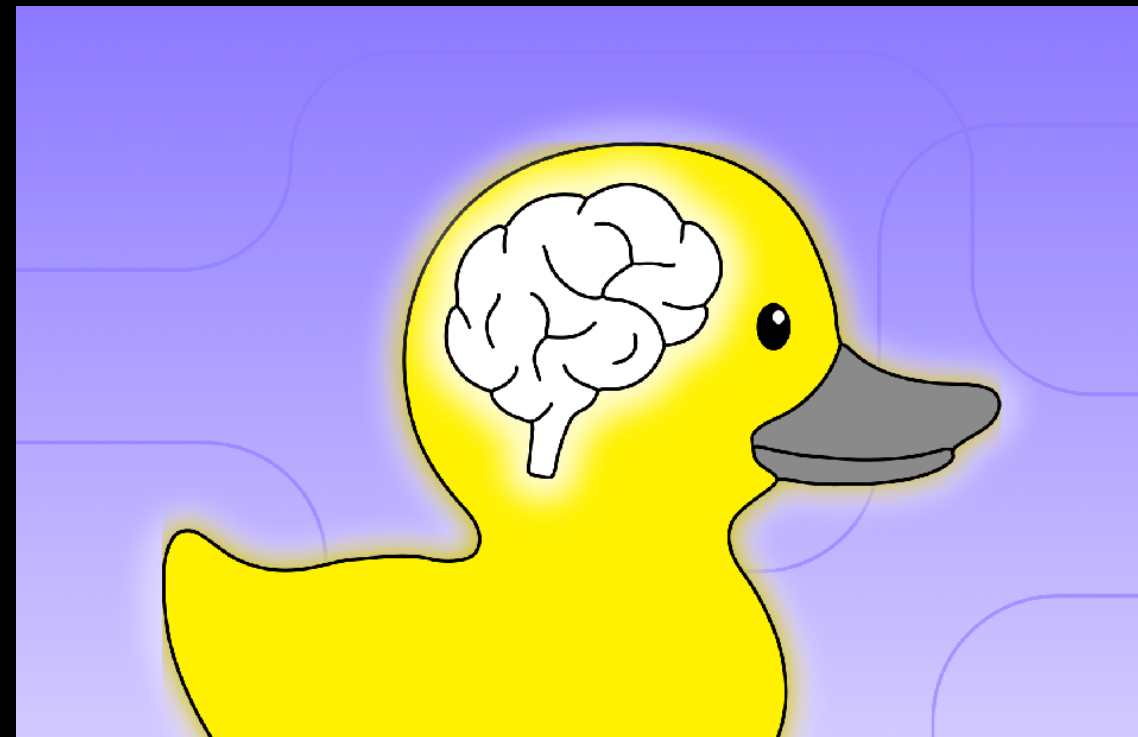
● Insights in DuckDB

● Summary

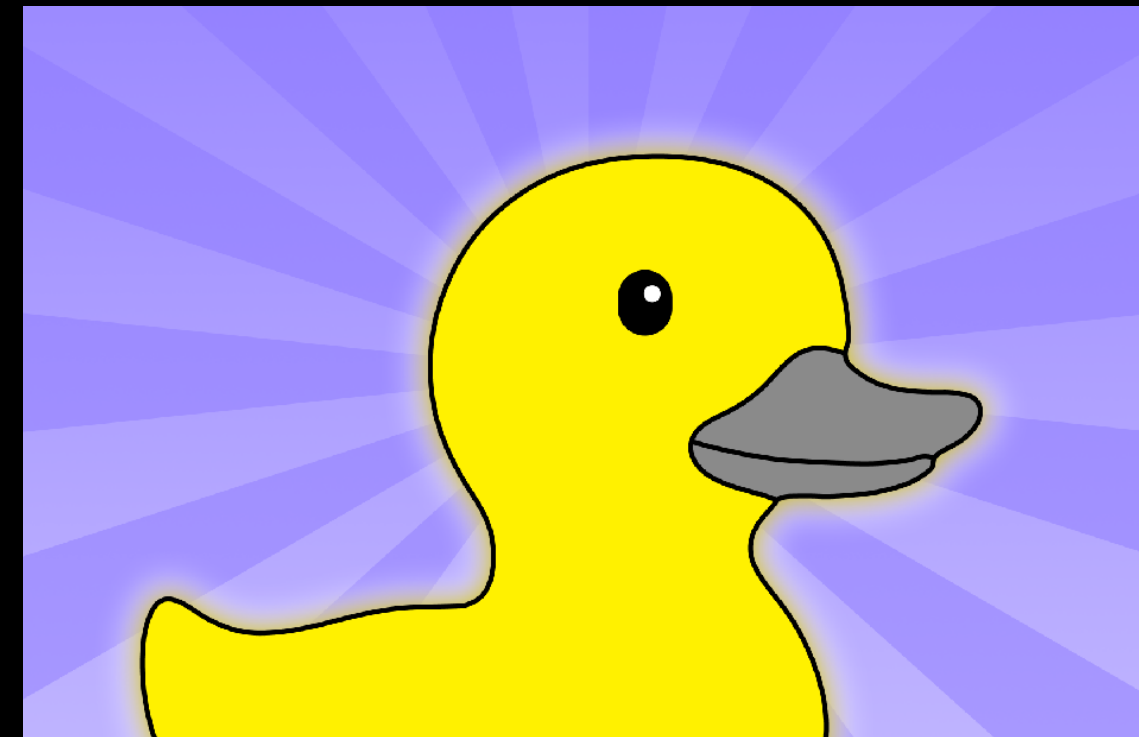




Text files



Binary files



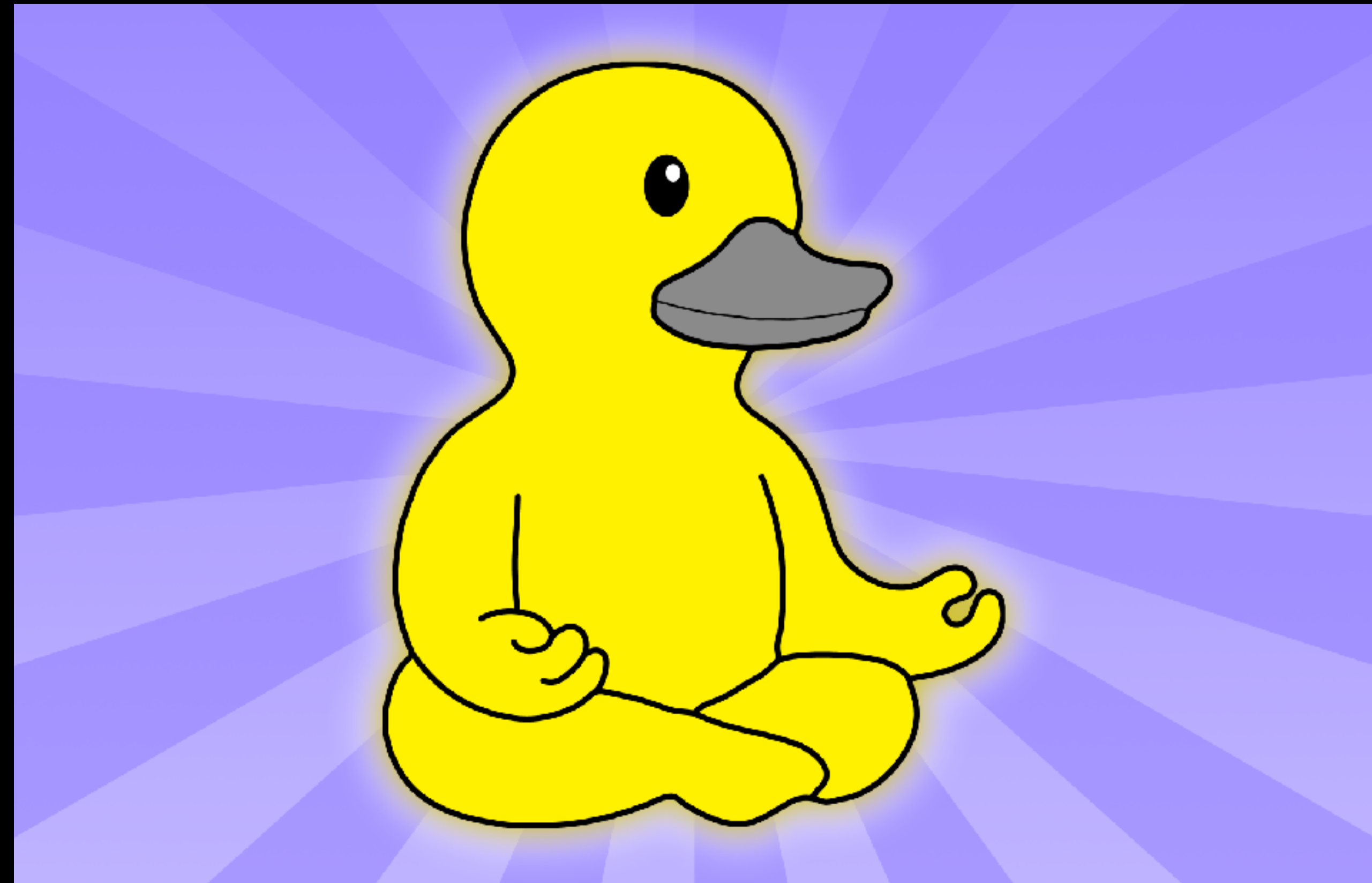
Database



Data lake

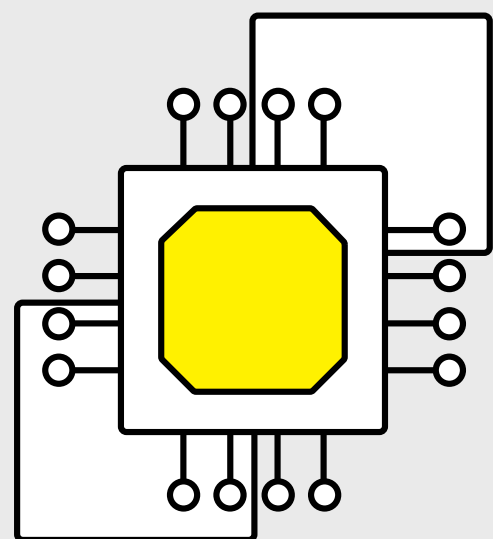
Data lake

*But first, some history
on cloud architecture!*



2005

Separation of compute and storage

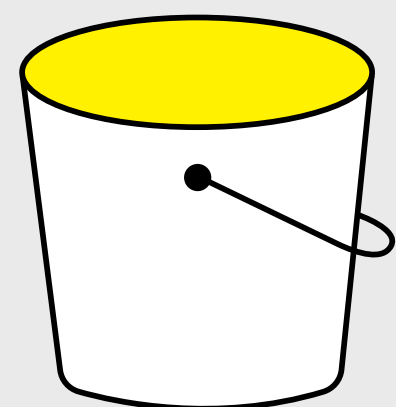
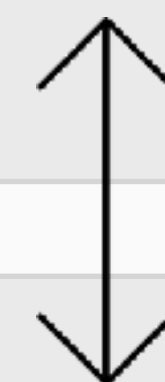
