

DuckLake

The SQL-Powered Lakehouse Format for the Rest of Us





Distributed?

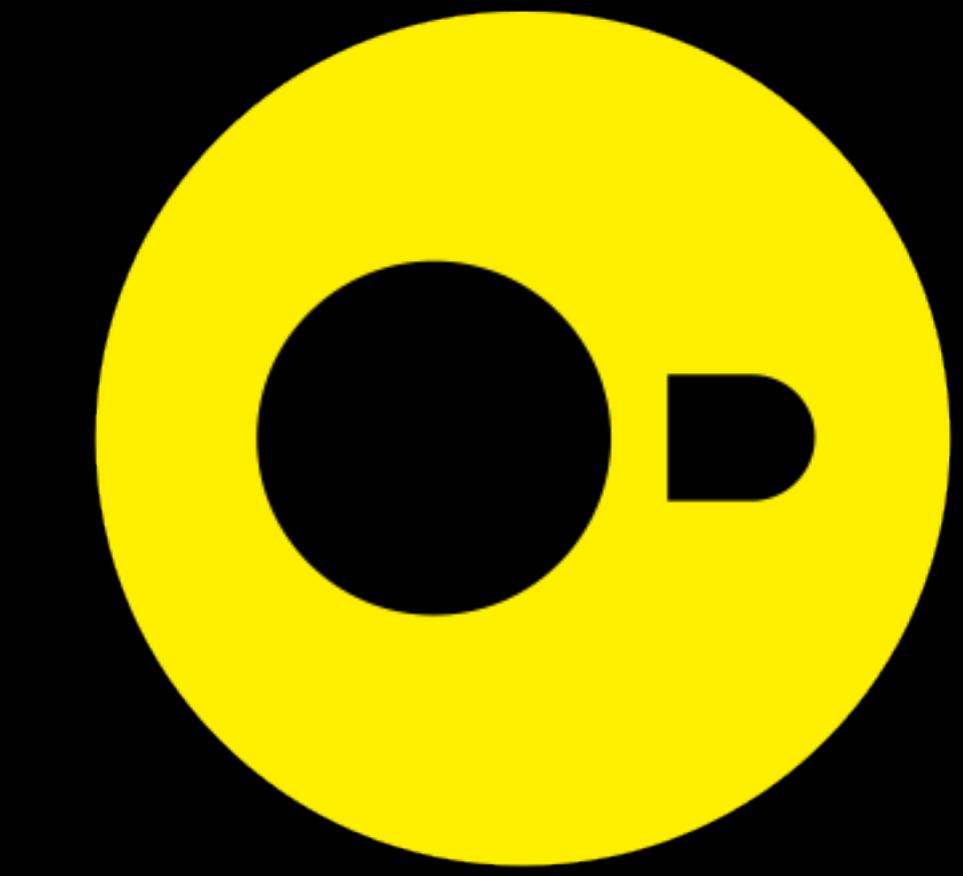
Architecture-Independent Distributed Query Processing

Hannes Fabian Mühleisen

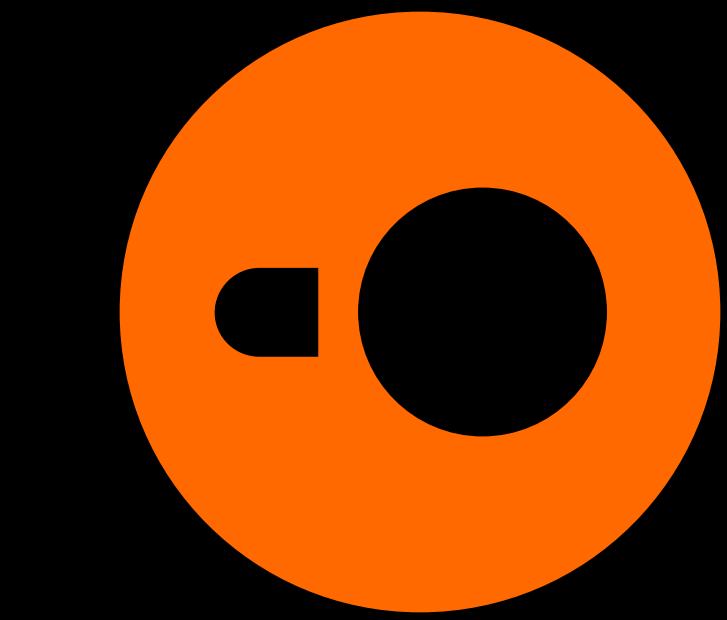


Dissertation submitted to the
Department of Mathematics and Computer Science
Freie Universität Berlin
MMXII

