

zoology 101: python, pandas and duckdb

presenter:

Jaroslav Bezdek

date:

16 November 2024



01

motivation



MY_TOOLBOX

MY_TOOLBOX



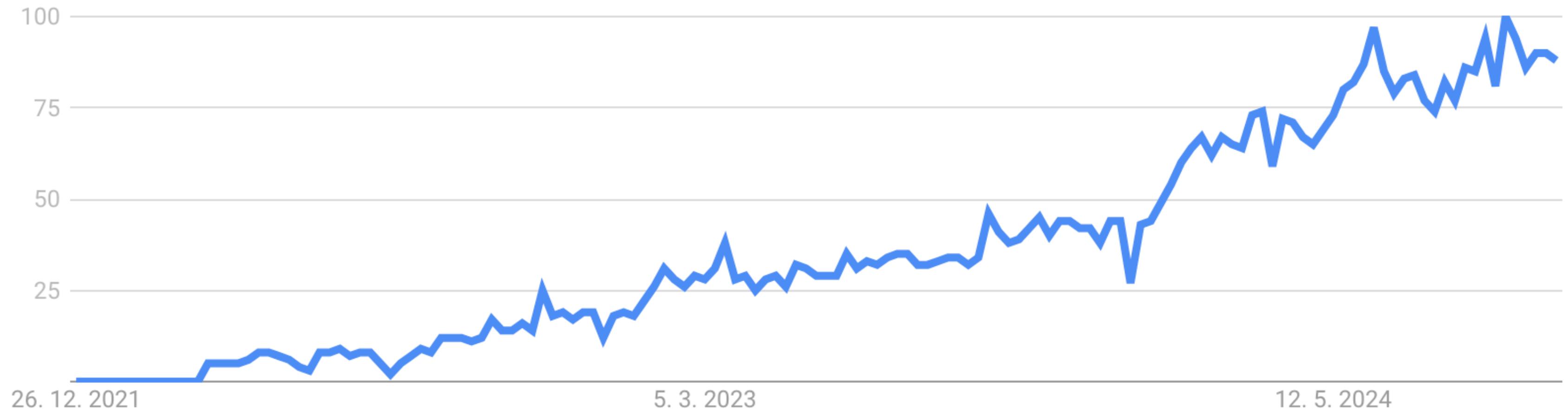
MY_TOOLBOX



MY_TOOLBOX

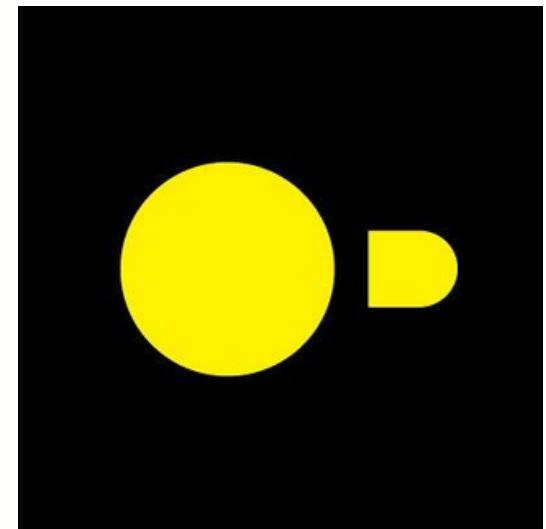


DUCKDB_INTEREST_OVER_TIME

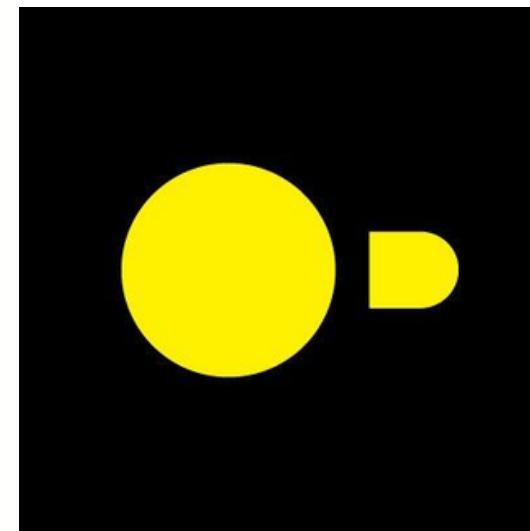
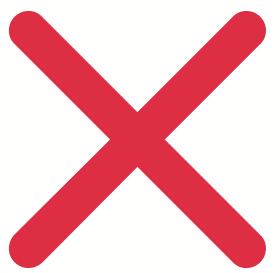


Source: Google Trends

MY_TOOLBOX_WITH_DUCKDB



PANDAS_VS_DUCKDB



02

data



```
1 >> print(pinguins)
2
3 shape: (344, 6)
4
5  pinguin_id    date      species   bill_length   bill_depth   sex
6  ---          ---       ---        ---           ---          ---
7  str          str       str        f64           f64          str
8
9  N1A1         2007-11-11 adelie     39.1          18.7         MALE
10 N1A2         2007-11-11 adelie     39.5          17.4         FEMALE
11 N2A1         2007-11-16 adelie     40.3          18.0         FEMALE
12 N2A2         2007-11-16 adelie     null          null         null
13 N3A1         2007-11-16 adelie     36.7          19.3         FEMALE
14 ...
15 N38A2        2009-12-01 gentoo    null          null         null
16 N39A1        2009-11-22 gentoo    46.8          14.3         FEMALE
17 N39A2        2009-11-22 gentoo    50.4          15.7         MALE
18 N43A1        2009-11-22 gentoo    45.2          14.8         FEMALE
19 N43A2        2009-11-22 gentoo    49.9          16.1         MALE
20
```