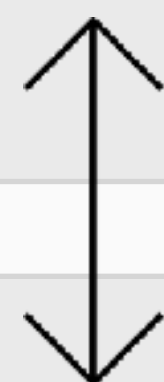


compute node

compute node

...

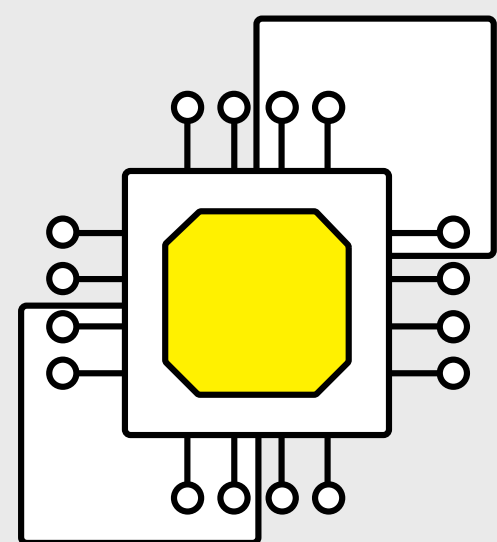
compute node



object storage ("blob storage")

2005

Separation of compute and storage

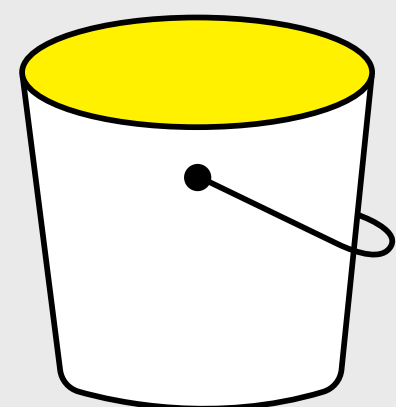
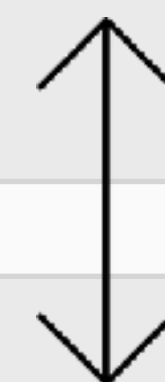
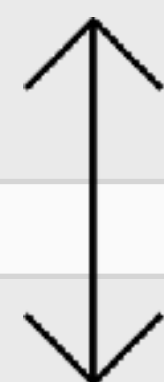


compute node

compute node

...

compute node

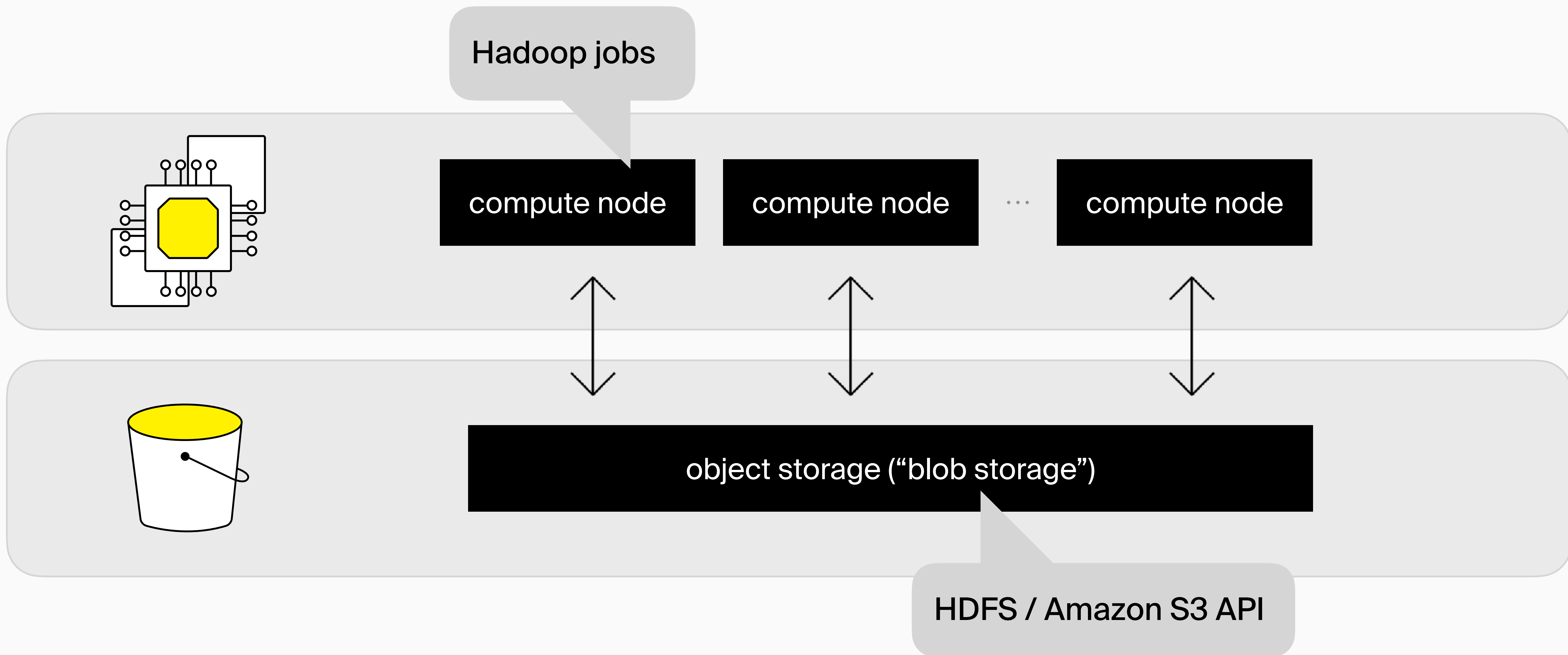


object storage ("blob storage")