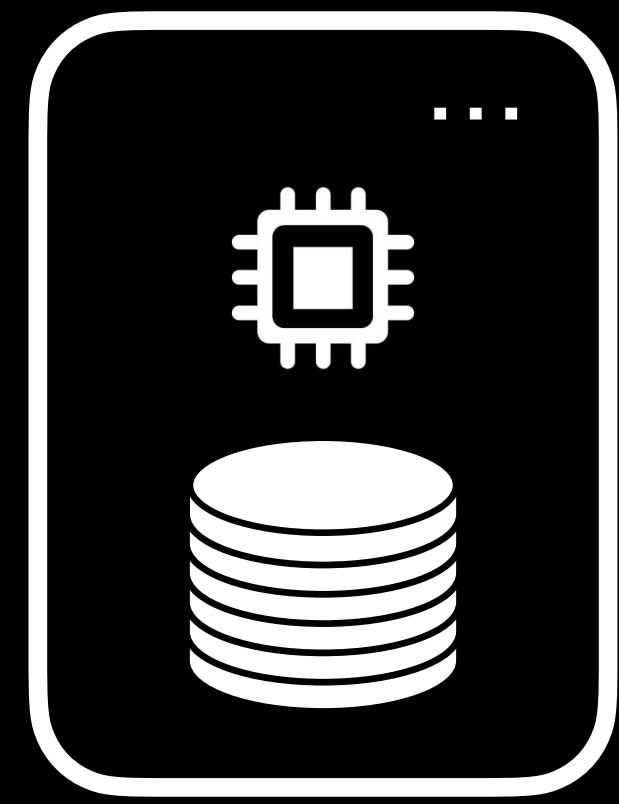


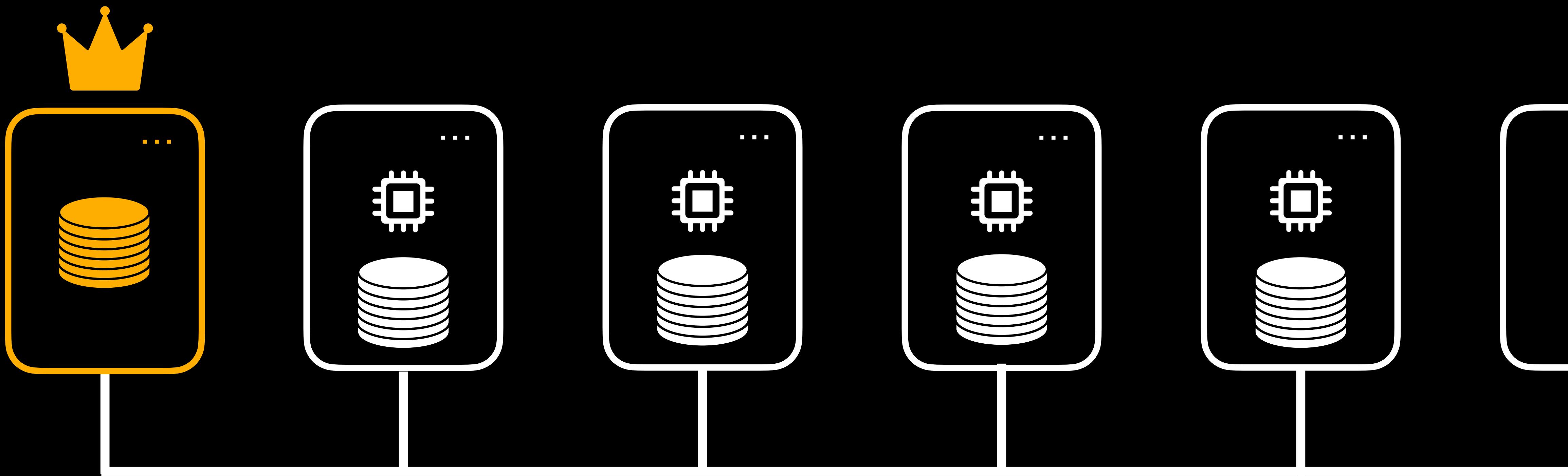


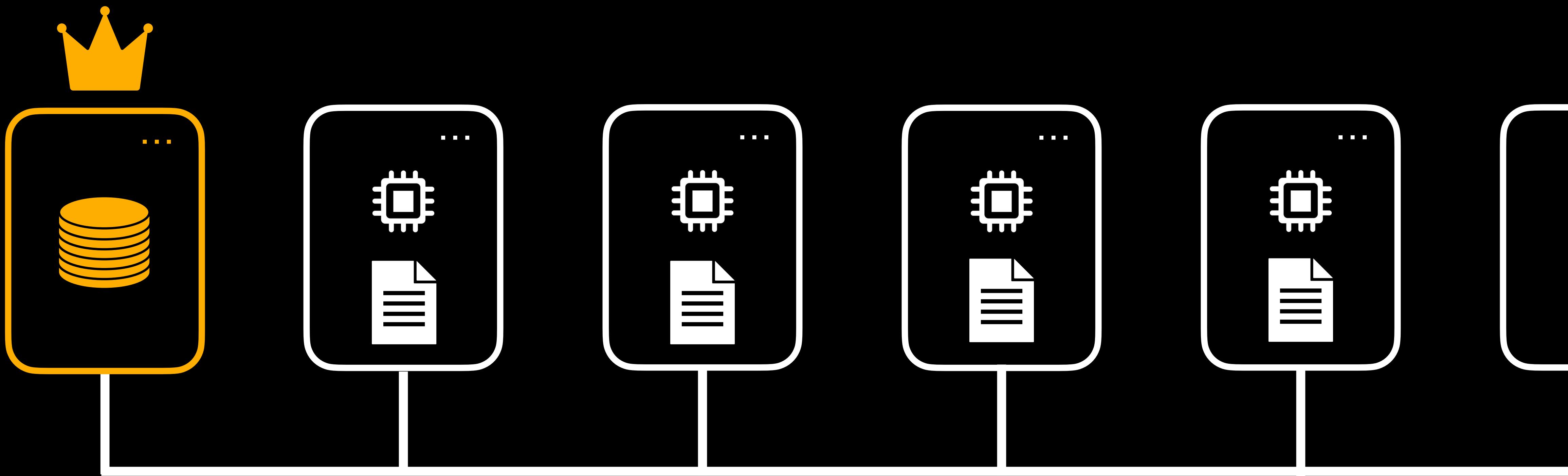
DuckDB

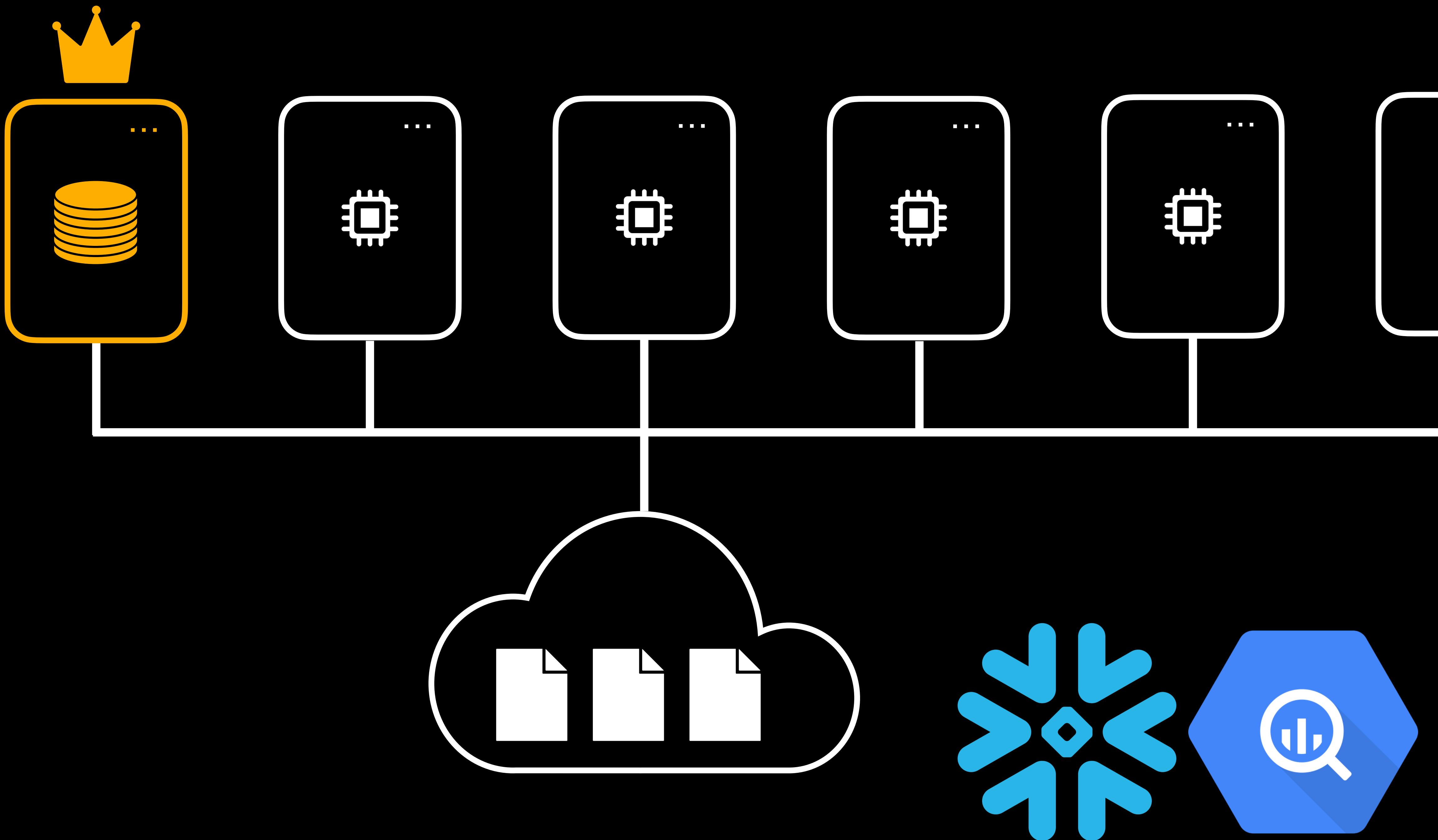


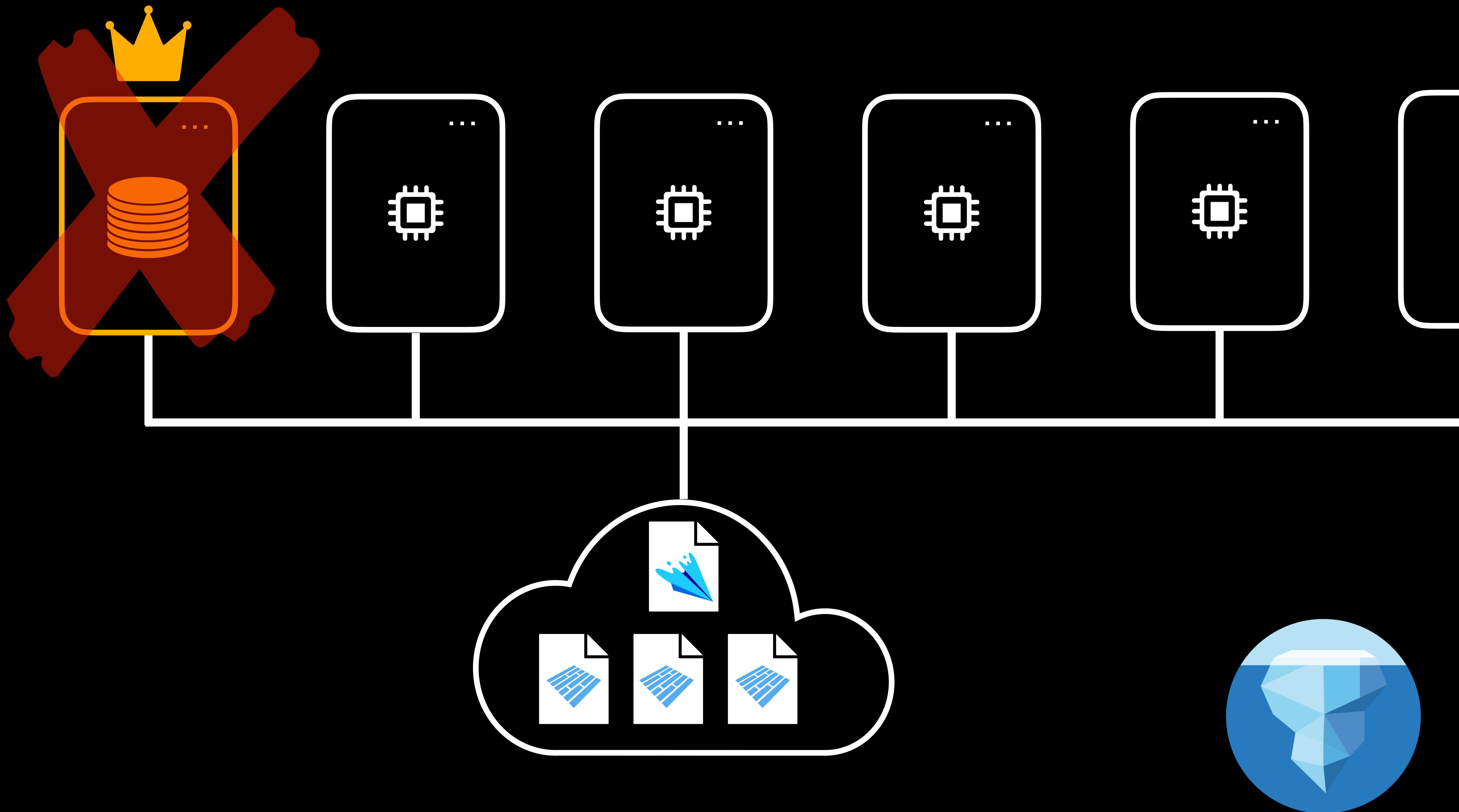


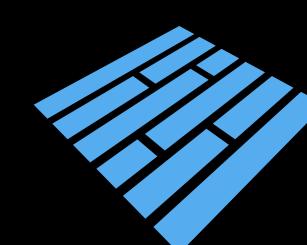
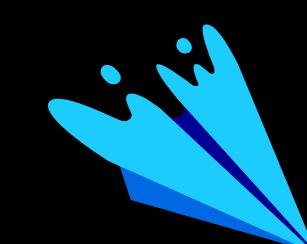
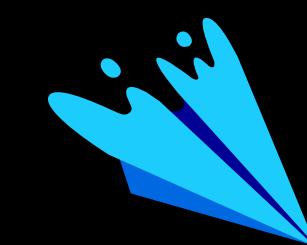
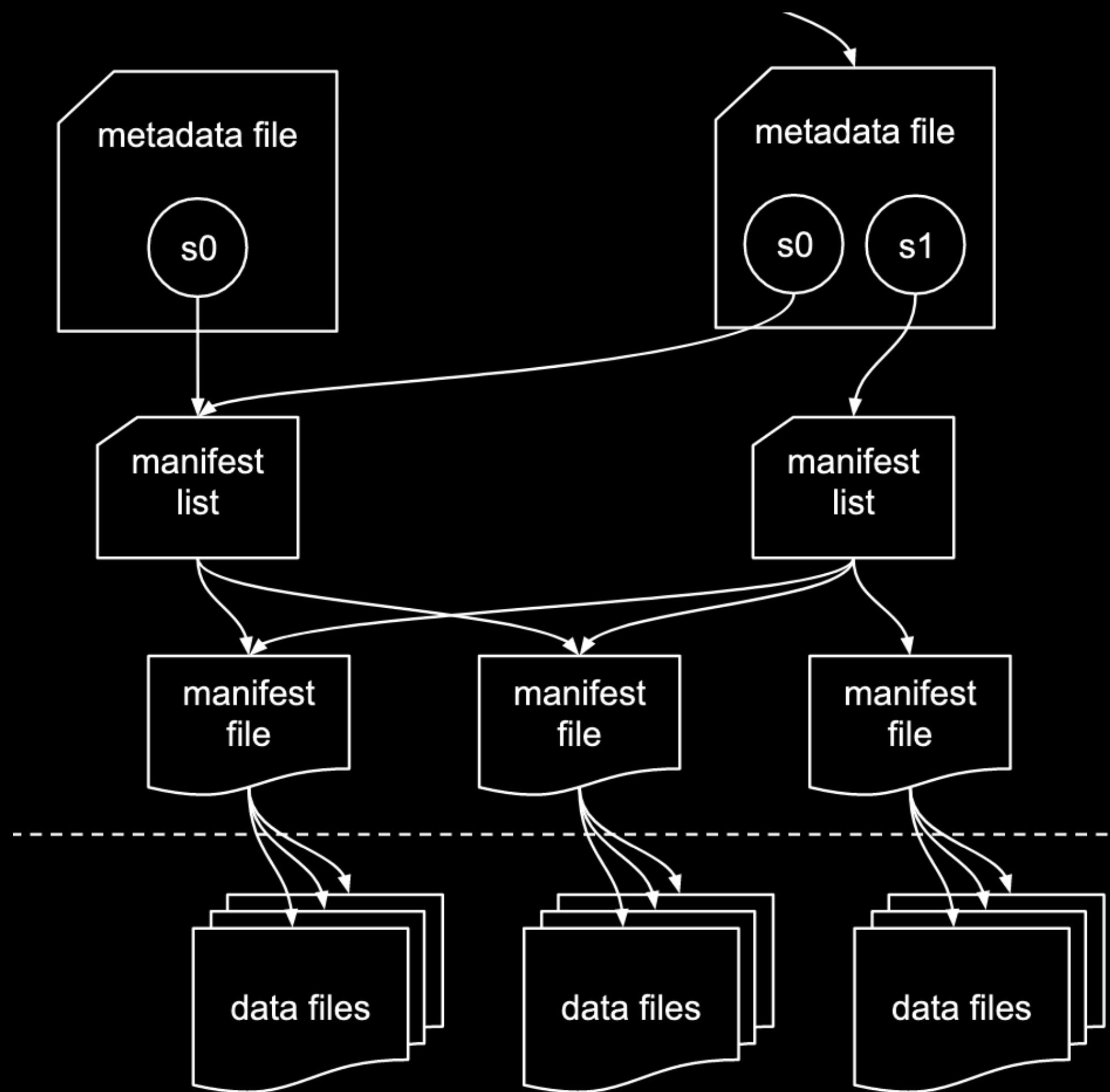