03

import pandas

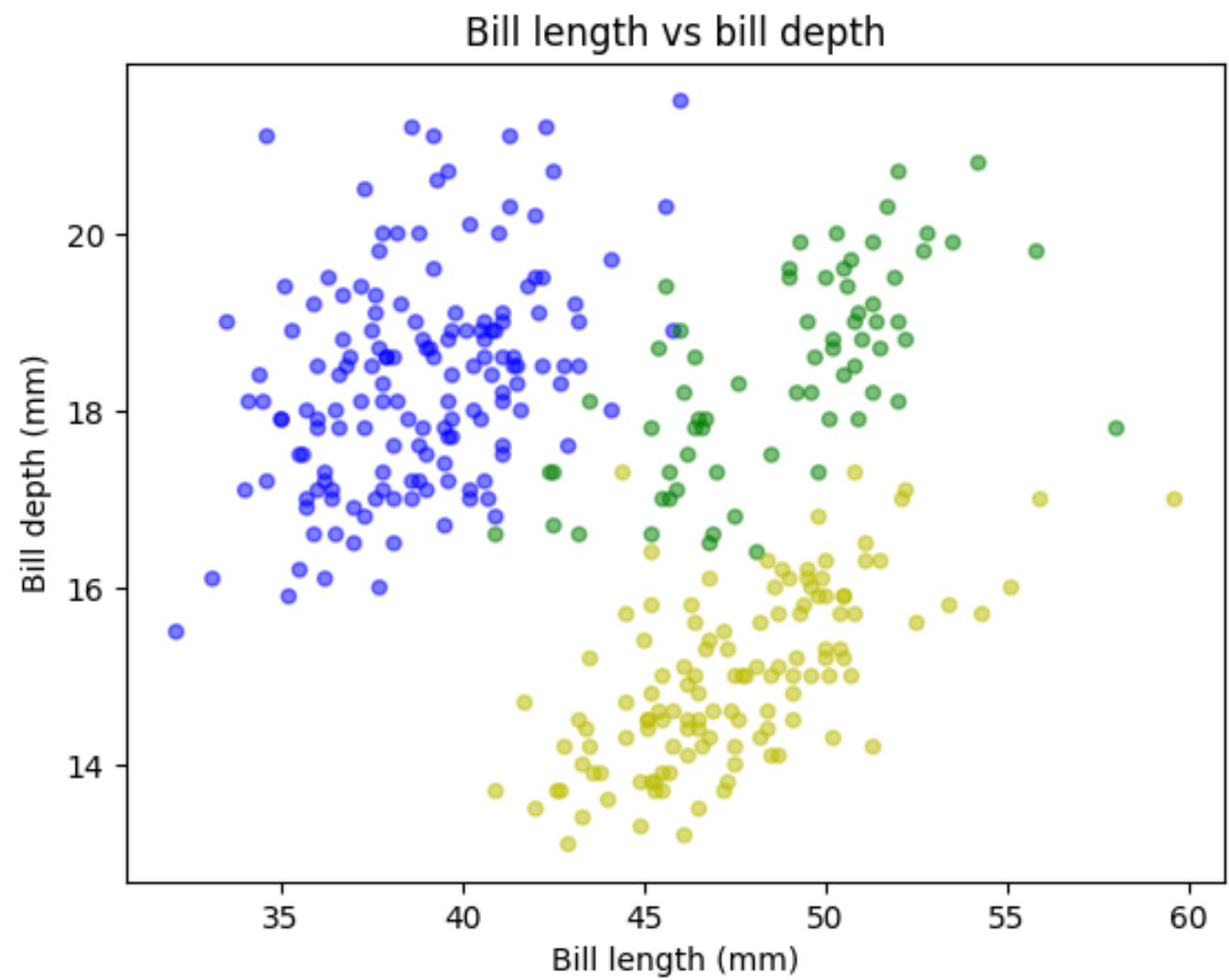


```
1 df = pd.read_csv(  
2     filepath_or_buffer=DATA_PATH,  
3 )  
4  
5 print(df)  
6  
7 #      pinguin_id        date   species  bill_length  bill_depth    sex  
8 # 0      N1A1  2007-11-11  adelie       39.1       18.7  MALE  
9 # 1      N1A2  2007-11-11  adelie       39.5       17.4 FEMALE  
10 # 2      N2A1  2007-11-16  adelie       40.3       18.0 FEMALE  
11 # 3      N2A2  2007-11-16  adelie       NaN         NaN  NaN  
12 # 4      N3A1  2007-11-16  adelie       36.7       19.3 FEMALE  
13 # ...    ...      ...      ...      ...      ...      ...  
14 # 339    N38A2  2009-12-01  gentoo      NaN         NaN  NaN  
15 # 340    N39A1  2009-11-22  gentoo       46.8       14.3 FEMALE  
16 # 341    N39A2  2009-11-22  gentoo       50.4       15.7  MALE  
17 # 342    N43A1  2009-11-22  gentoo       45.2       14.8 FEMALE  
18 # 343    N43A2  2009-11-22  gentoo       49.9       16.1  MALE  
19 #  
20 # [344 rows x 6 columns]
```

```
1 print(  
2     df  
3     .assign(  
4         sex=lambda _df: _df.sex.str.lower(),  
5     )  
6     .groupby(by="species")  
7     .agg(  
8         samples_cnt=("pinguin_id", "count"),  
9         male_cnt=("sex", lambda x: x[x == "male"].count()),  
10        bill_length_avg=("bill_length", "mean"),  
11    )  
12    .round(1)  
13 )  
14  
15 #          samples_cnt  male_cnt  bill_length_avg  
16 # species  
17 # adelie           152       73      38.8  
18 # chinstrap         68        34      48.8  
19 # gentoo           124       61      47.5
```

```
1 # convert string to datetime
2 df["date"] = pd.to_datetime(df.date)
3
4 # get unique sample dates from year 2007
5 unique_dates_2007 = (
6     df
7     .date
8     .loc[df.date.dt.year == 2007]
9     .drop_duplicates()
10    .sort_values()
11    .reset_index(drop=True)
12 )
13
14 # get the difference between sample days
15 print(unique_dates_2007 - unique_dates_2007.shift())
16
17 # 0      NaT
18 # 1      1 days
19 # 2      1 days
20 # ..
21 # 14     1 days
22 # 15     1 days
23 # 16     3 days
```

```
1  (
2    df
3    .plot(
4      kind="scatter",
5      x="bill_length",
6      y="bill_depth",
7      color=df.species.map({
8        "adelie": "b",
9        "chinstrap": "g",
10       "gentoo": "y",
11     }),
12     alpha=0.5,
13     xlabel="Bill length (mm)",
14     ylabel="Bill depth (mm)",
15     title="Bill length vs bill depth",
16   )
17 )
```