HDFS / Amazon S3 API

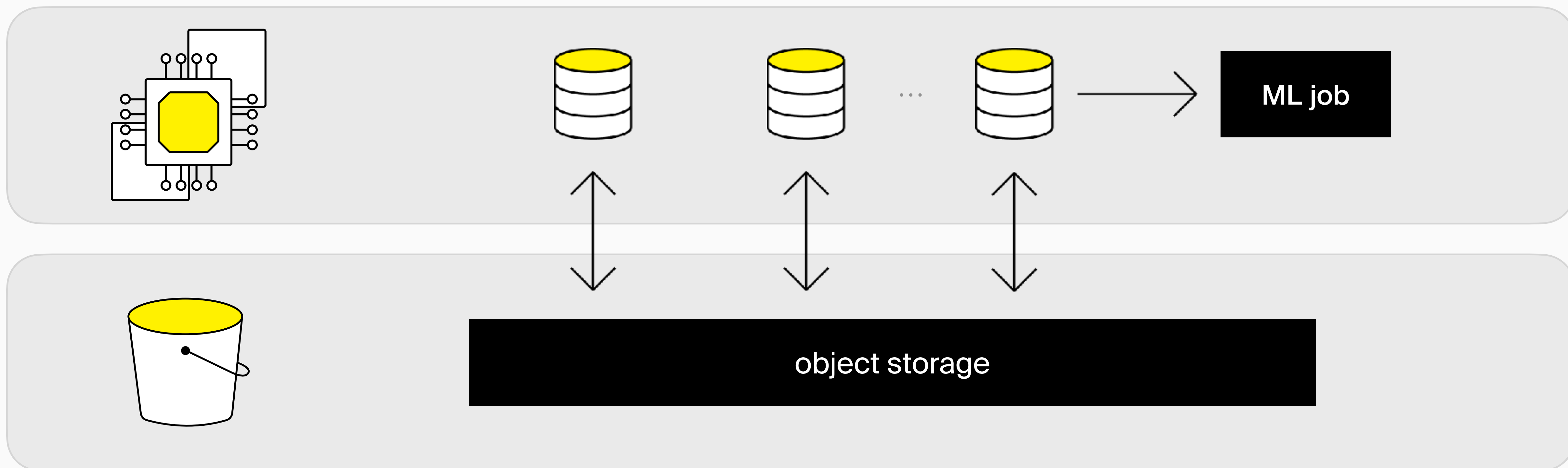
2005

Separation of compute and storage



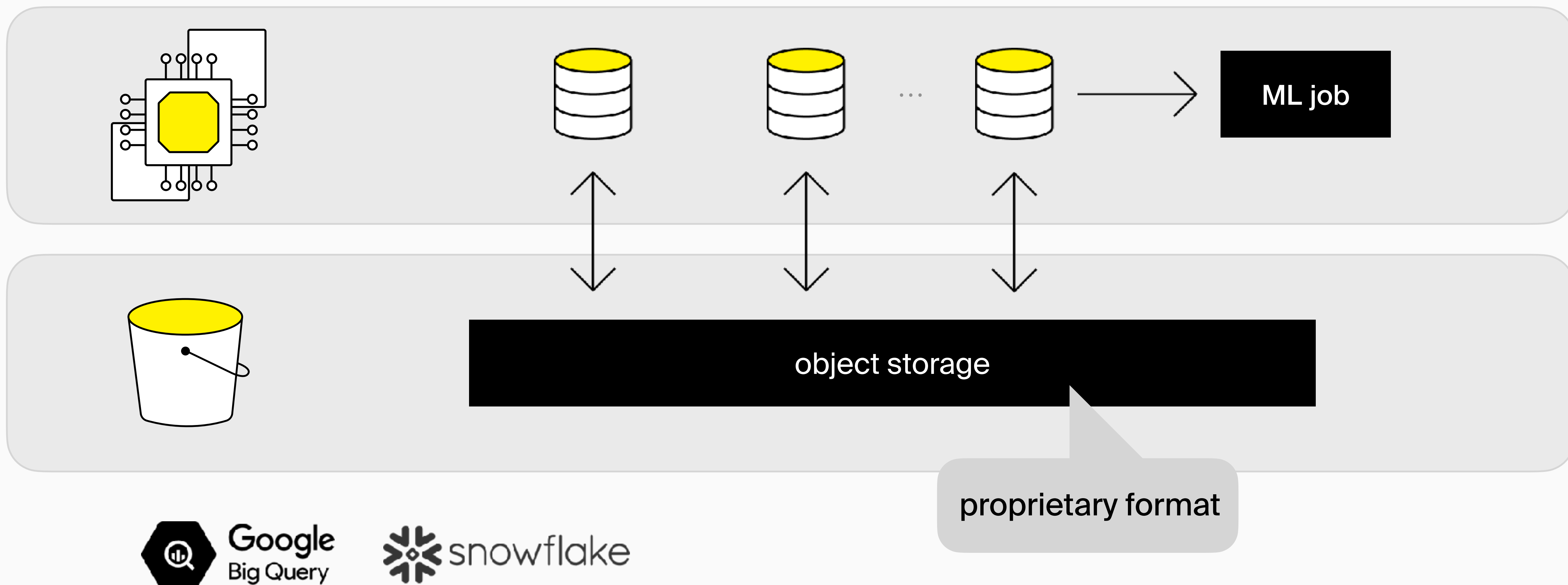
2010

Cloud data warehouses



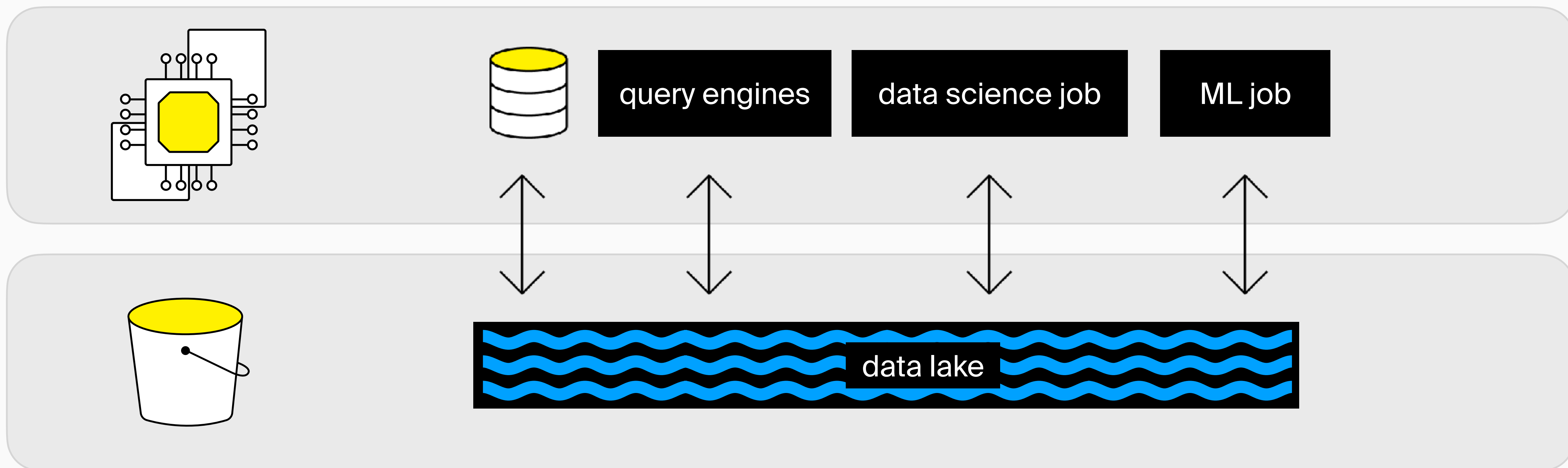
2010

Cloud data warehouses



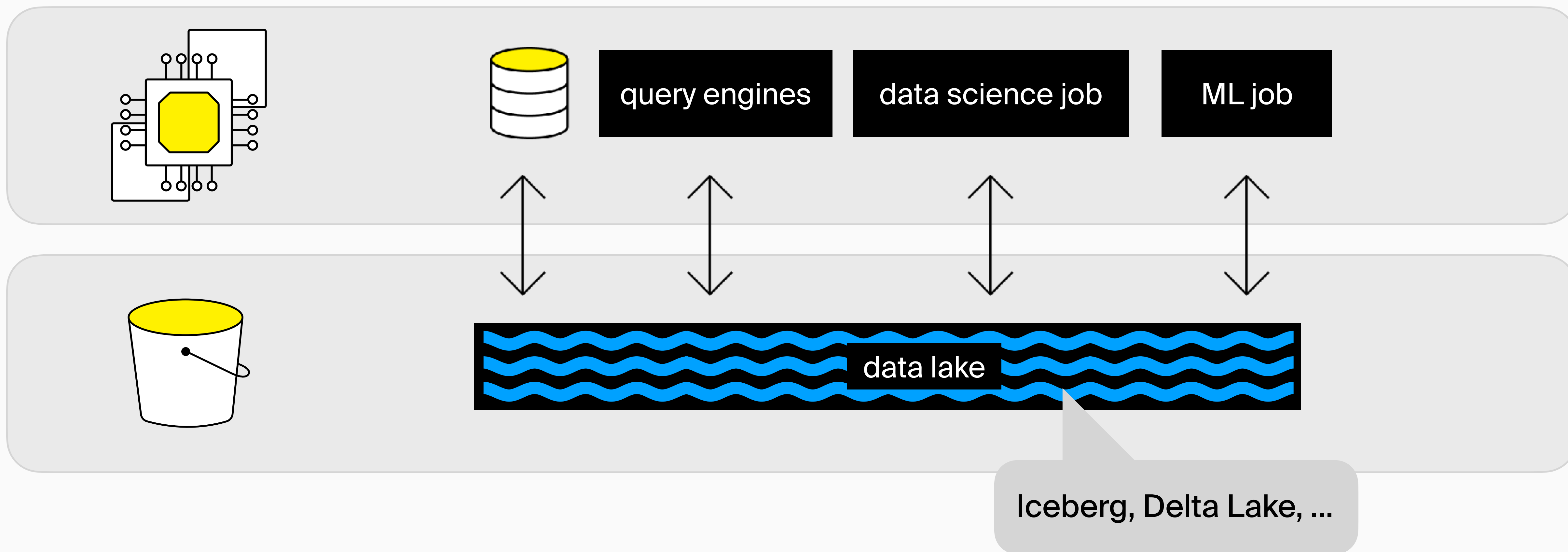
2015

Data lake architecture