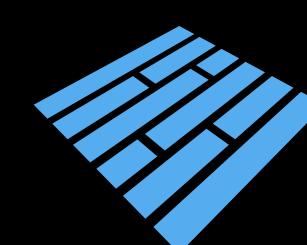
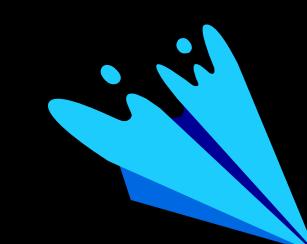
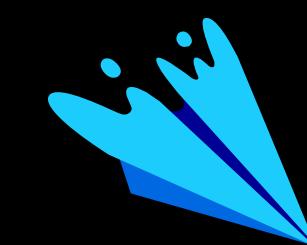
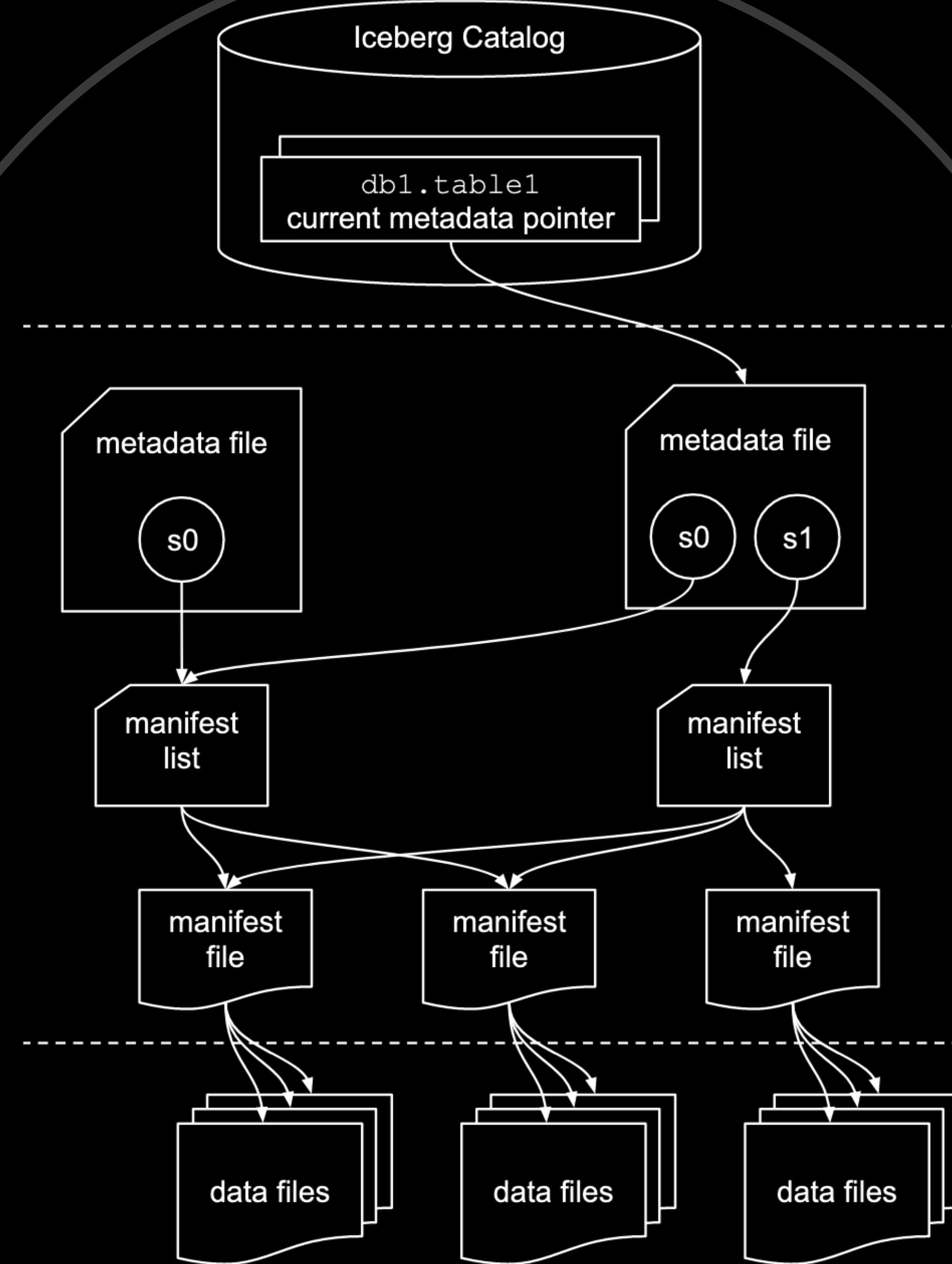
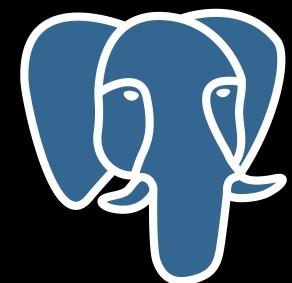




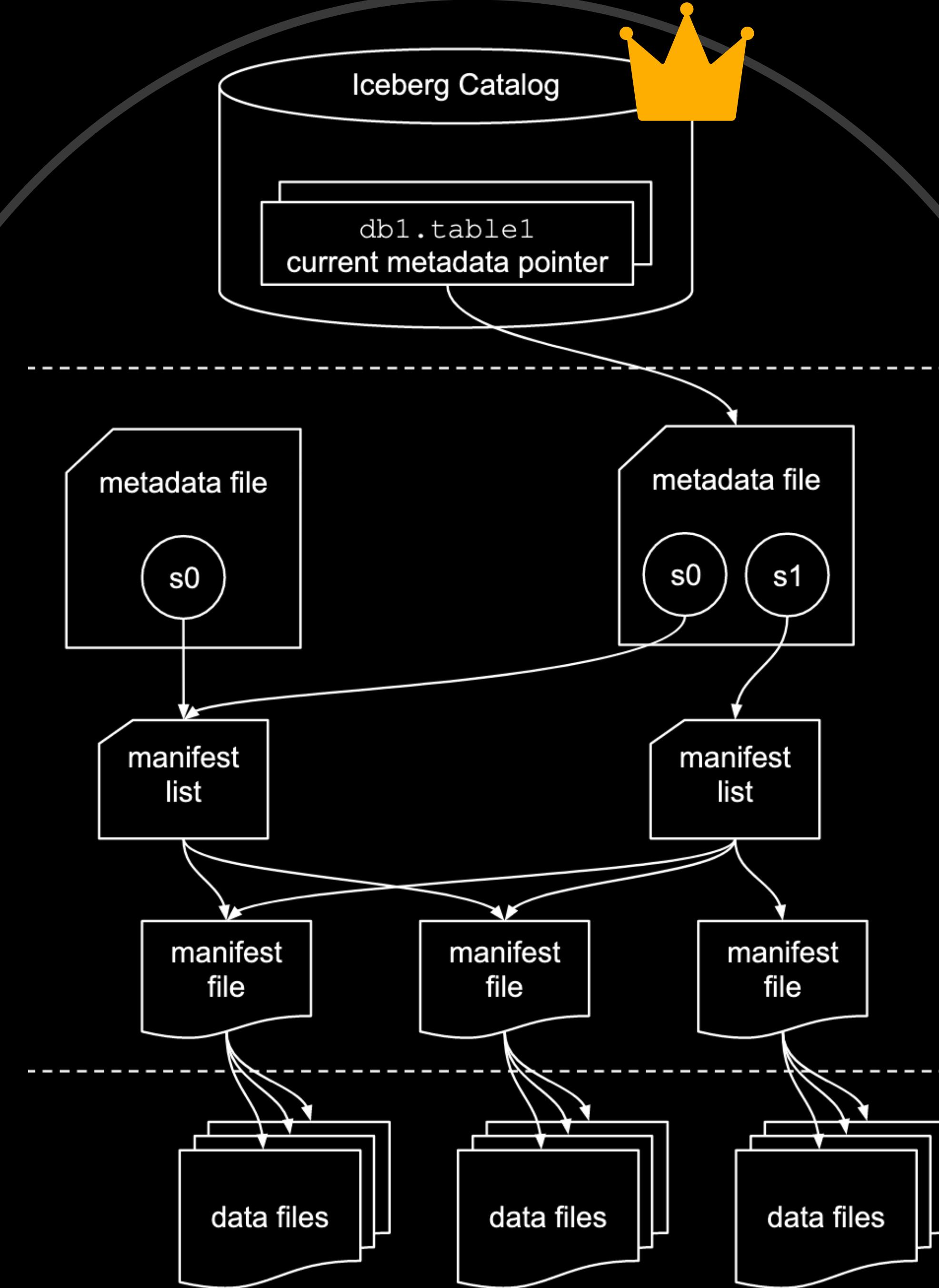
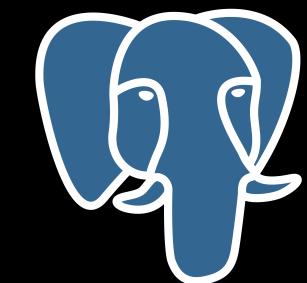