STRENGTHS_AND_WEAKNESSES



- **Swiss army knife of data analytics**

STRENGTHS_AND_WEAKNESSES



- **Swiss army knife of data analytics**
- **Integration with other Python libraries**

STRENGTHS_AND_WEAKNESSES



- **Swiss army knife of data analytics**
- **Integration with other Python libraries**
- **Community and documentation**

STRENGTHS_AND_WEAKNESSES



- **Swiss army knife of data analytics**
- **Integration with other Python libraries**
- **Community and documentation**



- **Performance with large datasets**

STRENGTHS_AND_WEAKNESSES



- **Swiss army knife of data analytics**
- **Integration with other Python libraries**
- **Community and documentation**



- **Performance with large datasets**
- **Focus on structured data**

STRENGTHS_AND_WEAKNESSES



- **Swiss army knife of data analytics**
- **Integration with other Python libraries**
- **Community and documentation**



- **Performance with large datasets**
- **Focus on structured data**
- **Steeper learning curve for advanced features**

04

import duckdb

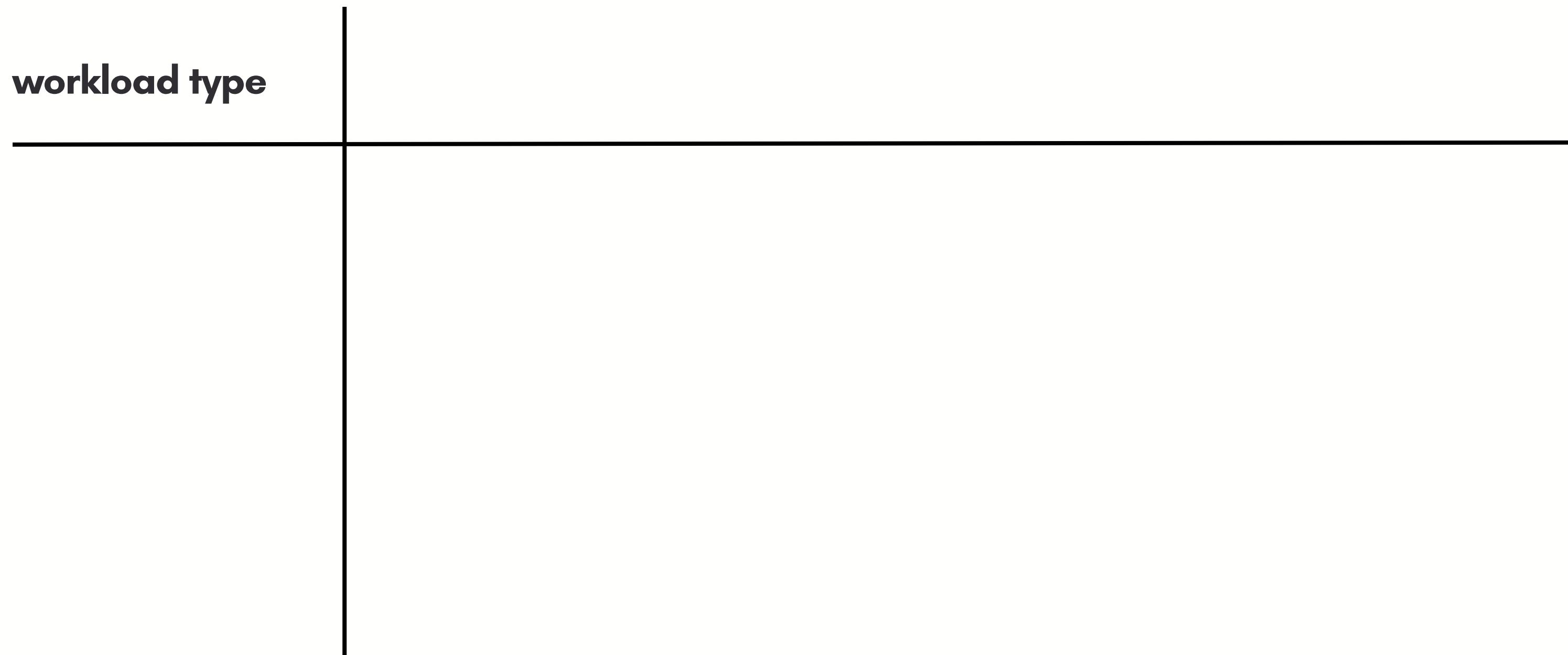


**“DuckDB is a fast in-process
analytical database.”**



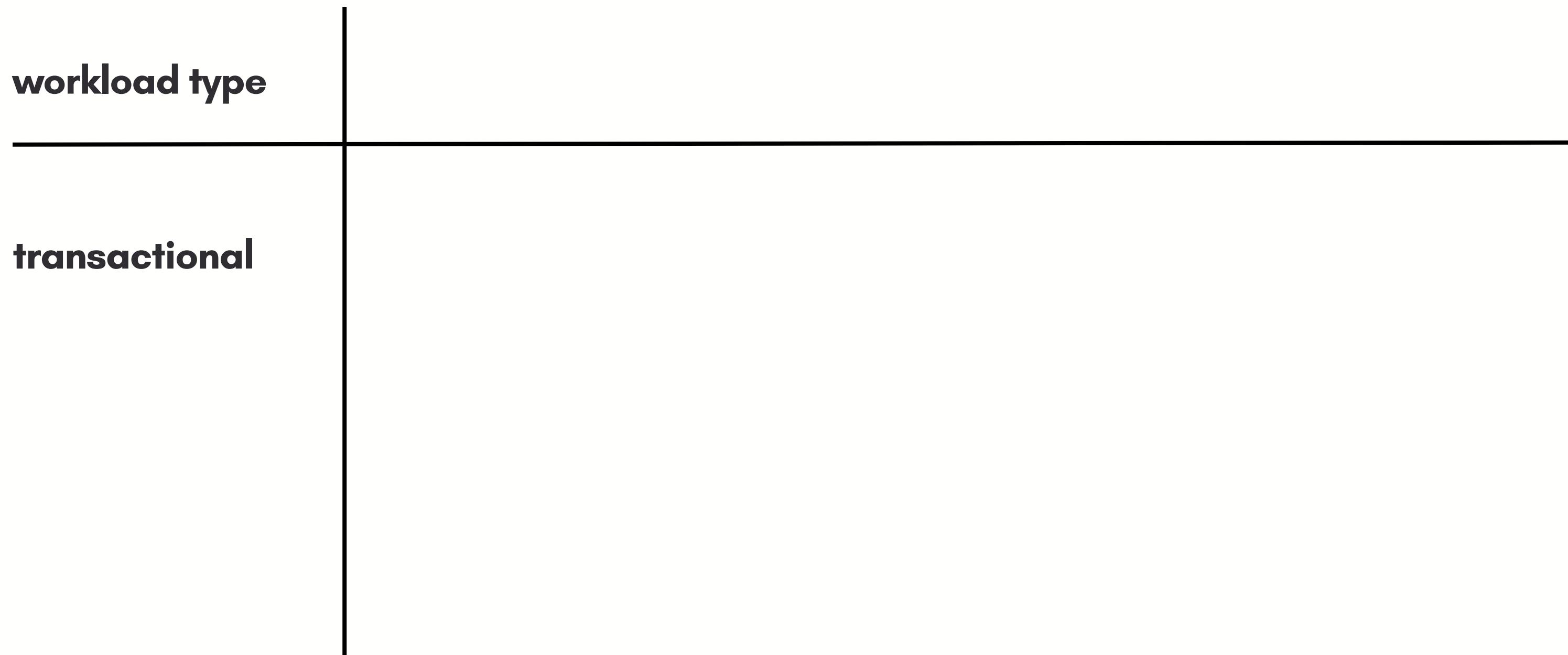