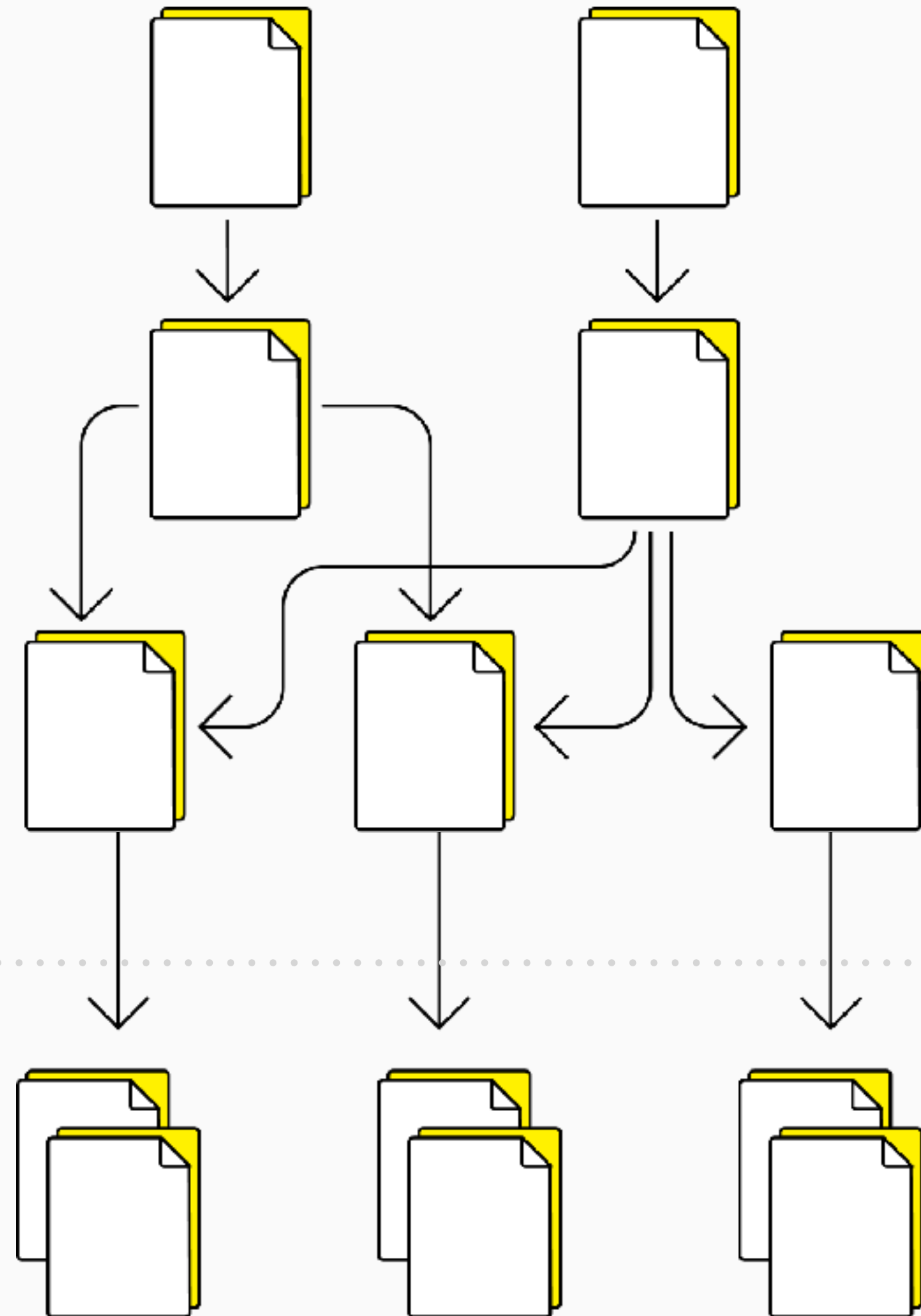


2015

Data lake architecture



metadata
(json / avro)
statistics



data
(parquet)

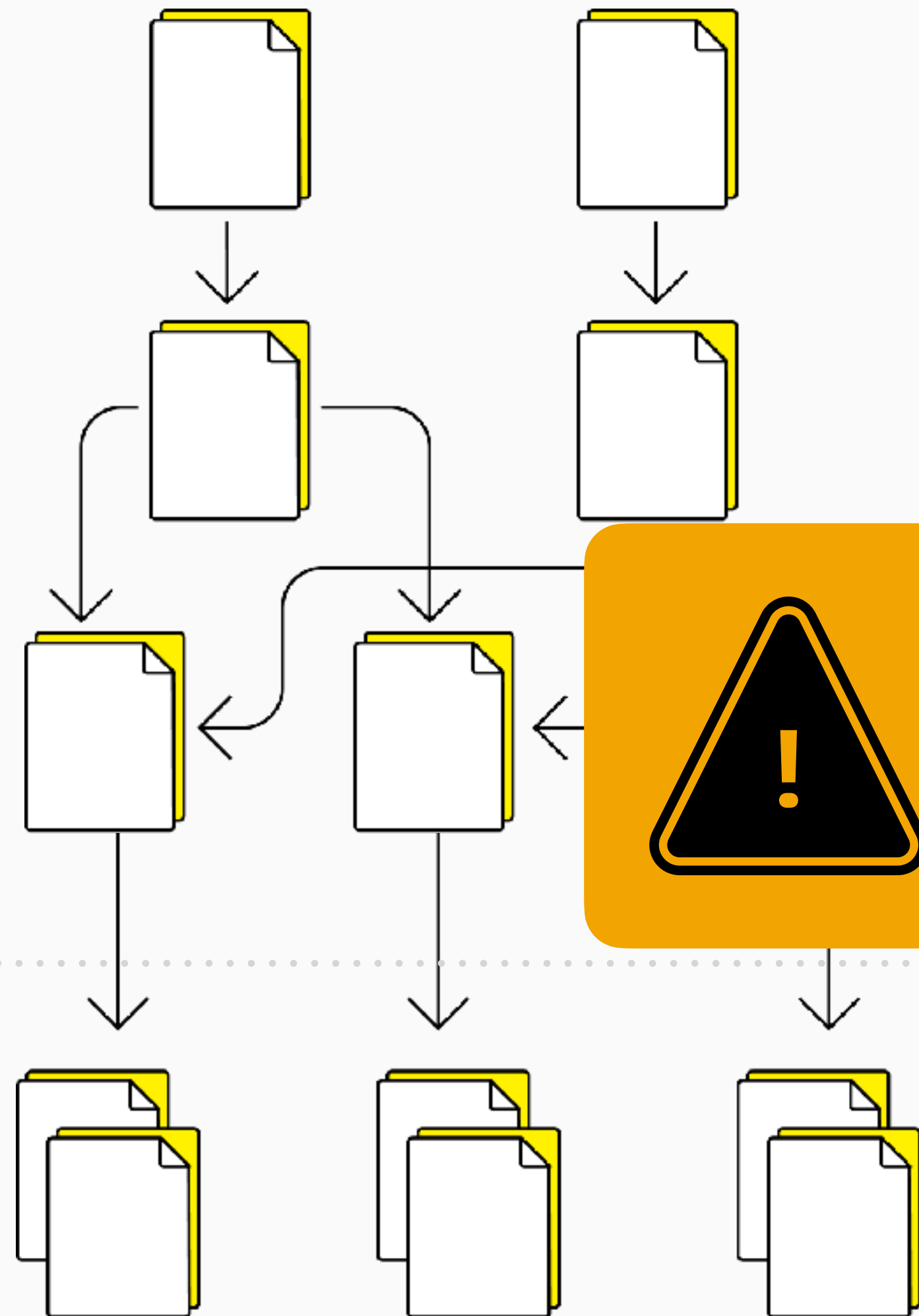


- File-only architecture
- Updatable table
- Schema evolution



- Performance problems
- “Small files” problem
- Only a single table

metadata
(json / avro)
statistics



data
(parquet)