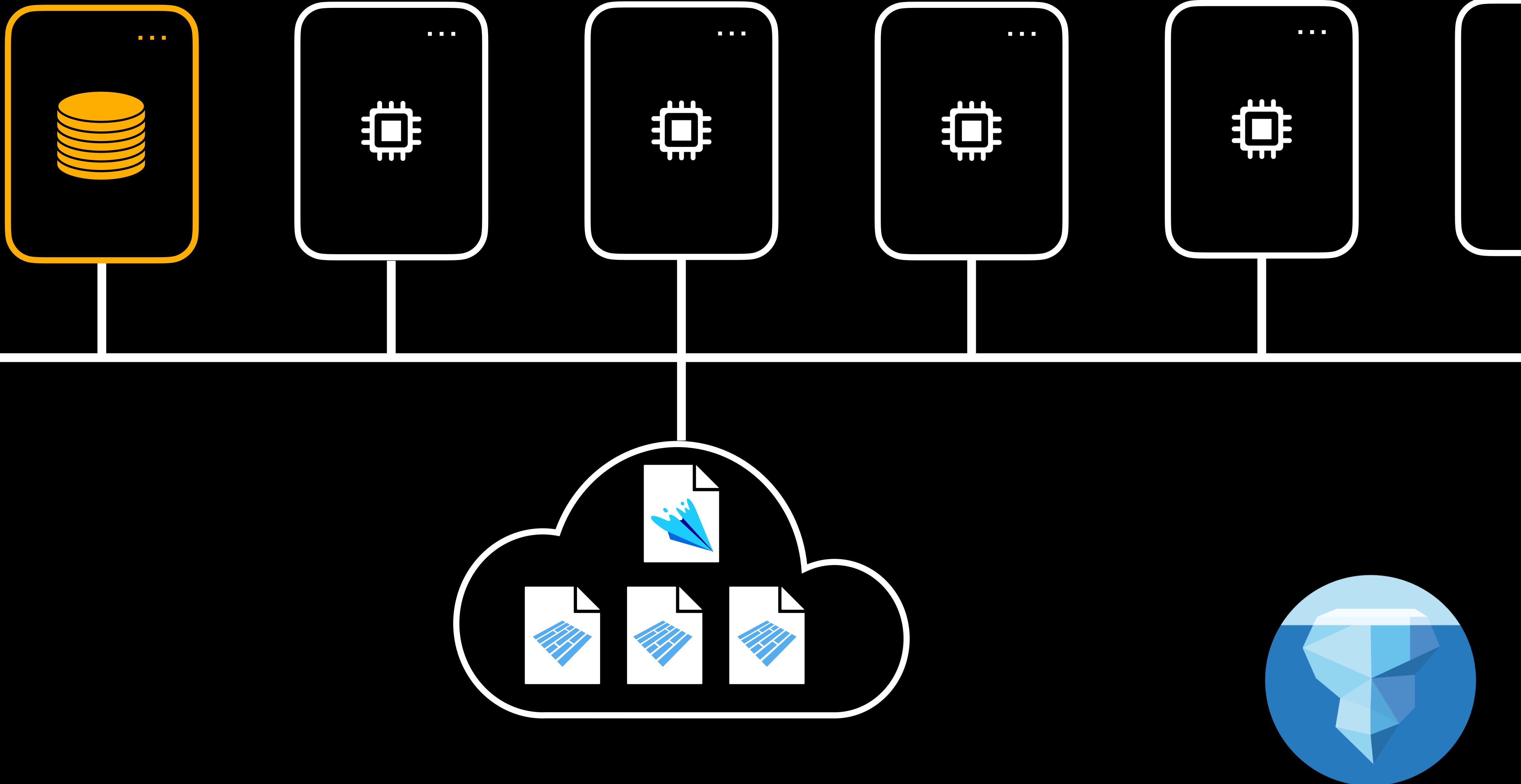


Java
REST
Postgres

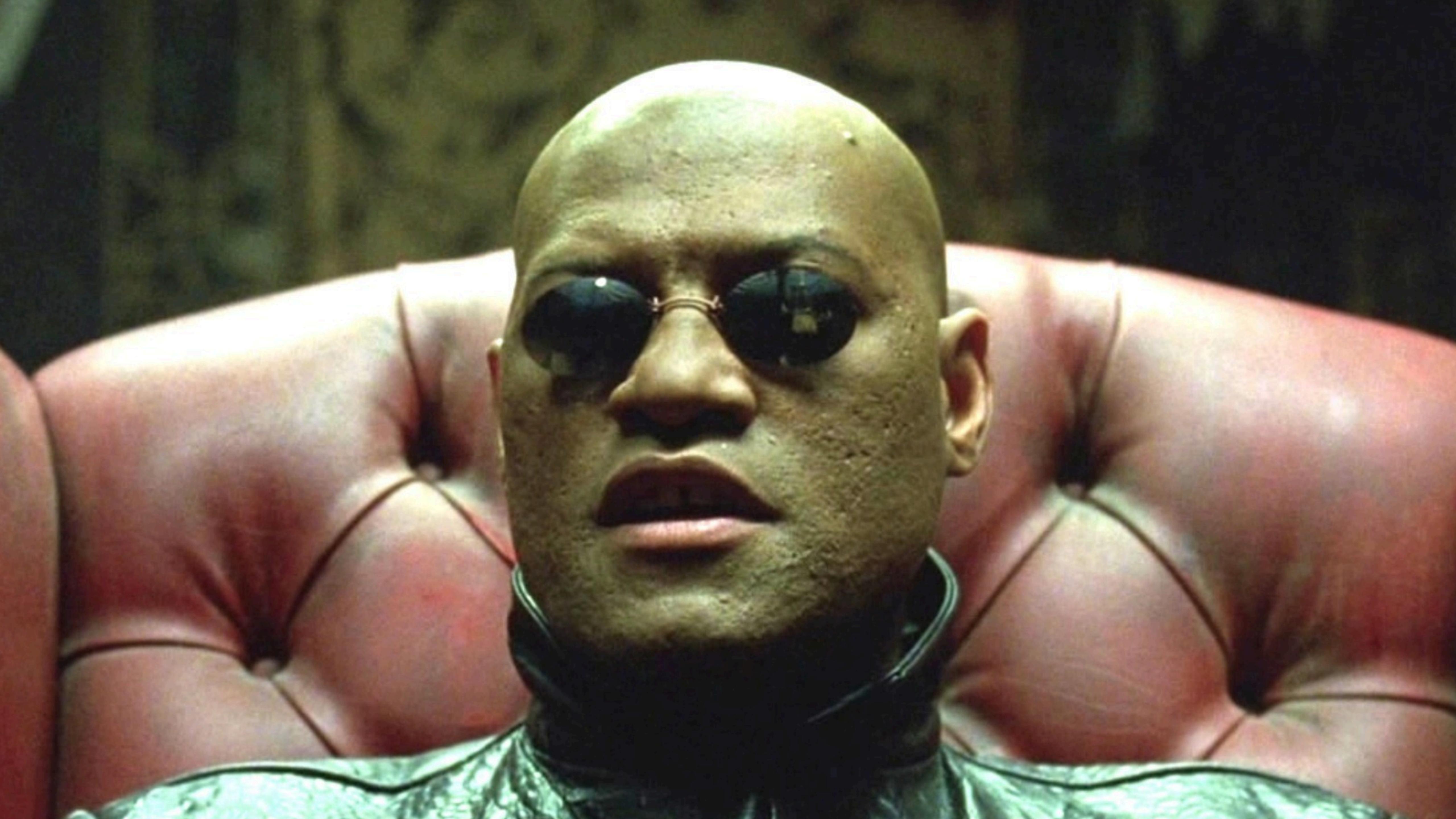


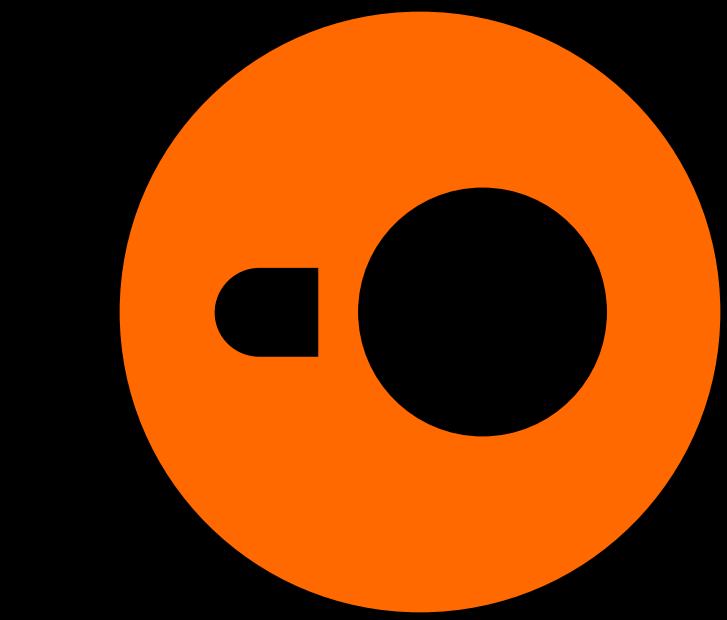
Java REST Postgres