duckdb.org

DATABASE_WORKLOAD_TYPES



Source: Kojo Osei

DATABASE_WORKLOAD_TYPES



Source: Kojo Osei

DATABASE_WORKLOAD_TYPES

workload type	example query
transactional	update users table with new user object after signup

Source: Kojo Osei

DATABASE_WORKLOAD_TYPES

workload type	example query	example database
transactional	update users table with new user object after signup	Postgres, SQLite

Source: Kojo Osei

DATABASE_WORKLOAD_TYPES

workload type	example query	example database
transactional	update users table with new user object after signup	Postgres, SQLite
analytical		

Source: Kojo Osei

DATABASE_WORKLOAD_TYPES

workload type	example query	example database
transactional	update users table with new user object after signup	Postgres, SQLite
analytical	how many users have signed up in the last 2 weeks, broken by age and location?	

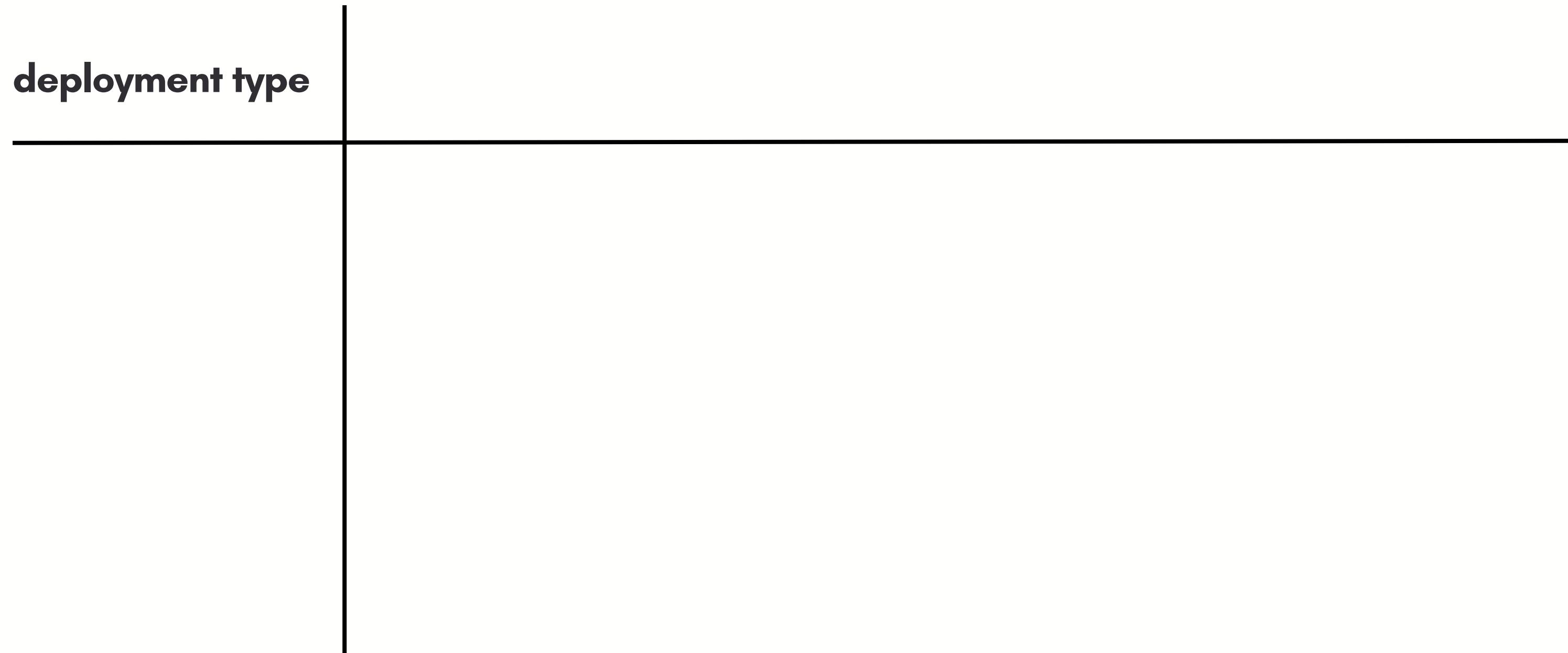
Source: Kojo Osei

DATABASE_WORKLOAD_TYPES

workload type	example query	example database
transactional	update users table with new user object after signup	Postgres, SQLite
analytical	how many users have signed up in the last 2 weeks, broken by age and location?	Snowflake, BigQuery