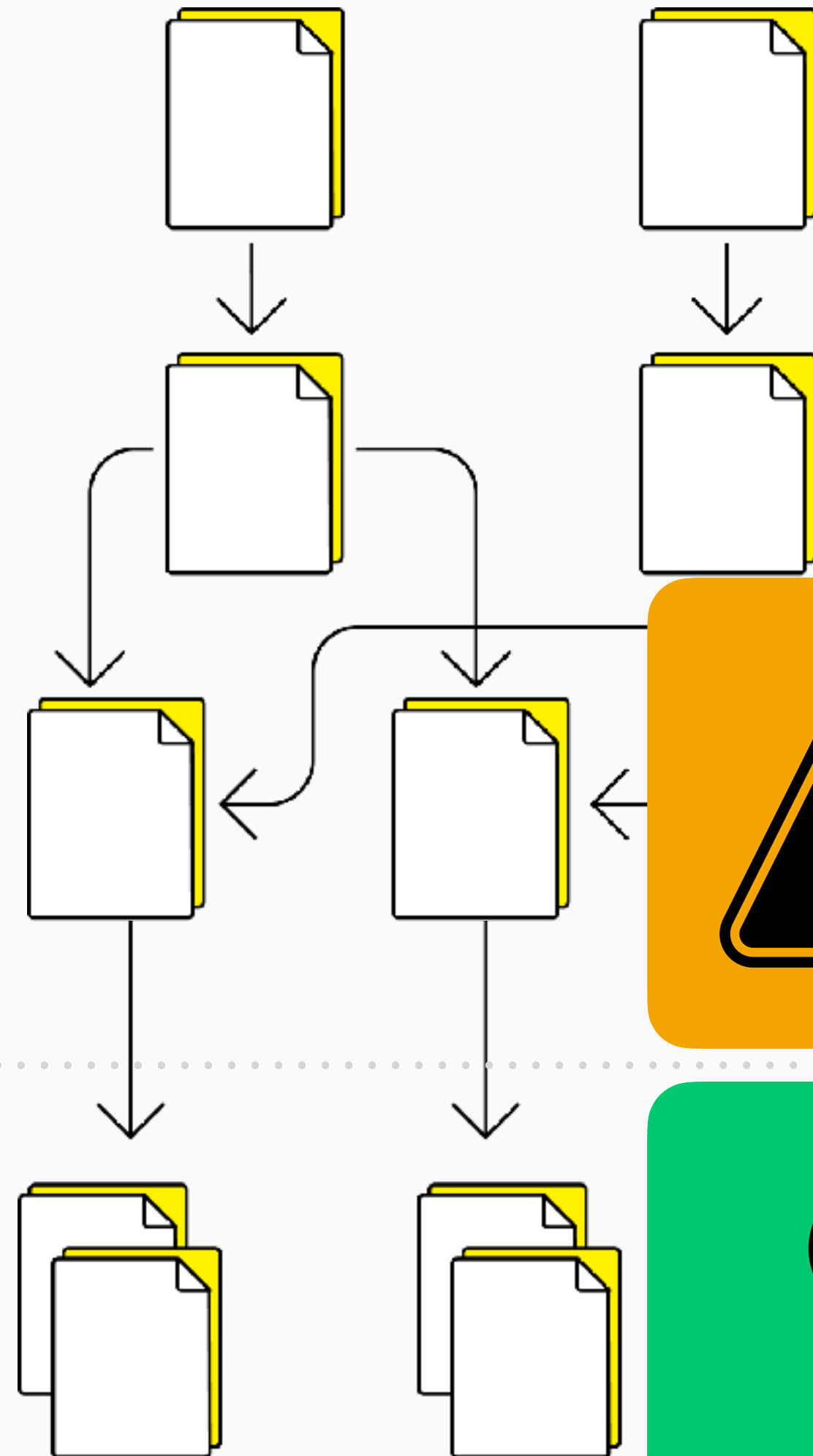


- File-only architecture
- Updatable table

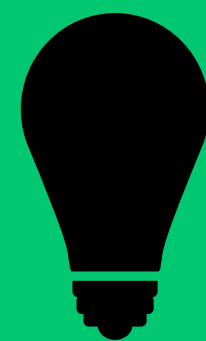
**150M record inserted one-by-one:
600M files (metadata + data)**

- Performance problems
- “Small files” problem
- Only a single table

metadata
(json / avro)
statistics



150M record inserted one-by-one:
600M files (metadata + data)



Compaction algorithms

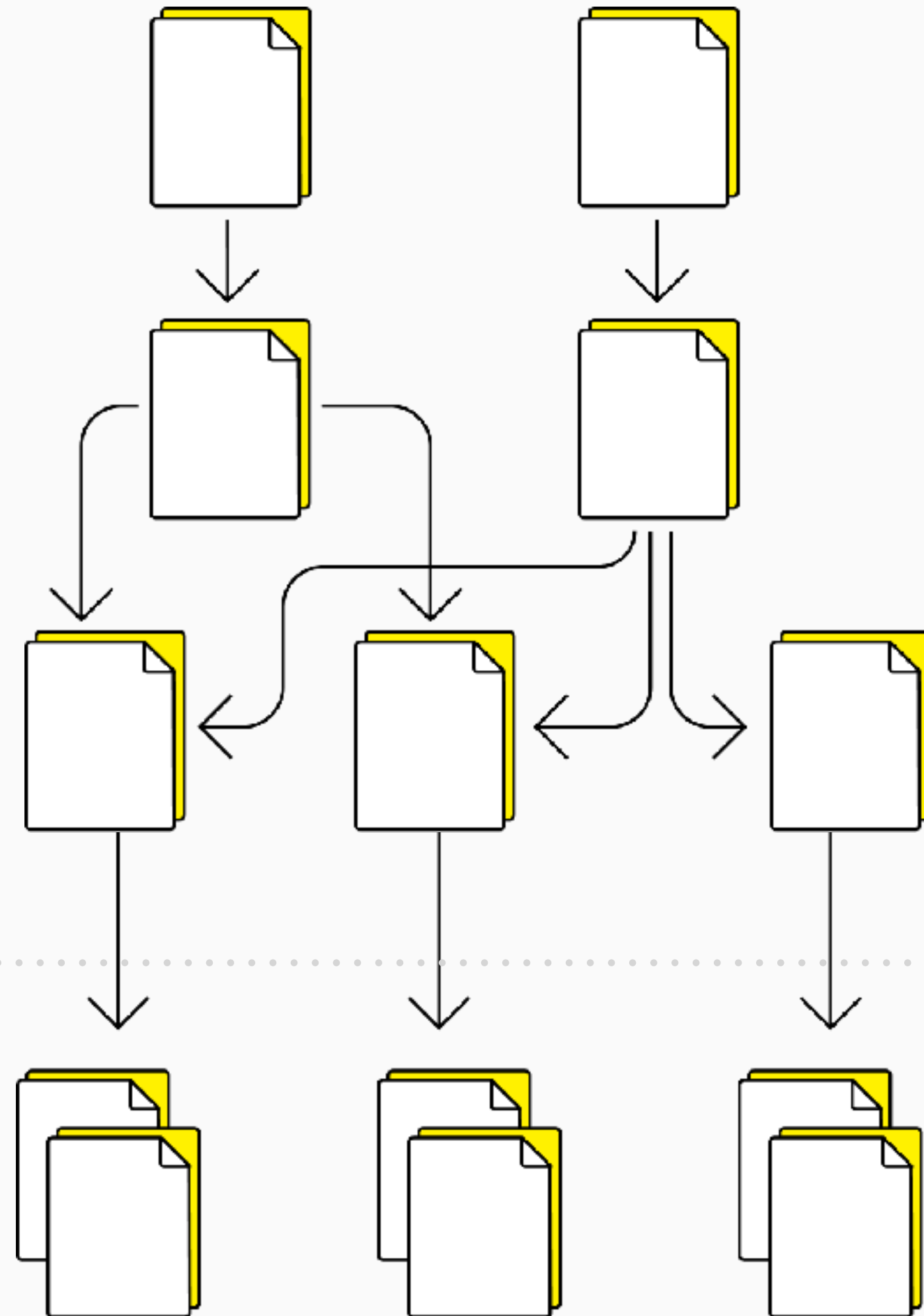


- File-only architecture
- Updatable table

data
(parquet)



metadata
(json / avro)
statistics



data
(parquet)

DuckDB has full Iceberg support!

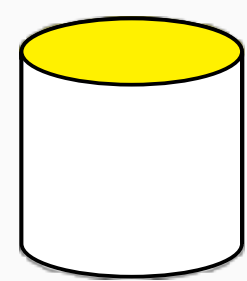
```
INSTALL iceberg;
```

It's an alternative to the pyiceberg package, integrated into DuckDB.

2020

Lakehouse architecture

catalog
(database)



Lakekeeper



Polaris



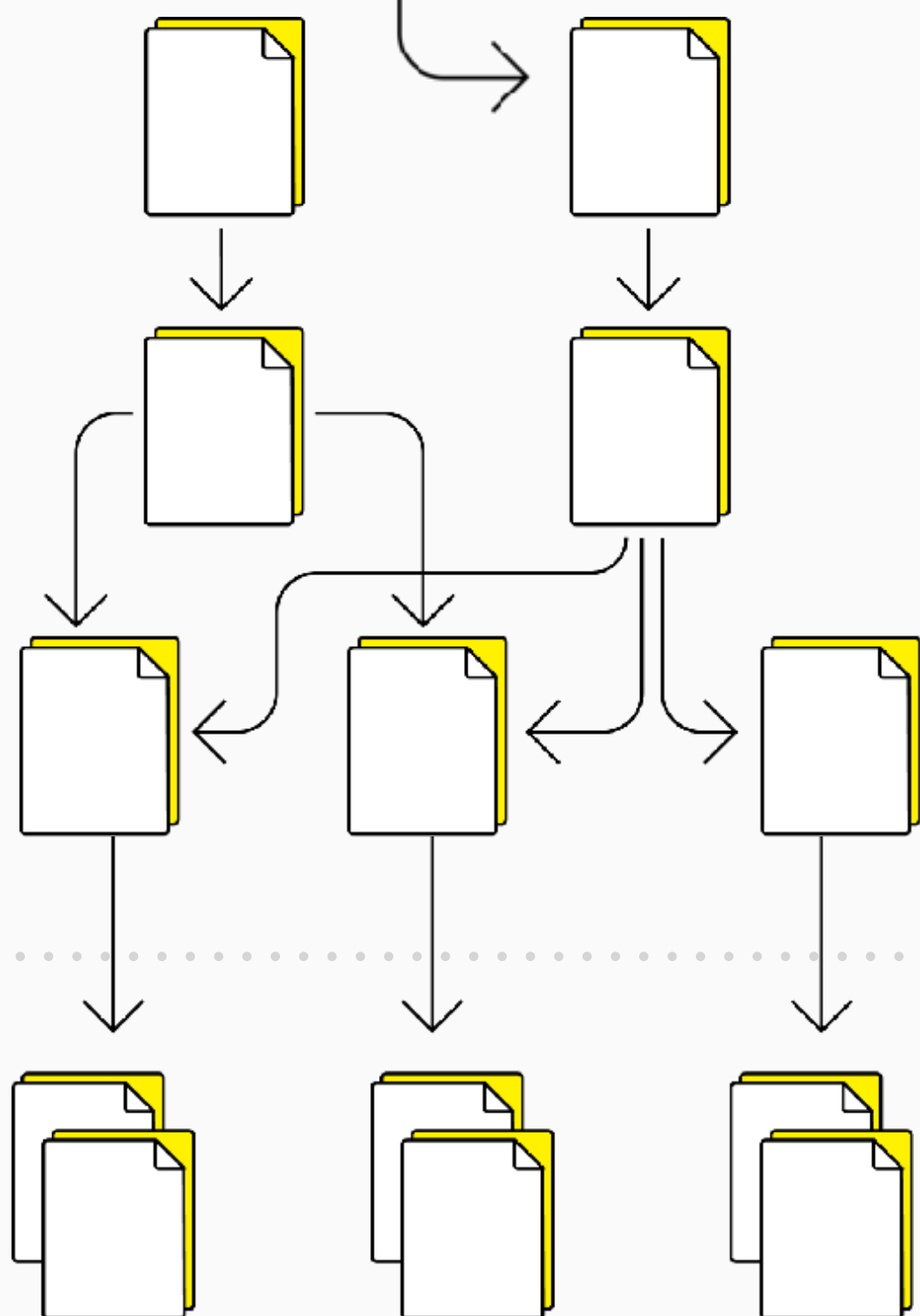
Nessie



Amazon
S3 Tables

metadata
(json / avro)

statistics



data
(parquet)