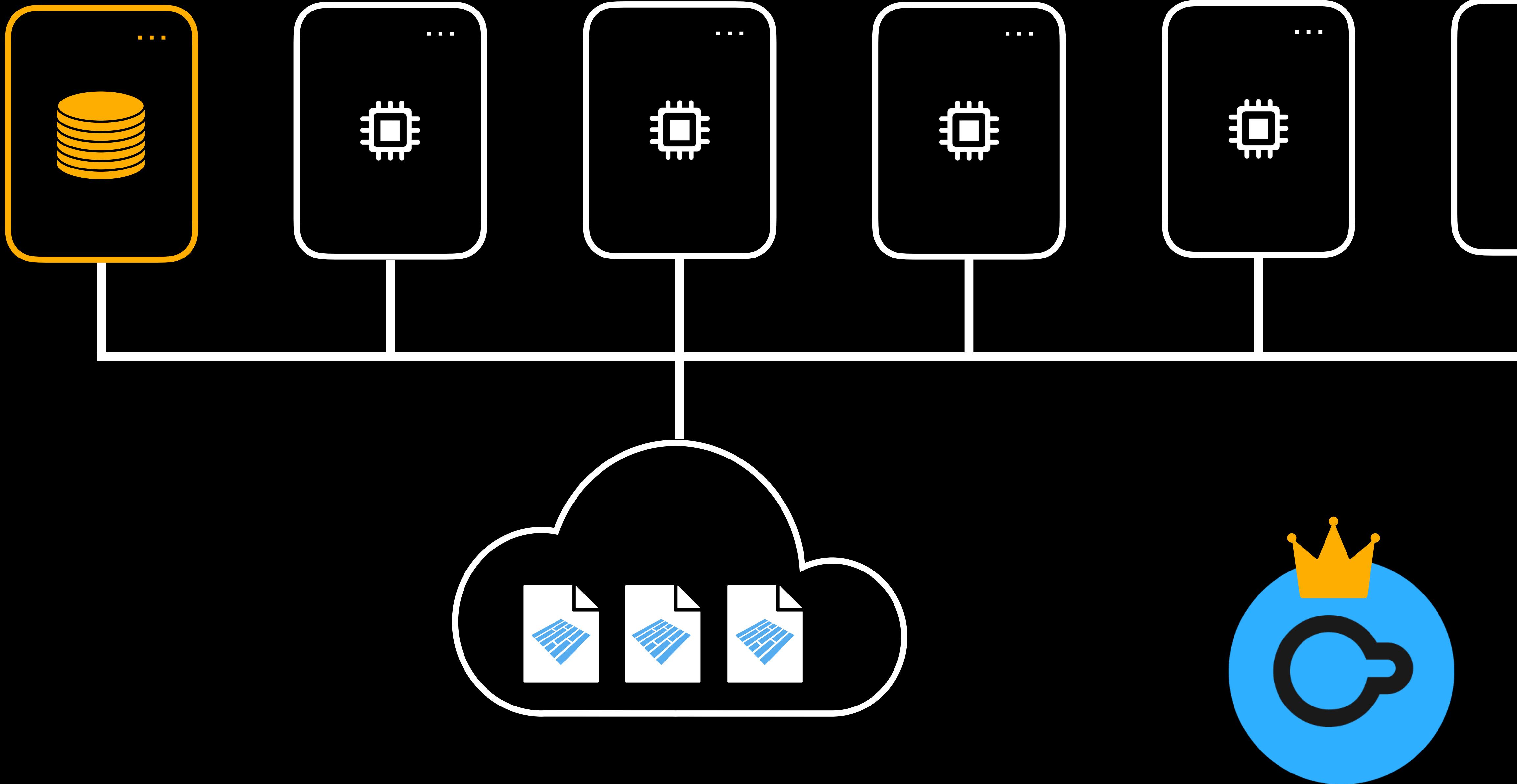








DuckLake



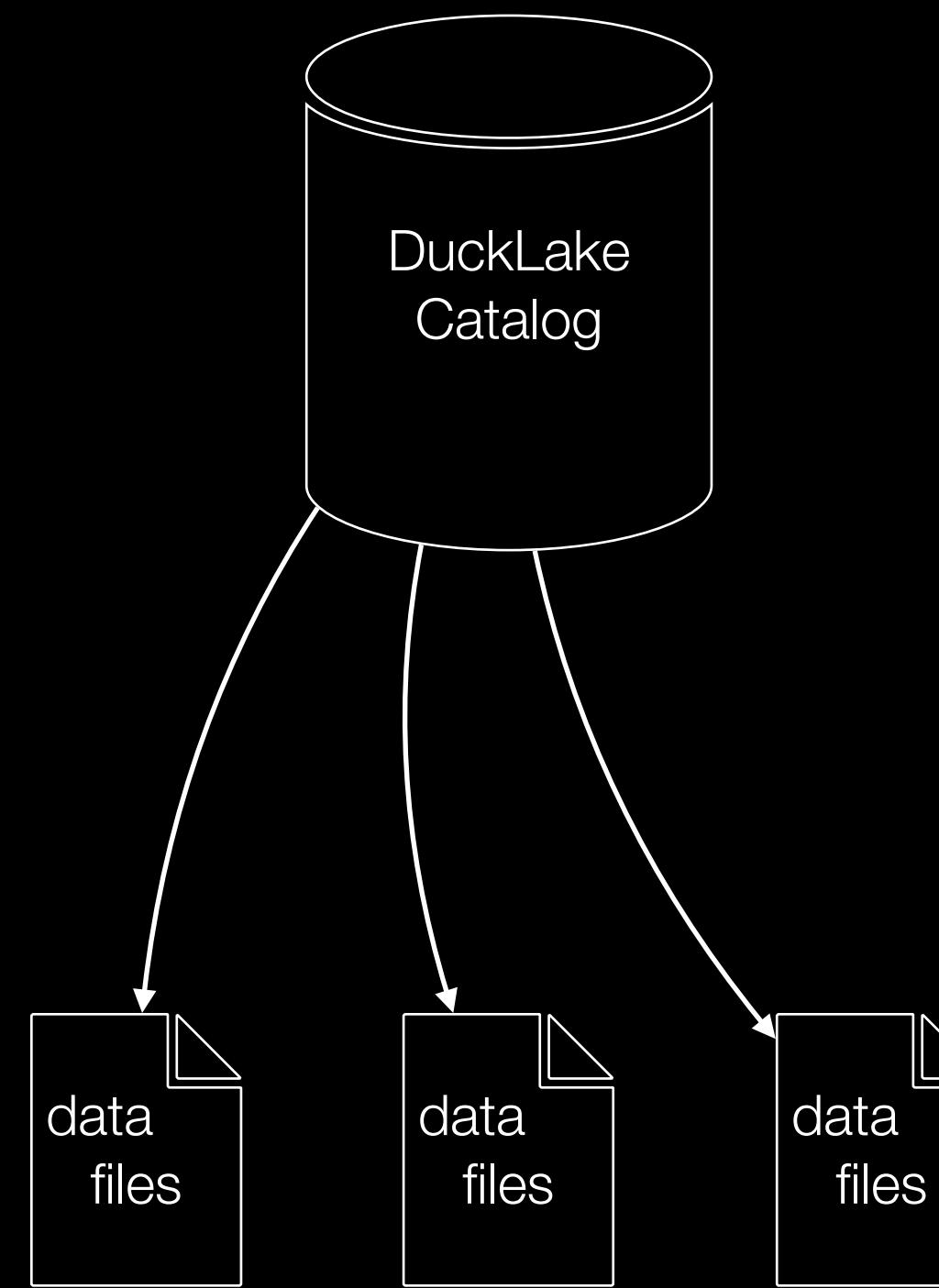
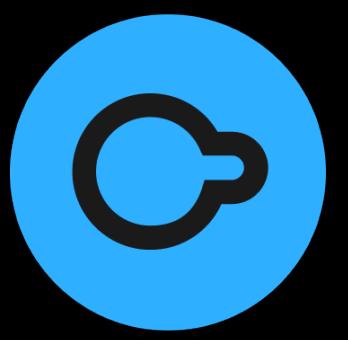
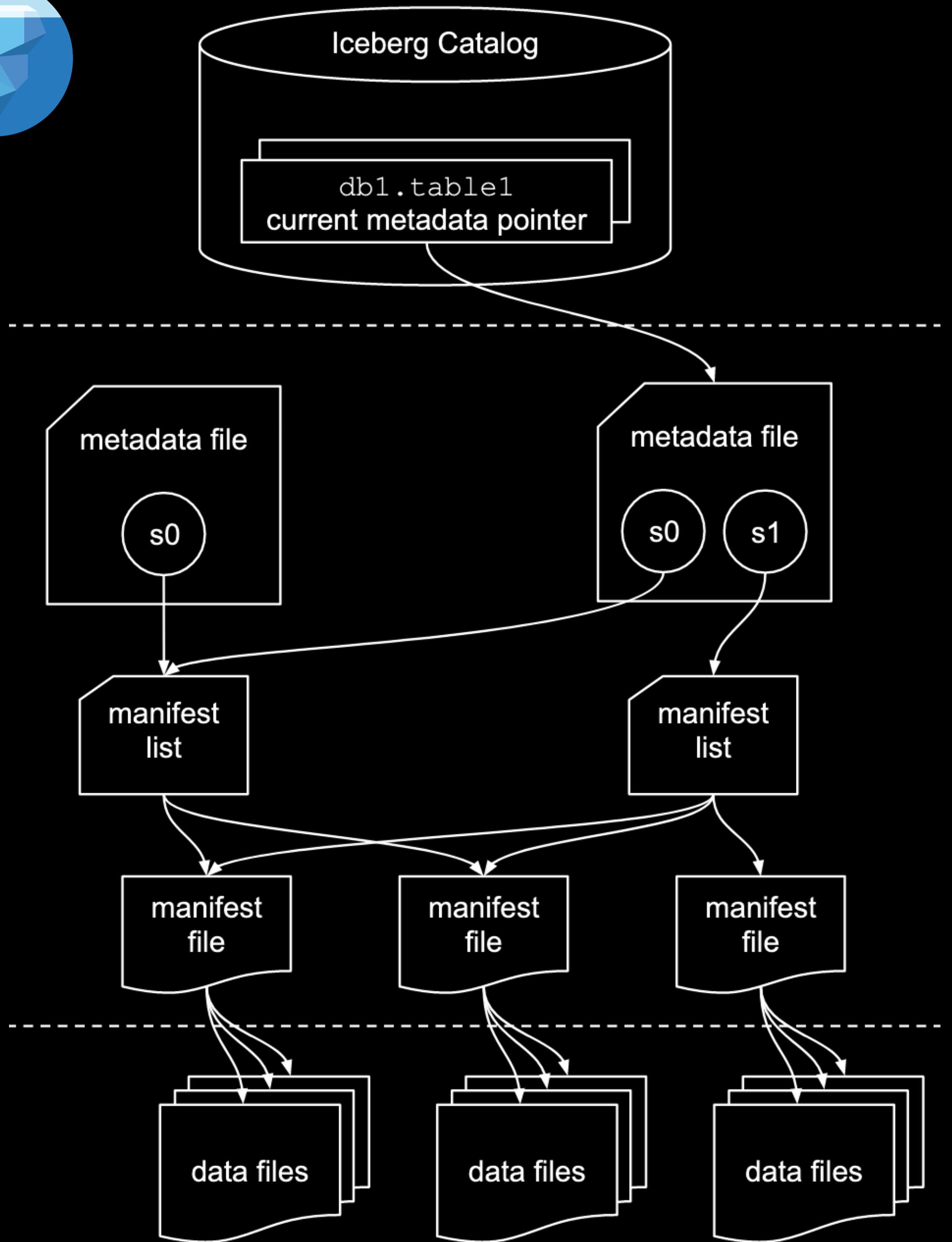
Simplicity
Scalability
Speed