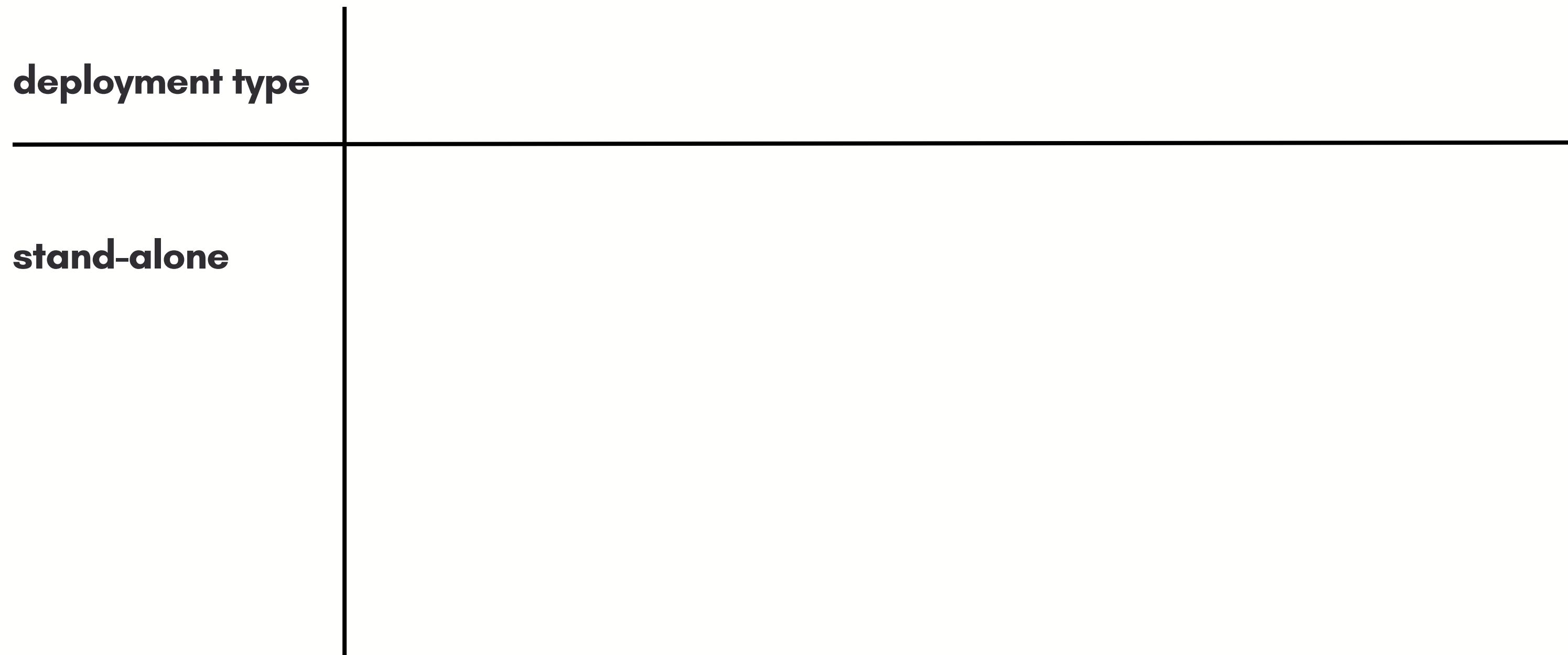
Source: Kojo Osei

DATABASE_DEPLOYMENT_TYPES



Source: Kojo Osei

DATABASE_DEPLOYMENT_TYPES



Source: Kojo Osei

DATABASE_DEPLOYMENT_TYPES

deployment type	example database
stand-alone	Snowflake, Postgres, BigQuery

Source: Kojo Osei

DATABASE_DEPLOYMENT_TYPES

deployment type	example database
stand-alone	Snowflake, Postgres, BigQuery
embedded	