- Cost-effective
- Small catalog database
- Full database features

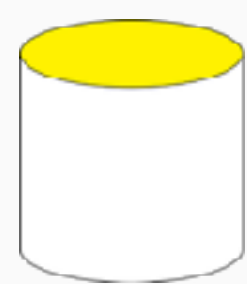


- Performance problems
- Complex operations

2020

Lakehouse architecture

catalog
(database)



Lakekeeper



Polaris



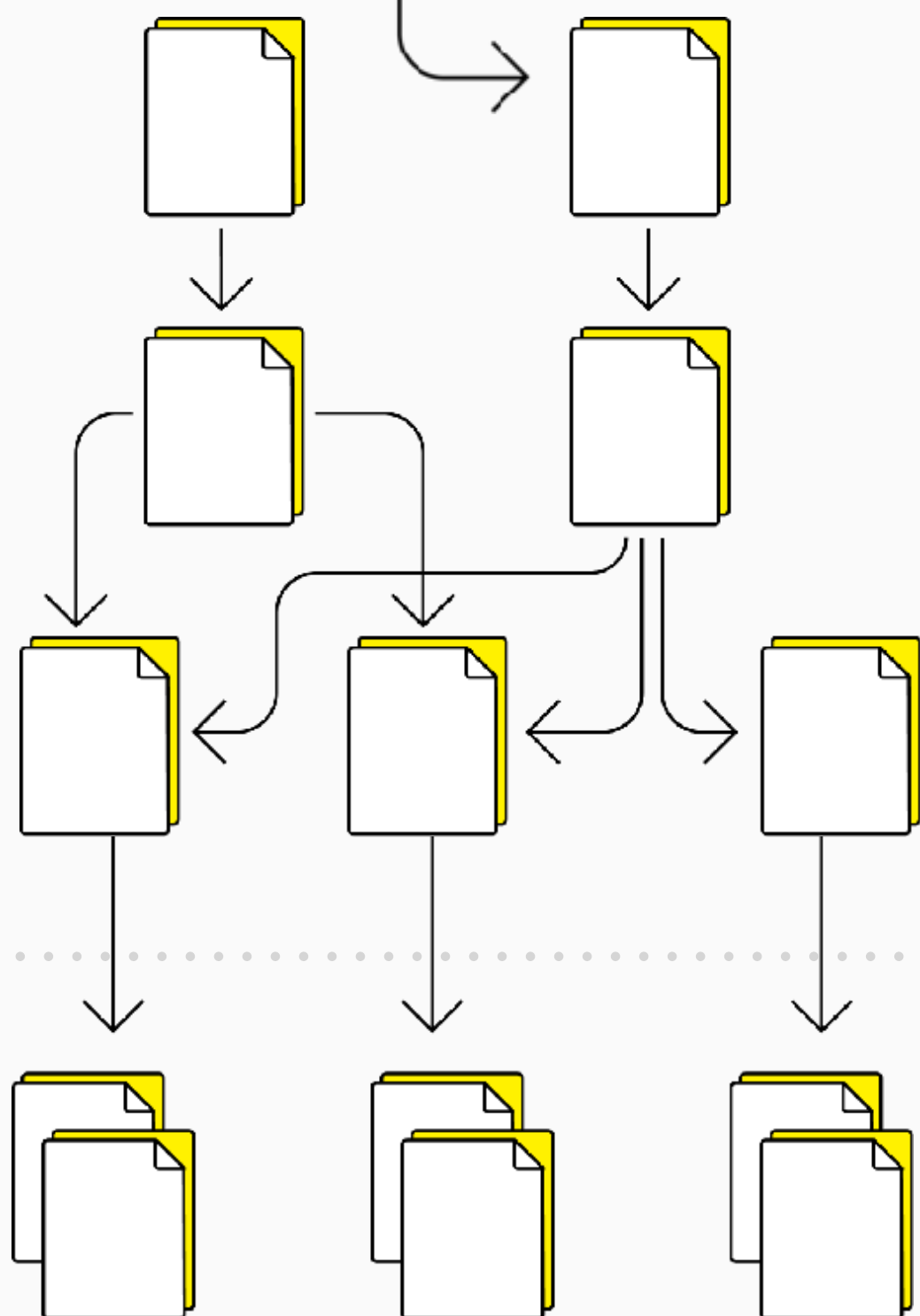
Nessie



Amazon
S3 Tables

metadata
(json / avro)

statistics



data
(parquet)

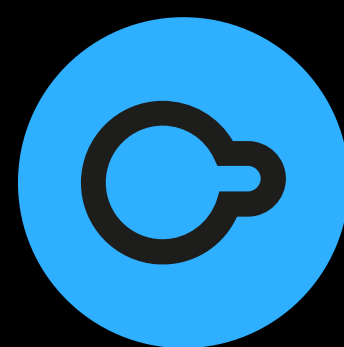


- Cost-effective
- Small catalog database
- Full database features



- Performance problems
- Complex operations

2025



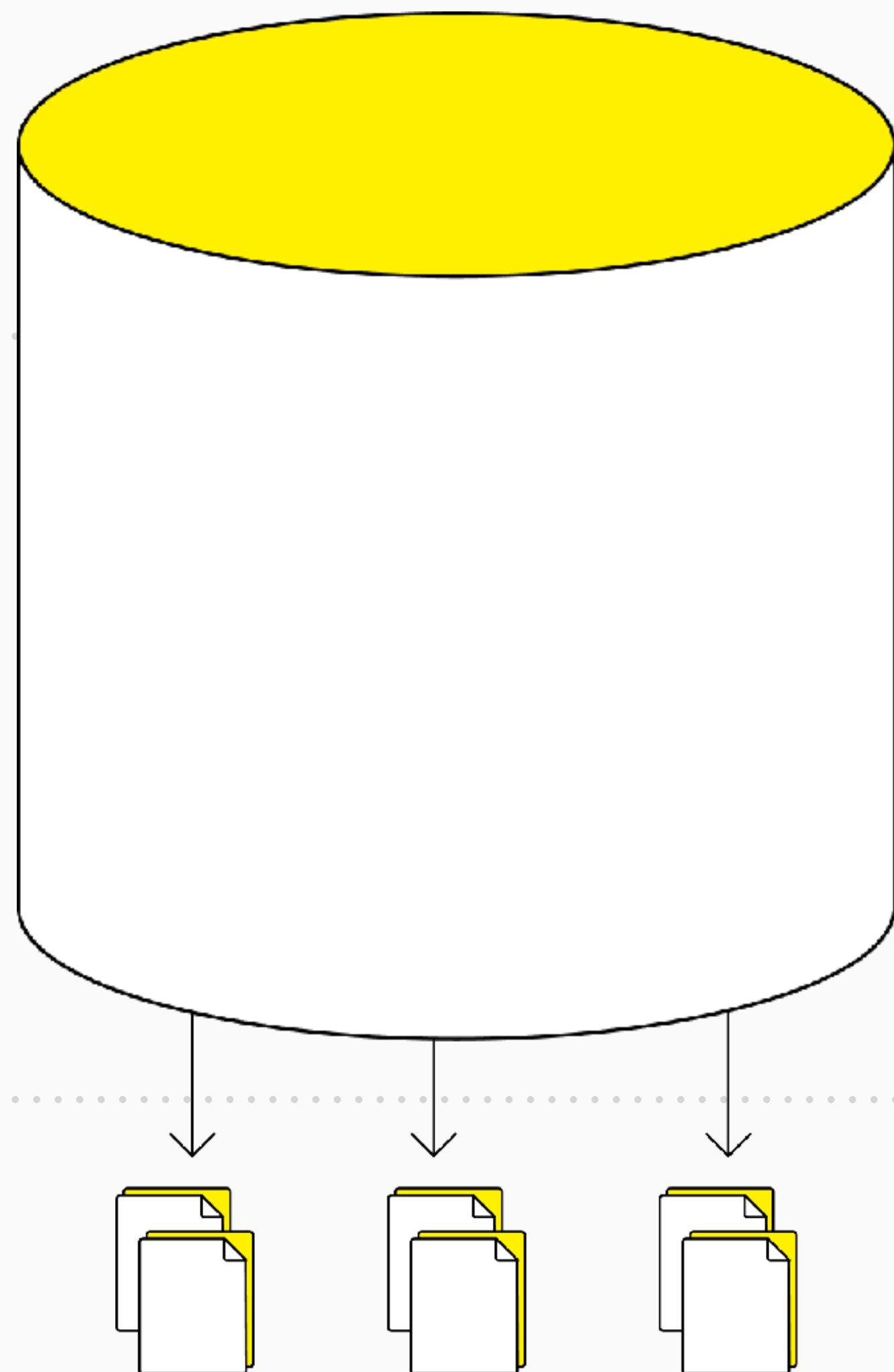
DuckLake

catalog
(database)