Simplicity

```
$ curl install.duckdb.org | bash
```

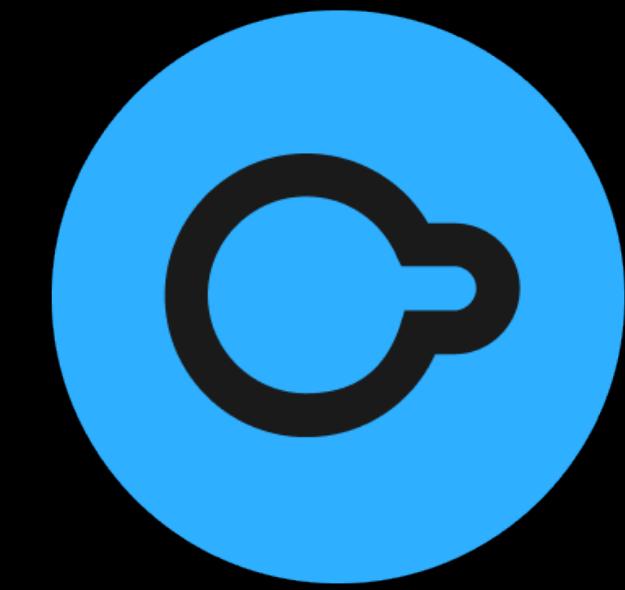
```
$ curl install.duckdb.org | bash  
$ ~/.duckdb/cli/latest/duckdb
```

```
$ curl install.duckdb.org | bash  
$ ~/.duckdb/cli/latest/duckdb  
▶ ATTACH 'ducklake:metadata.ducklake' AS dl1;
```

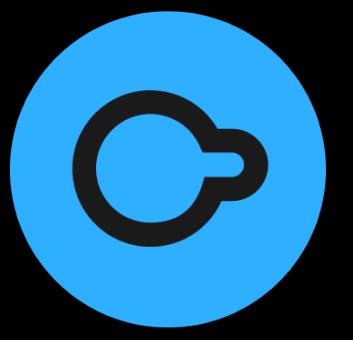


Database





Standard ◊
Implementation



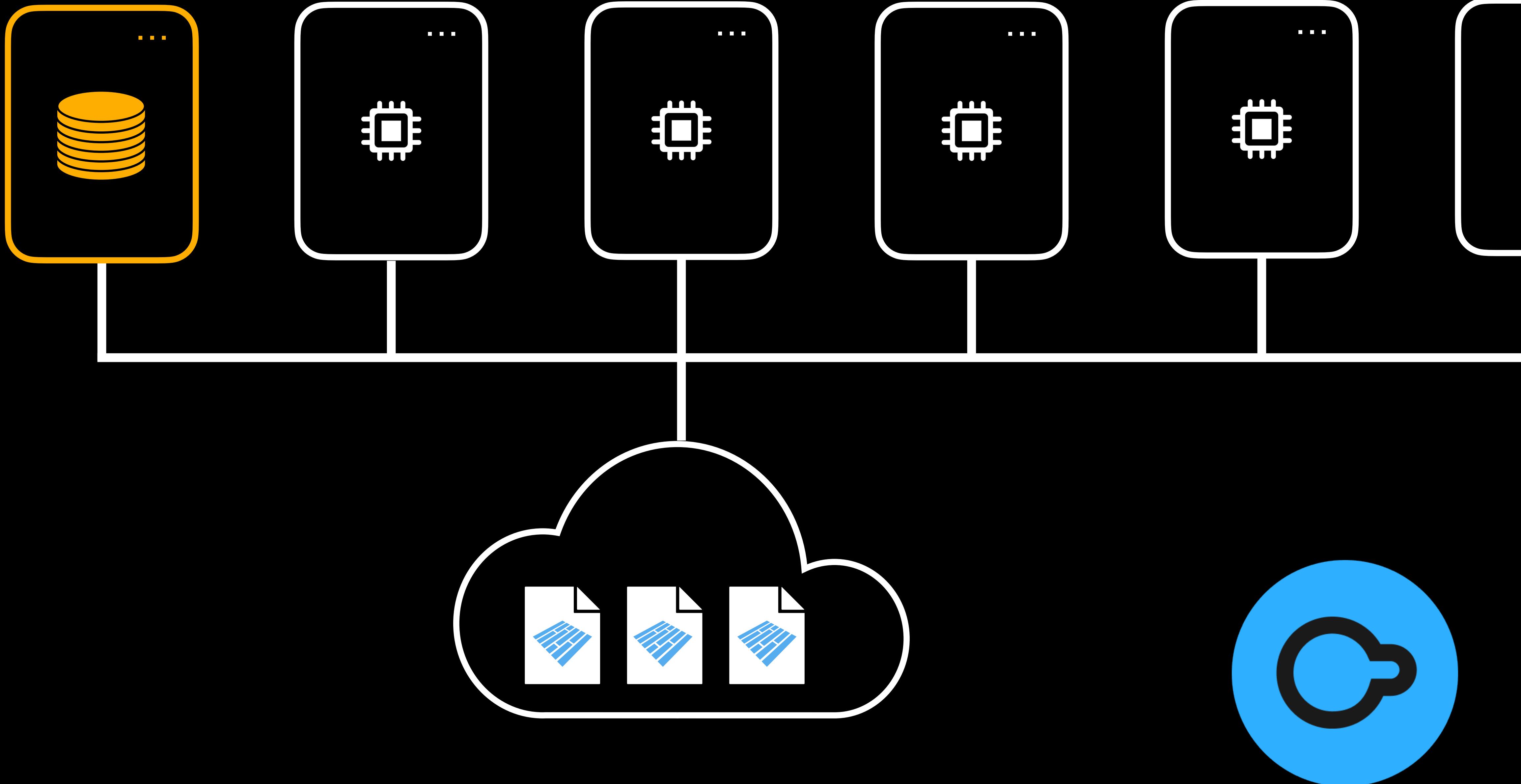
- Postgres
- Docker
- Polaris
- Java™
- REST
- JSON
- AVRO
- Parquet
- S3
- Spark

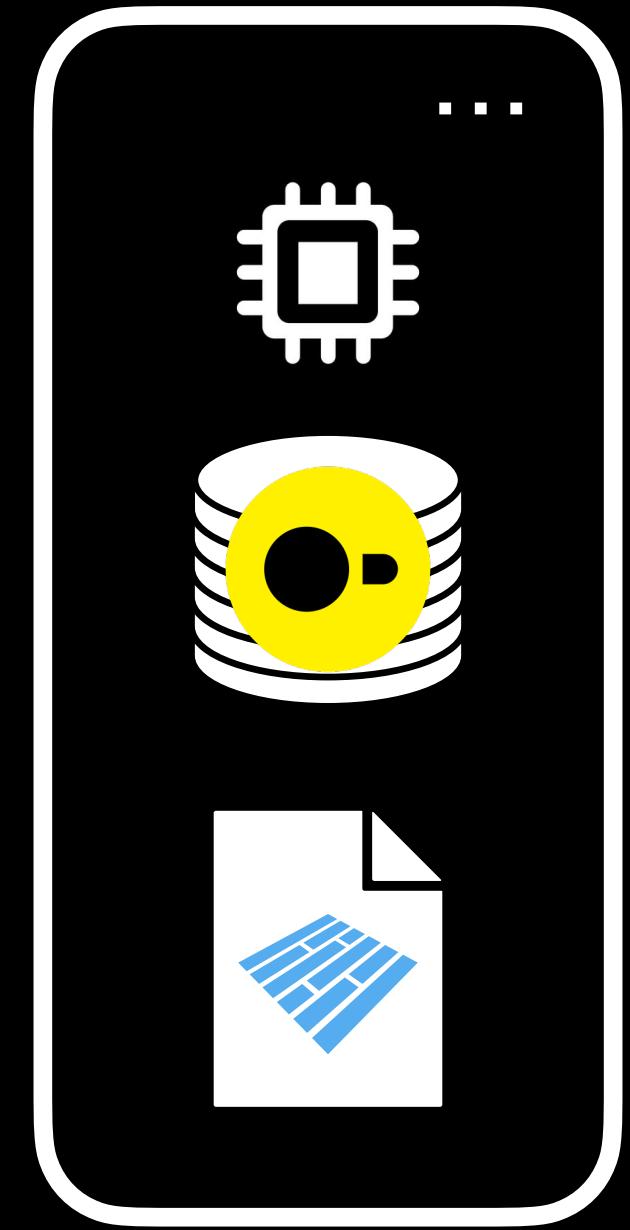
- DuckDB
- (Postgres)
- (S3)