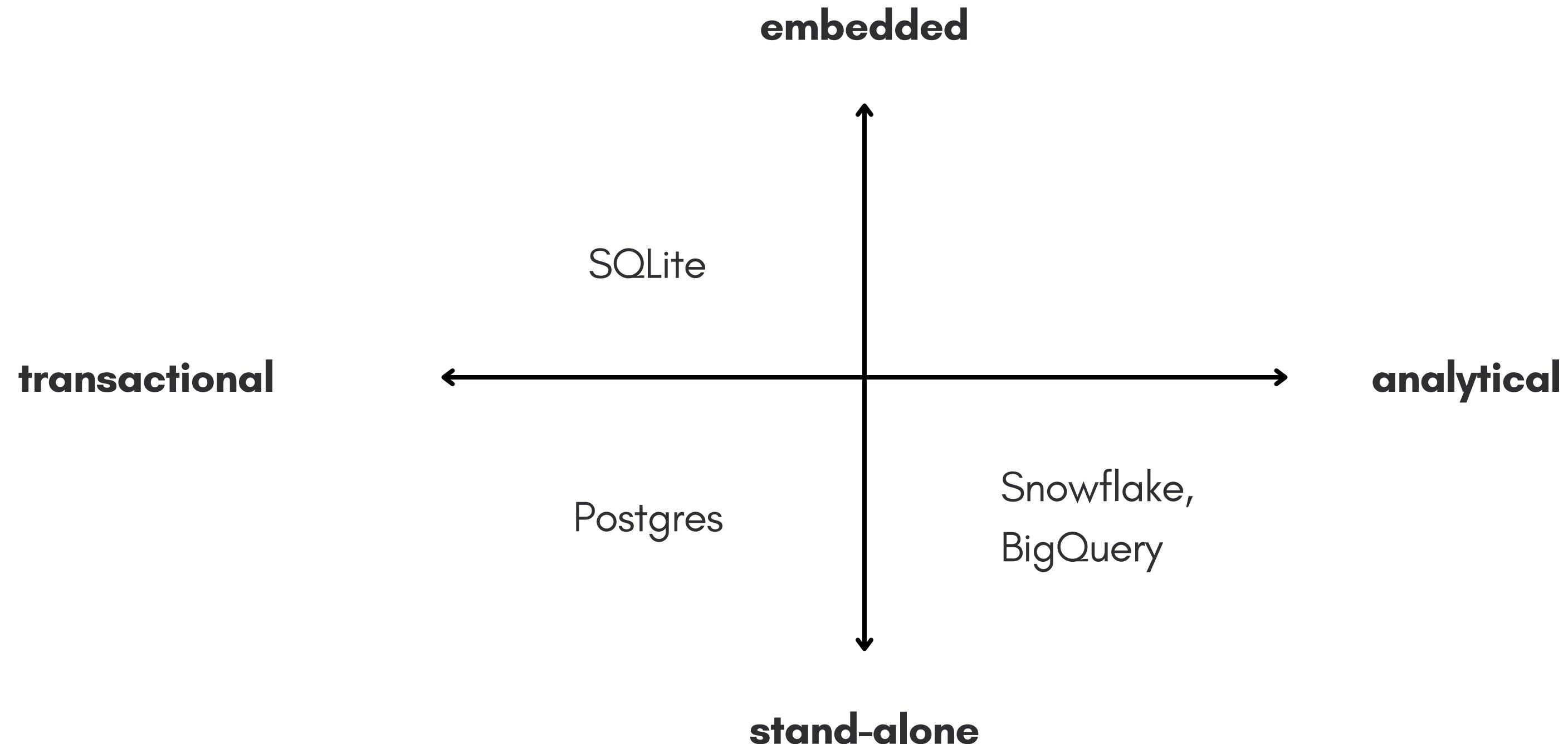
Source: Kojo Osei

DATABASE_DEPLOYMENT_TYPES

deployment type	example database
stand-alone	Snowflake, Postgres, BigQuery
embedded	SQLite