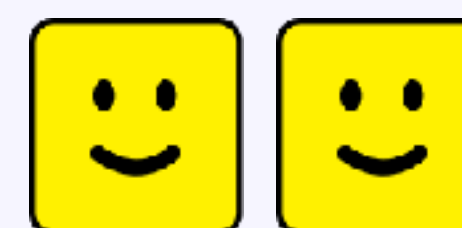
metadata
(also in
the database)

statistics

data
(parquet)



- Cost-effective
- Full database features
- Simple operations
- Scalable



- Snapshot mechanism
- Multi-table transactions
- High-performance

```
ATTACH 'ducklake:trains.ducklake' AS trains;  
USE trains;
```

```
CREATE TABLE services AS  
FROM 'services.csv';
```

```
UPDATE services
SET train_number =
    train_number + (random() * 100 + 1)::INT
WHERE delay > 150;
```

```
UPDATE services
SET train_number =
    train_number + (random() * 100 + 1)::INT
WHERE delay > 150;
```

```
SELECT snapshot_id, changes
FROM ducklake_snapshots('trains');
```

```
UPDATE services
SET train_number =
    train_number + (random() * 100 + 1)::INT
WHERE delay > 150;
```

```
SELECT snapshot_id, changes
FROM ducklake_snapshots('trains');
```

snapshot_id int64	changes map(varchar, varchar[])
0	{schemas_created=[main]}
1	{tables_created=[main.services], tables_inserted_into=[1]}
2	{tables_inserted_into=[1], tables_deleted_from=[1]}

```
SELECT date, train_number, delay
FROM services AT (VERSION => 1)
WHERE date = '2025-02-27'
      AND station = 'Amsterdam Centraal'
      AND train_number = 420;
```

```
SELECT date, train_number, delay
FROM services AT (VERSION => 1)
WHERE date = '2025-02-27'
      AND station = 'Amsterdam Centraal'
      AND train_number = 420;
```

```
SELECT date, train_number, delay
FROM services AT (VERSION => 1)
WHERE date = '2025-02-27'
      AND station = 'Amsterdam Centraal'
      AND train_number = 420;
```

date	train_number	delay
2025-02-27	420	174

```
SELECT date, train_number, delay
FROM services AT (VERSION => 1)
WHERE date = '2025-02-27'
      AND station = 'Amsterdam Centraal'
      AND train_number = 420;
```

date	train_number	delay
2025-02-27	420	174



```
SELECT date, train_number, delay
FROM services AT (VERSION => 1)
WHERE date = '2025-02-27'
      AND station = 'Amsterdam Centraal'
      AND train_number = 420;
```

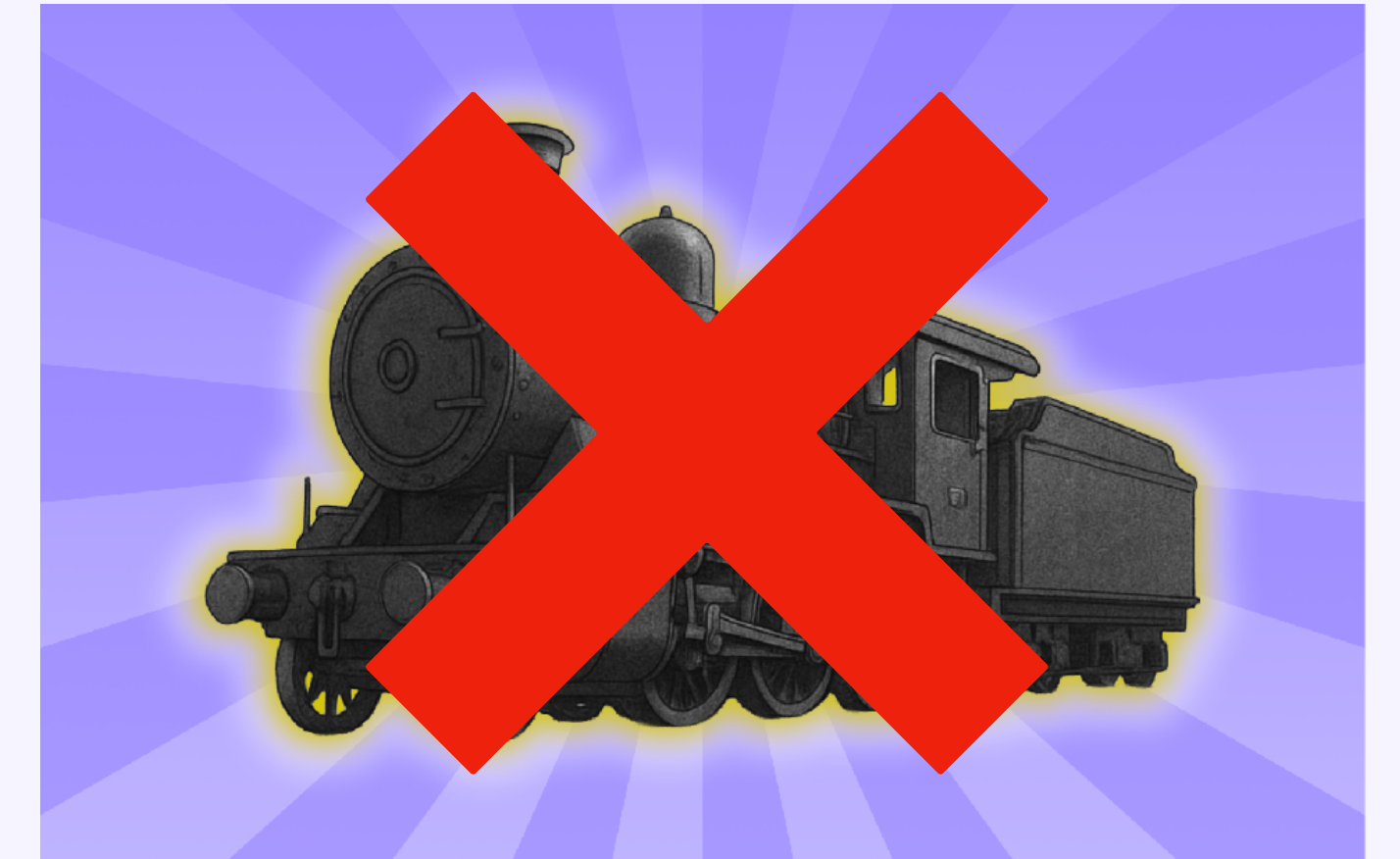
date	train_number	delay
2025-02-27	420	174



Passengers - Railway co.
1 - 0

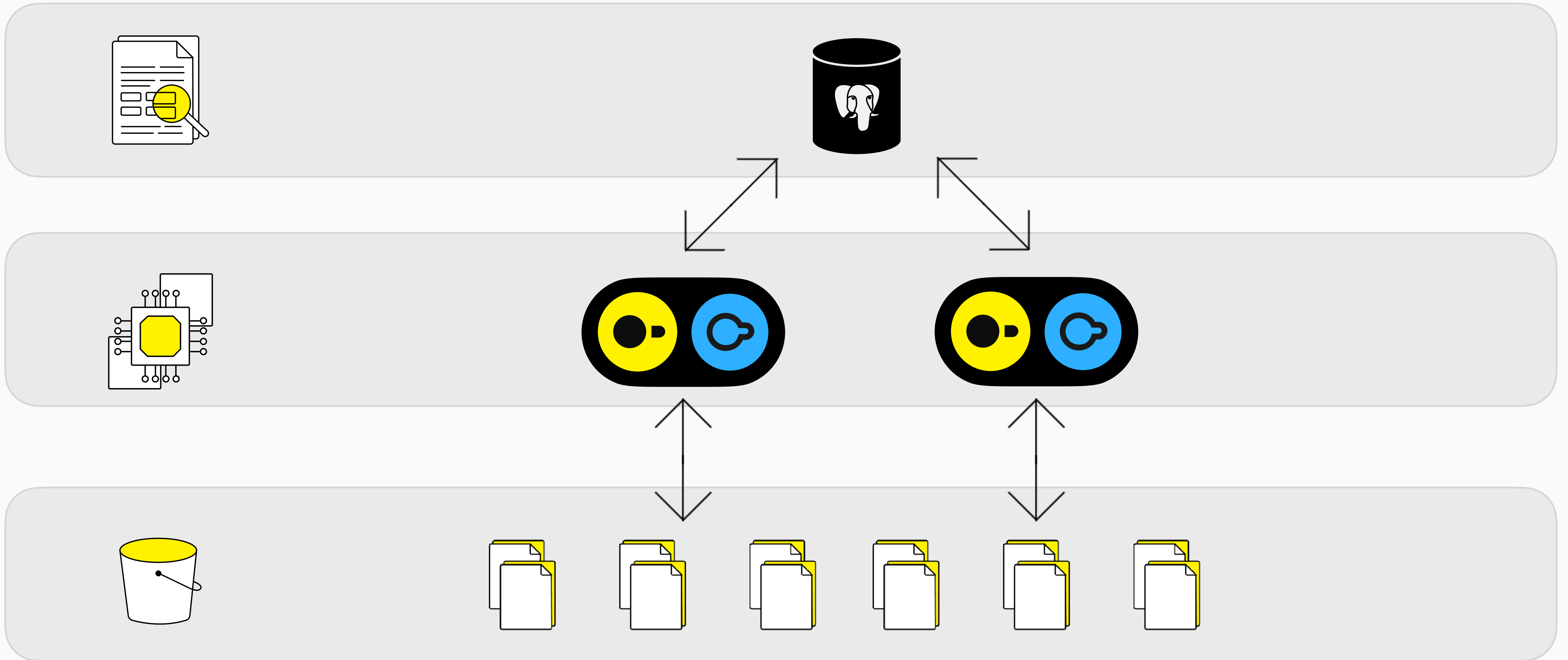
```
SELECT date, train_number, delay
FROM services AT (VERSION => 1)
WHERE date = '2025-02-27'
      AND station = 'Amsterdam Centraal'
      AND train_number = 420;
```