Scalability







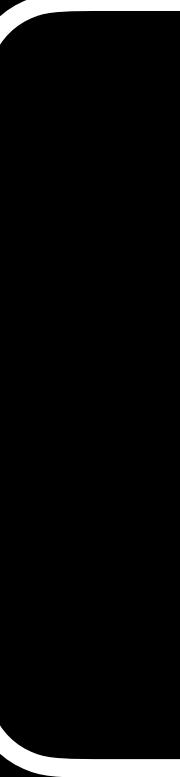
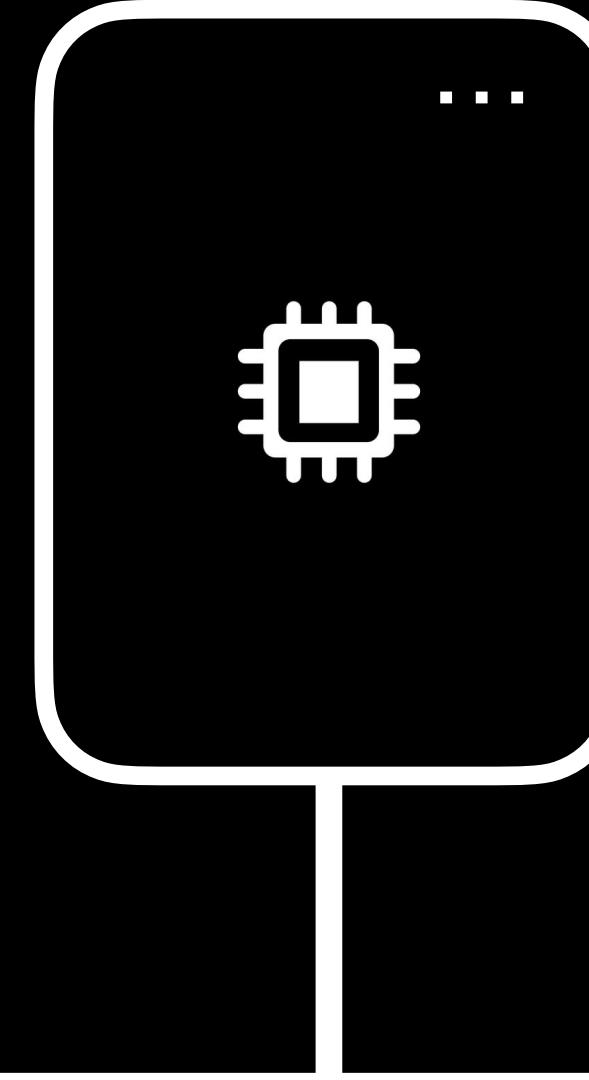
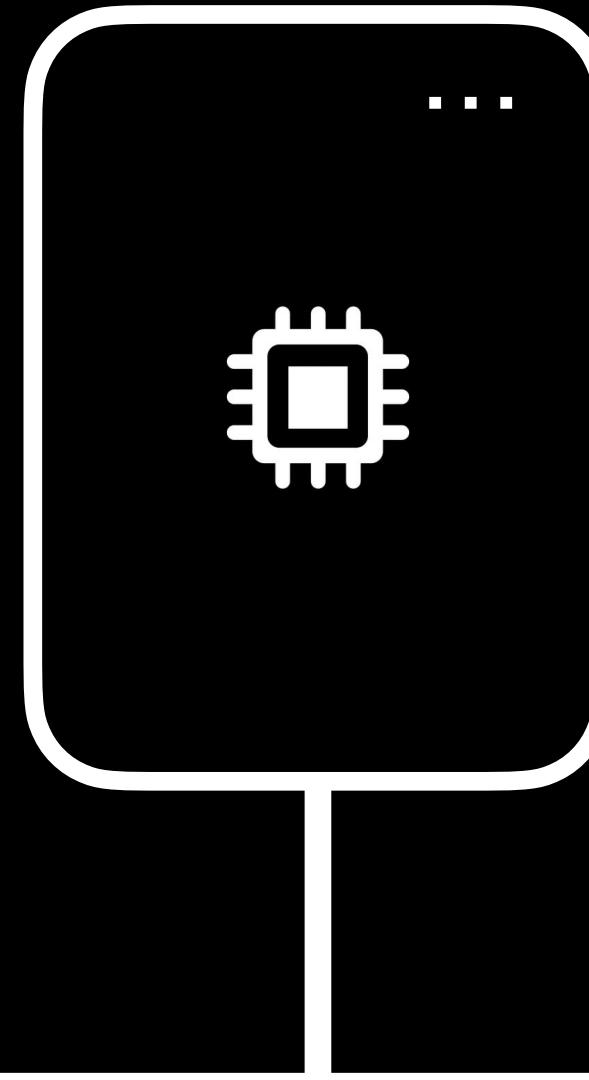
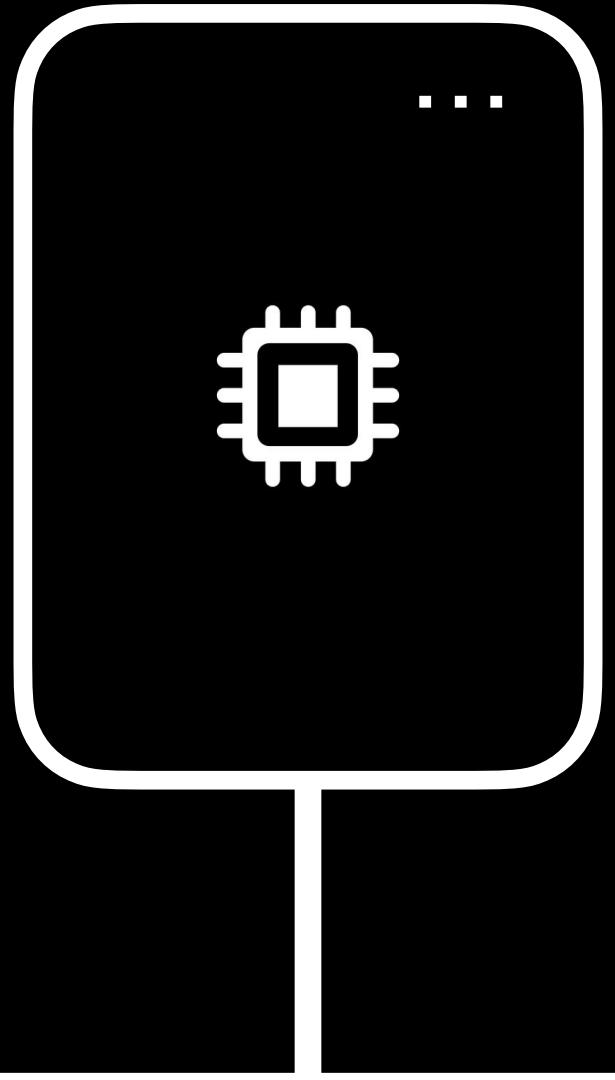
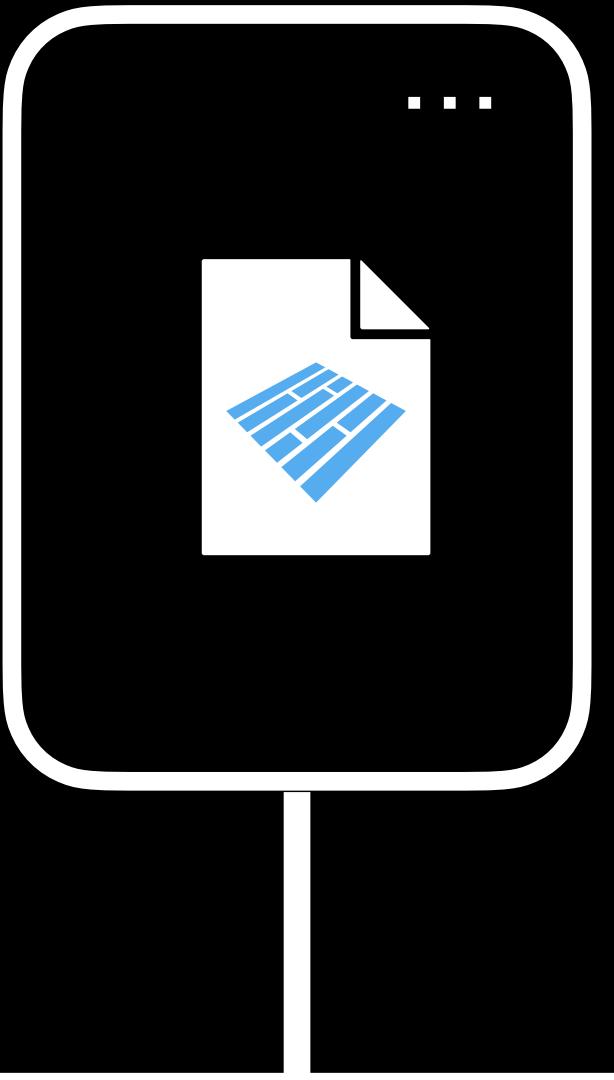


ATTACH 'ducklake:metadata1.ducklake' AS dl1;