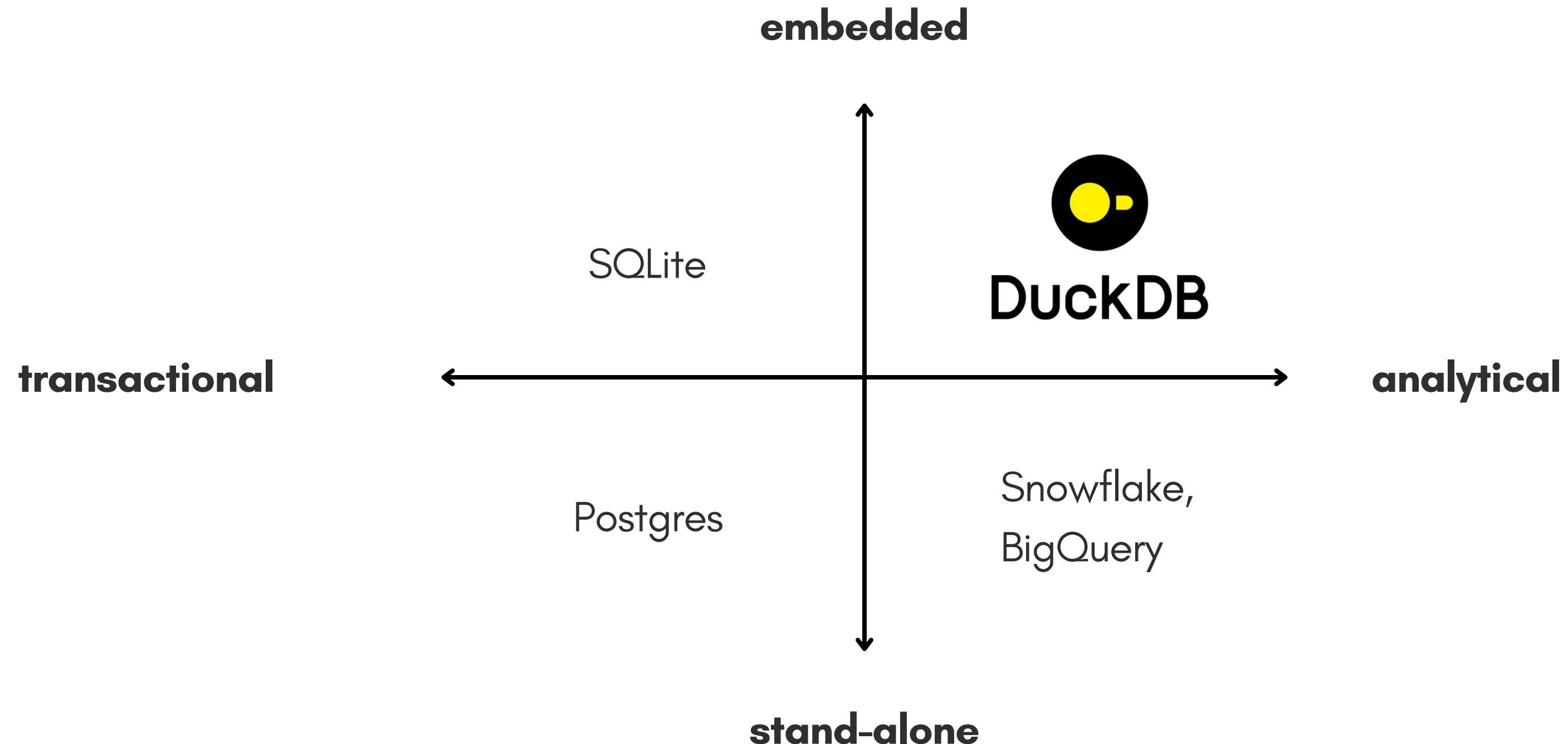
Source: Kojo Osei

DATABASE_TYPES



Source: Kojo Osei

DATABASE_TYPES



Source: [Kojo Osei](#)

```
1 import duckdb
2
3 data_raw = (
4     duckdb
5         .read_csv(name=DATA_PATH)
6         .select("species", "bill_length", "bill_depth", "sex")
7 )
8
9 print(data_raw)
10
11 #> #> species bill_length bill_depth sex
12 #> #> varchar    double      double    varchar
13 #> #> +-----+-----+-----+
14 #> #> |       |       |       |
15 #> #> | adelie | 39.1  | 18.7   | MALE
16 #> #> | adelie | 39.5  | 17.4   | FEMALE
17 #> #> | .      | .     | .     | .
18 #> #> | gentoo | 45.2  | 14.8   | FEMALE
19 #> #> | gentoo | 49.9  | 16.1   | MALE
20 #> #> +-----+-----+-----+
21 #> #> 344 rows (20 shown) 4 columns
22 #>
```

```
1  (
2      data_raw
3      .project("""
4          species,
5          lower(sex) as sex,
6          bill_length,
7          bill_depth
8      """)
9      .filter("sex = 'male'")
10     .aggregate("""
11         species,
12         avg(bill_length)::int as bill_length_avg,
13         avg(bill_depth)::int as bill_depth_avg
14     """)
15     .order("species")
16 )
17
18 #




```

```

1  (
2      data_raw
3      .project("""
4          species,
5          lower(sex) as sex,
6          bill_length,
7          bill_depth
8      """)

9      .filter("sex = 'male'")

10     .aggregate("""
11         species,
12         avg(bill_length)::int as bill_length_avg,
13         avg(bill_depth)::int as bill_depth_avg
14     """)

15     .order("species")
16 )
17
18 #




```

```

1 duckdb.sql(query=f"""
2     with cleaned as (
3         select
4             species,
5             lower(sex) as sex,
6             bill_length,
7             bill_depth
8         from data_raw
9         where sex is not null
10    ),
11    agg as (
12        select
13            species,
14            avg(bill_length)::int as bill_length_avg,
15            avg(bill_depth)::int as bill_depth_avg
16        from cleaned
17        where sex = 'male'
18        group by species
19        order by species
20    )
21
22    select * from agg;
23 """
24
25 )
26
27 #




```

```
1 df_agg = agg.to_df()
2 print(df_agg)
3
4 #      species  bill_length_avg  bill_depth_avg
5 # 0    adelie        40            19
6 # 1  chinstrap       51            19
7 # 2    gentoo        49            16
8
9 agg = duckdb.sql("select * from df_agg")
10 print(agg)
11
12 #




```

STRENGTHS_AND_WEAKNESSES



- **Performance**

STRENGTHS_AND_WEAKNESSES



- **Performance**
- **Reading multiple files at once**

STRENGTHS_AND_WEAKNESSES



- **Performance**
- **Reading multiple files at once**
- **SQL syntax**

STRENGTHS_AND_WEAKNESSES



- **Performance**
- **Reading multiple files at once**
- **SQL syntax**
- **APIs for major programming languages**

STRENGTHS_AND_WEAKNESSES



- **Performance**
- **Reading multiple files at once**
- **SQL syntax**
- **APIs for major programming languages**
- **“Only” analytical workloads**

STRENGTHS_AND_WEAKNESSES



- **Performance**
- **Reading multiple files at once**
- **SQL syntax**
- **APIs for major programming languages**



- **“Only” analytical workloads**
- **Community support**

STRENGTHS_AND_WEAKNESSES



- **Performance**
- **Reading multiple files at once**
- **SQL syntax**
- **APIs for major programming languages**



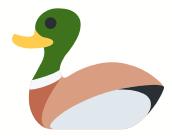
- **“Only” analytical workloads**
- **Community support**
- **Immaturity**