date	train_number	delay
2025-02-27	420	174



Passengers - Railway co.
1 - 0

Re-run improved scripts
on old datasets



We release DuckLake v1.0
last Monday


Focus #1:
Stable specification

Focus #2:
Production-readiness



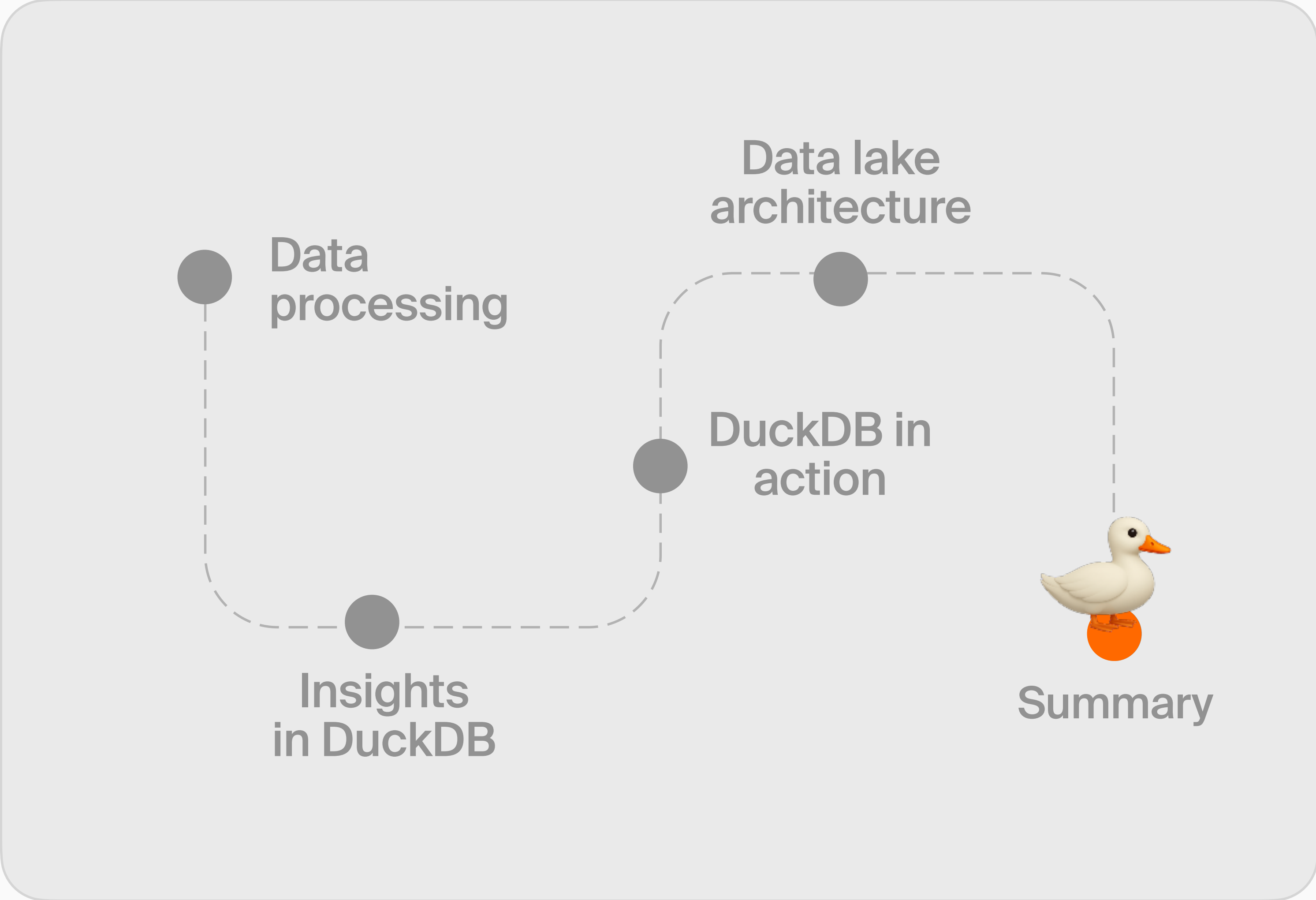
DuckDB uses RDBMS to attack classic 'small changes' problem in lakehouses

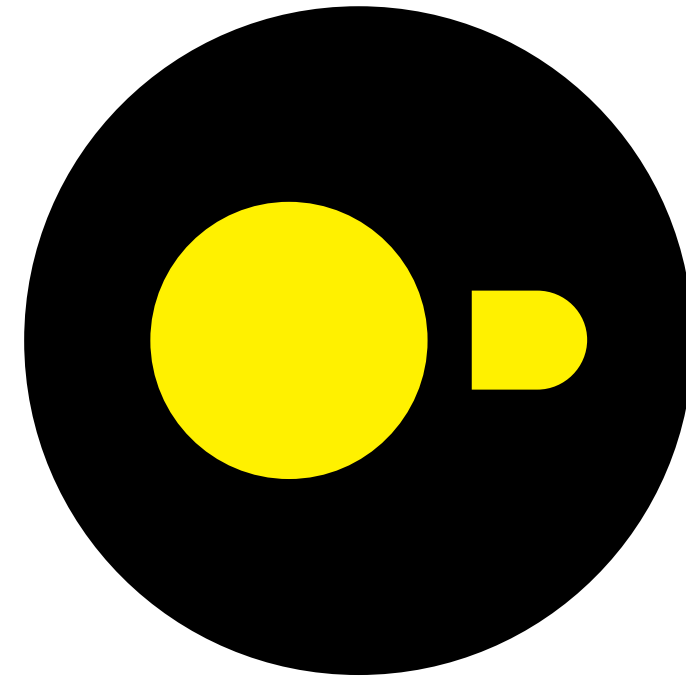
Batching teensy changes in chunks creates massive performance boost, DuckDB Labs team claims

 [Lindsay Clark](#)

Thu 16 Apr 2026 // 16:15 UTC

The team behind in-process OLAP database DuckDB has put forward a solution to the "small changes" problem that they say plagues lakehouse implementations of the kind based on technologies from Databricks, Snowflake, Google, and others.





DuckDB

**An open-source in-process analytical
database management system
that you can import as a library**

SQL and DuckDB internals from Universität Tübingen

Tabular Database Systems

Design and Implementation of DuckDB Internals

07

Inner Join Follows Foreign Keys

vehicles					
vid	vehicle	kind	seats	wheels?	pid
v ₁		car	5	true	p ₄
v ₂		SUV	3	true	p ₄
v ₃		bus	42	true	□
v ₄		bus	7	true	□
v ₅		bike	1	true	p ₂
v ₆		tank	0	false	p ₃
v ₇		cabrio	2	true	p ₄

peeps			
pid	pic	name	born
p ₁		Cleo	2013
p ₂		Bert	1968
p ₃		Drew	□
p ₄		Alex	2002

```
SELECT v.vehicle, p.name AS driver, p.pic
FROM vehicles AS v INNER JOIN peeps AS p ON (v.pid = p.pid)
                                "equi-join" ↗
```

- Some rows in table `peeps` find multiple *join partners* in `vehicles` (like `p4` 😊), some find none (`p1` 😞).
- General: A join between `t1` and `t2` may yield $0 \dots |t_1| \times |t_2|$ rows.

05

DuckDB: A Tabular DBMS Inside Your Own Process

The diagram illustrates the DuckDB architecture. On the left, a box labeled 'DuckDB CLI' contains a CPU icon and a memory icon. An arrow points from this box to a file icon labeled 'File *.db'. Another arrow points from the file icon to a box labeled 'Python 3', which also contains a CPU icon and a memory icon. A dashed arrow points from the Python 3 box to a file icon labeled '*. {csv,json}'.

- DuckDB kernel and client share a **single process** `□`.
 - Python: `import duckdb`, C/C++: link with `libduckdb`.
- Table data resides in a **single database file** `📄 *.db`, native DuckDB data formats `📄` in file and in RAM are similar.
- DuckDB sees in-process client data and can **directly read/write client data structures** `📄` using SQL (“zero copy”, replacement scans).
- DuckDB allocates sizable RAM buffers (but can use disk `📄` for temporary storage if required).