05

comparison



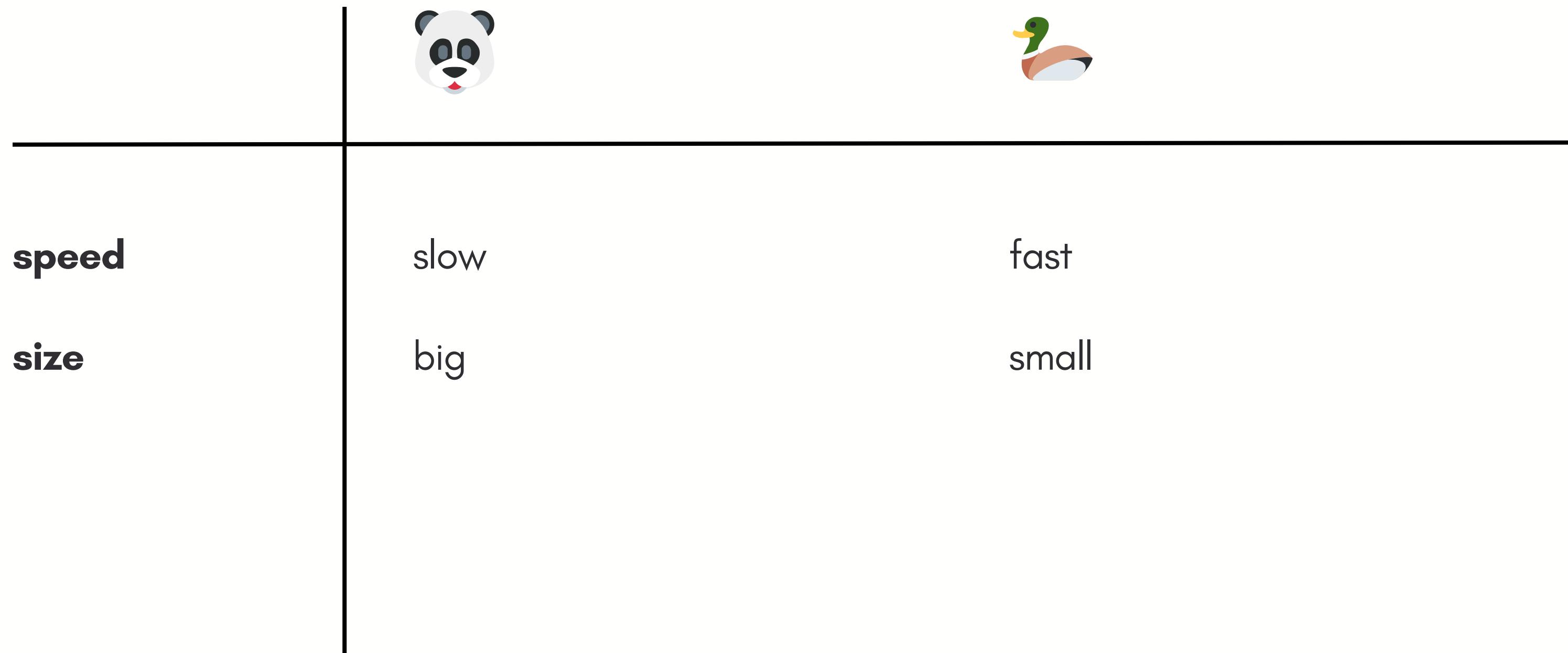
COMPARISON



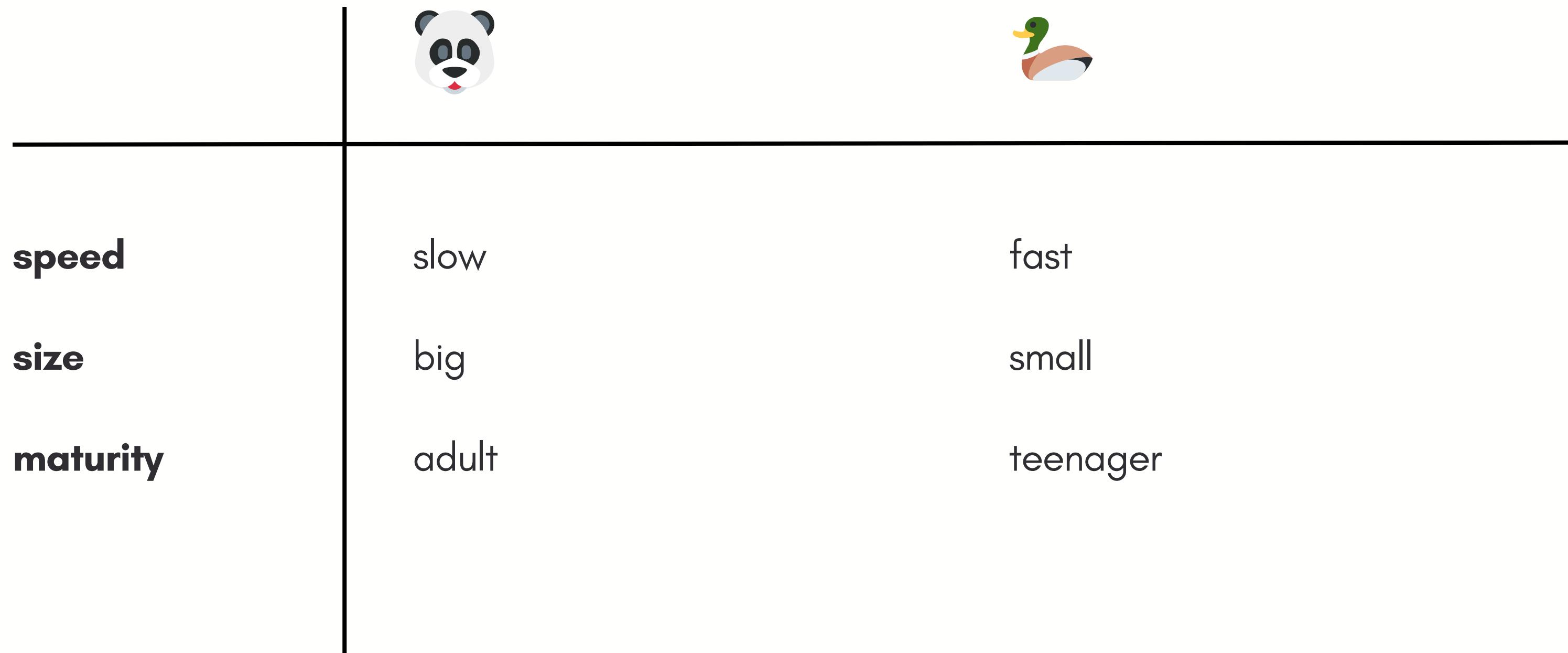
COMPARISON



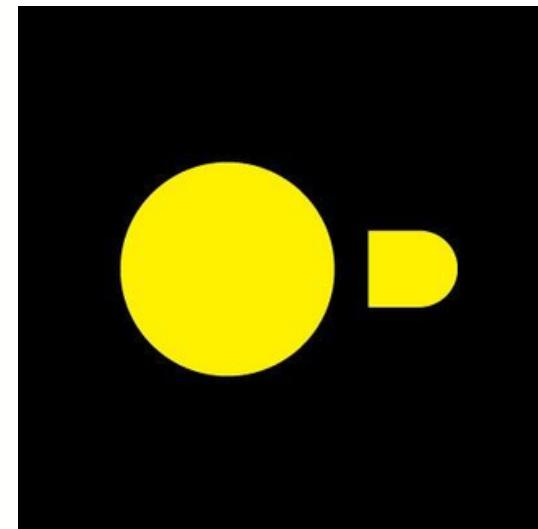
COMPARISON



COMPARISON



MY_TOOLBOX_ENRICHED_BY_DUCKDB



Thank You