metadata

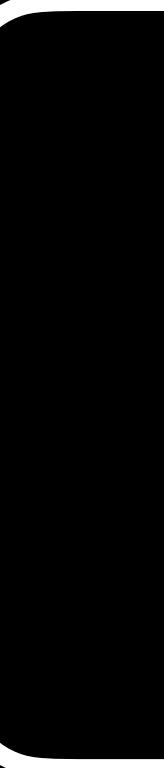
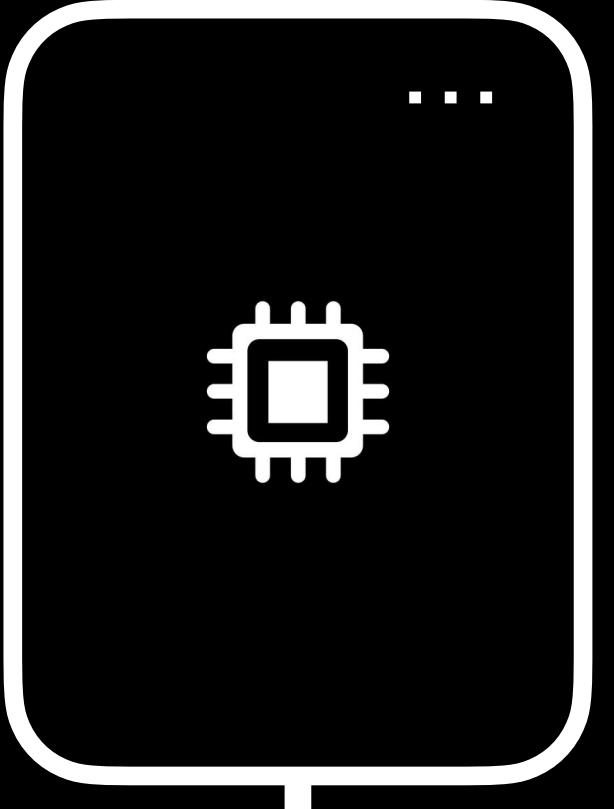
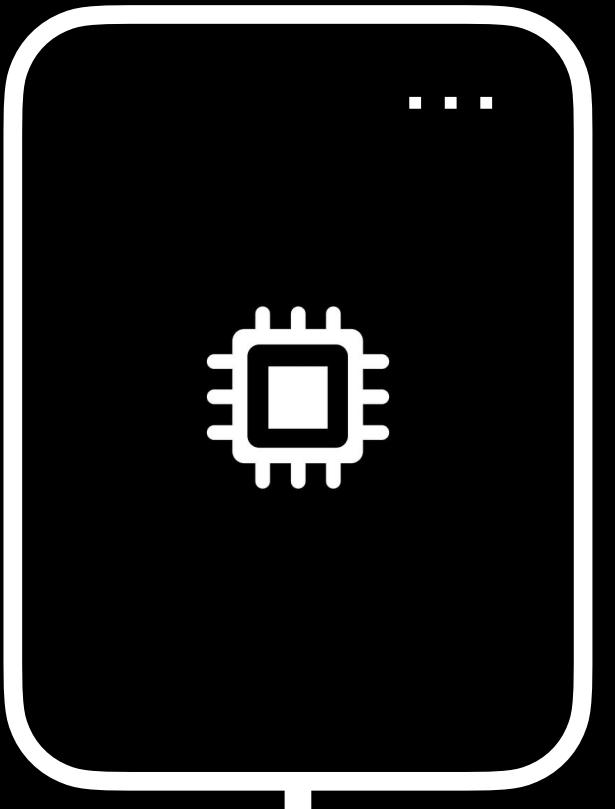
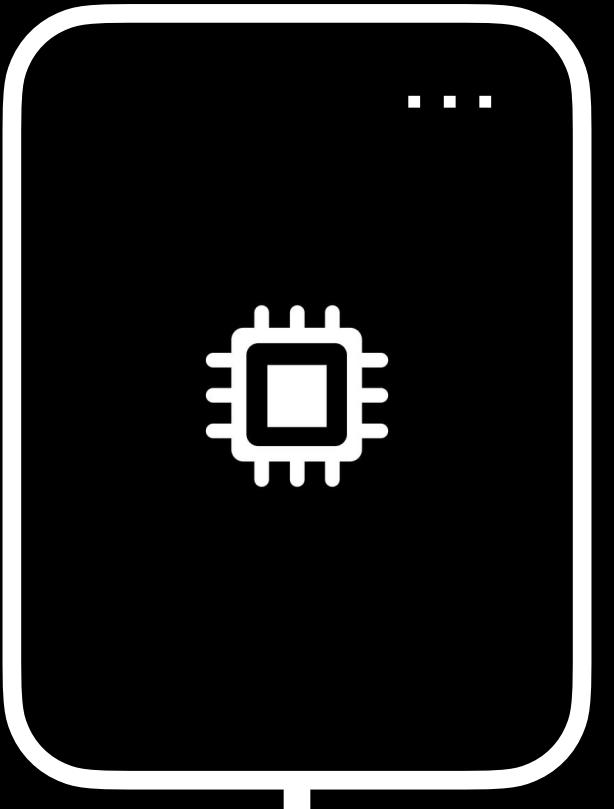
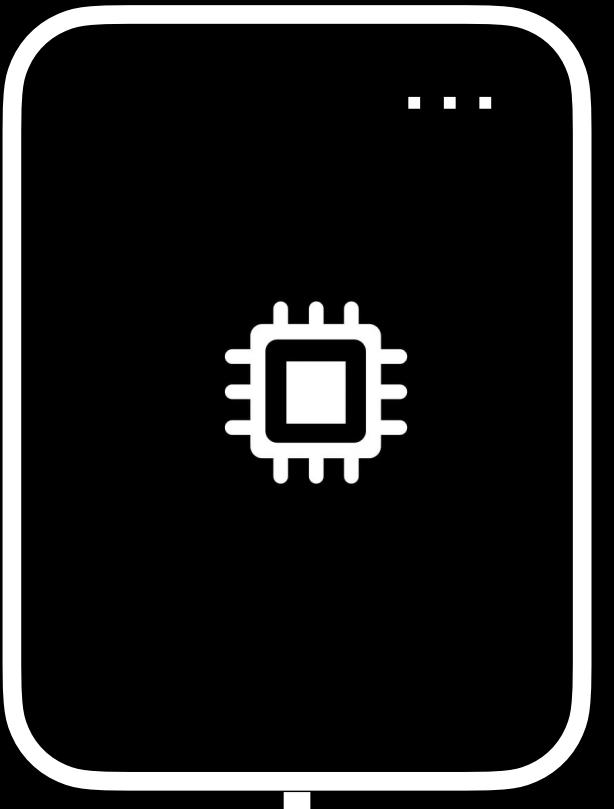
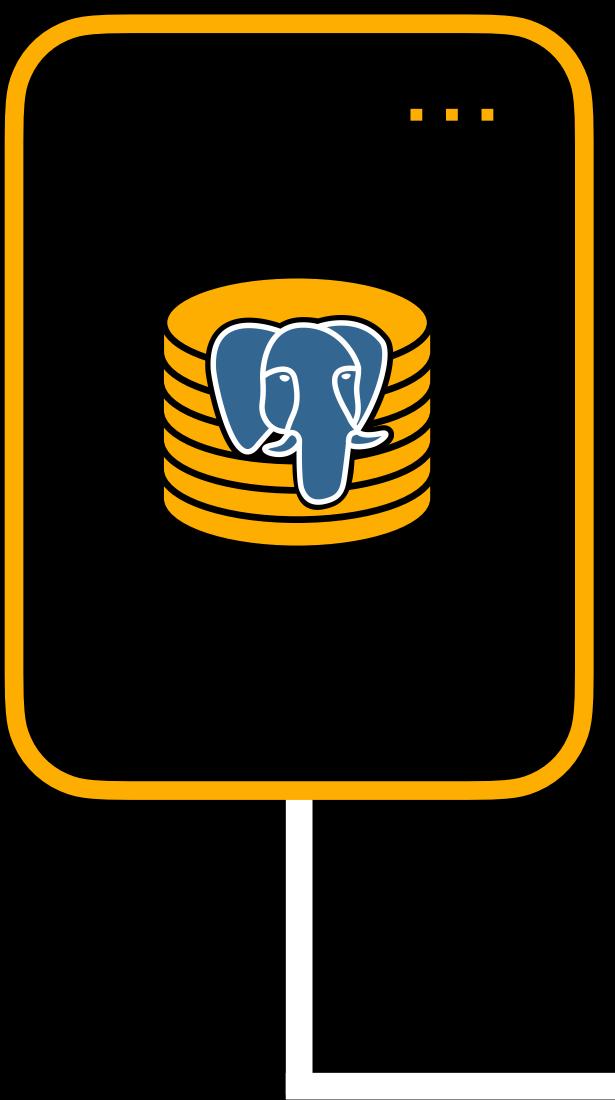


nfs

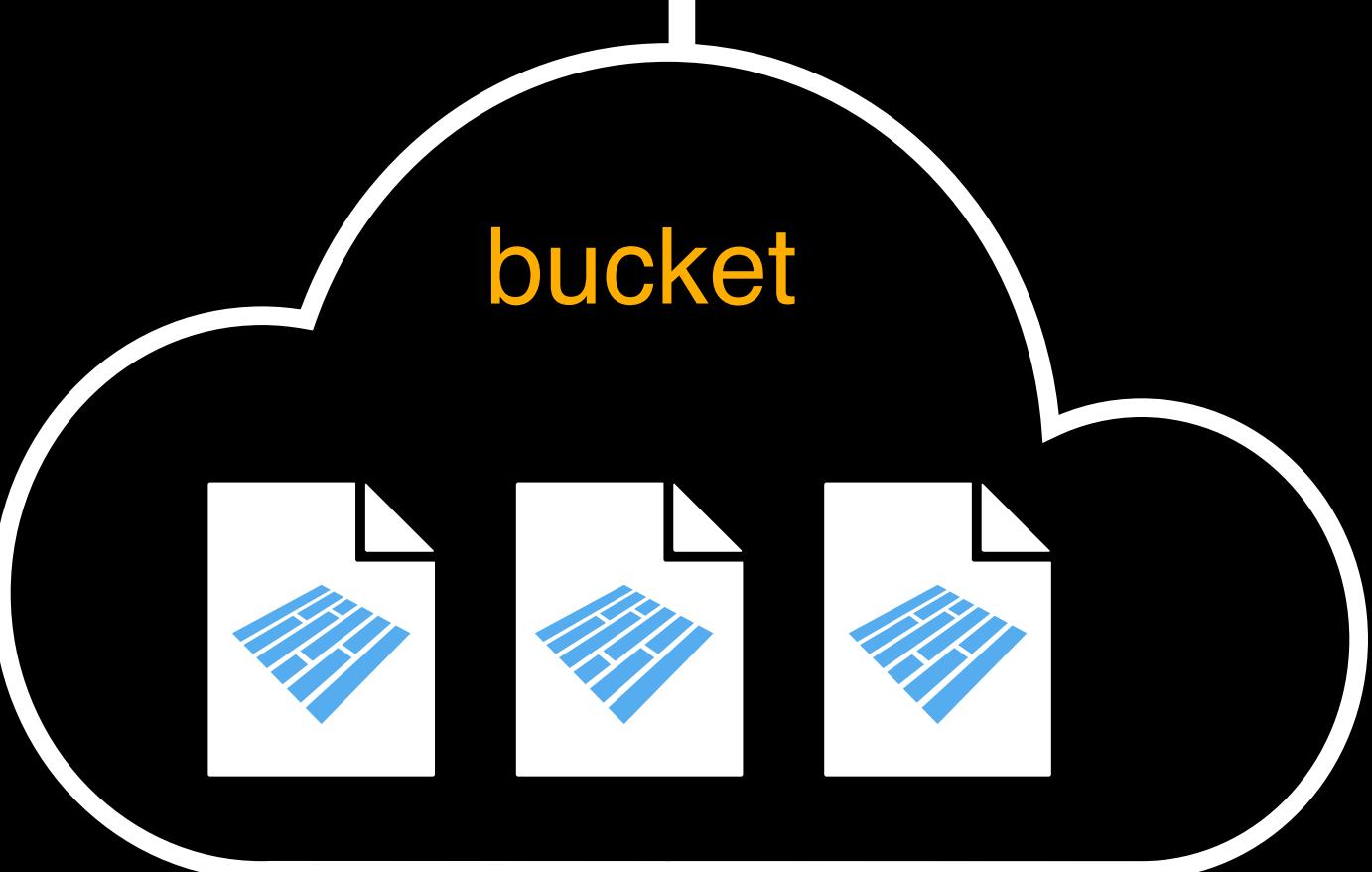


```
ATTACH 'ducklake:postgres:host=metadata' AS dl  
(DATA_PATH '/nfs/data_files/');
```

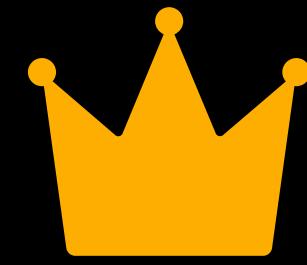
metadata



bucket



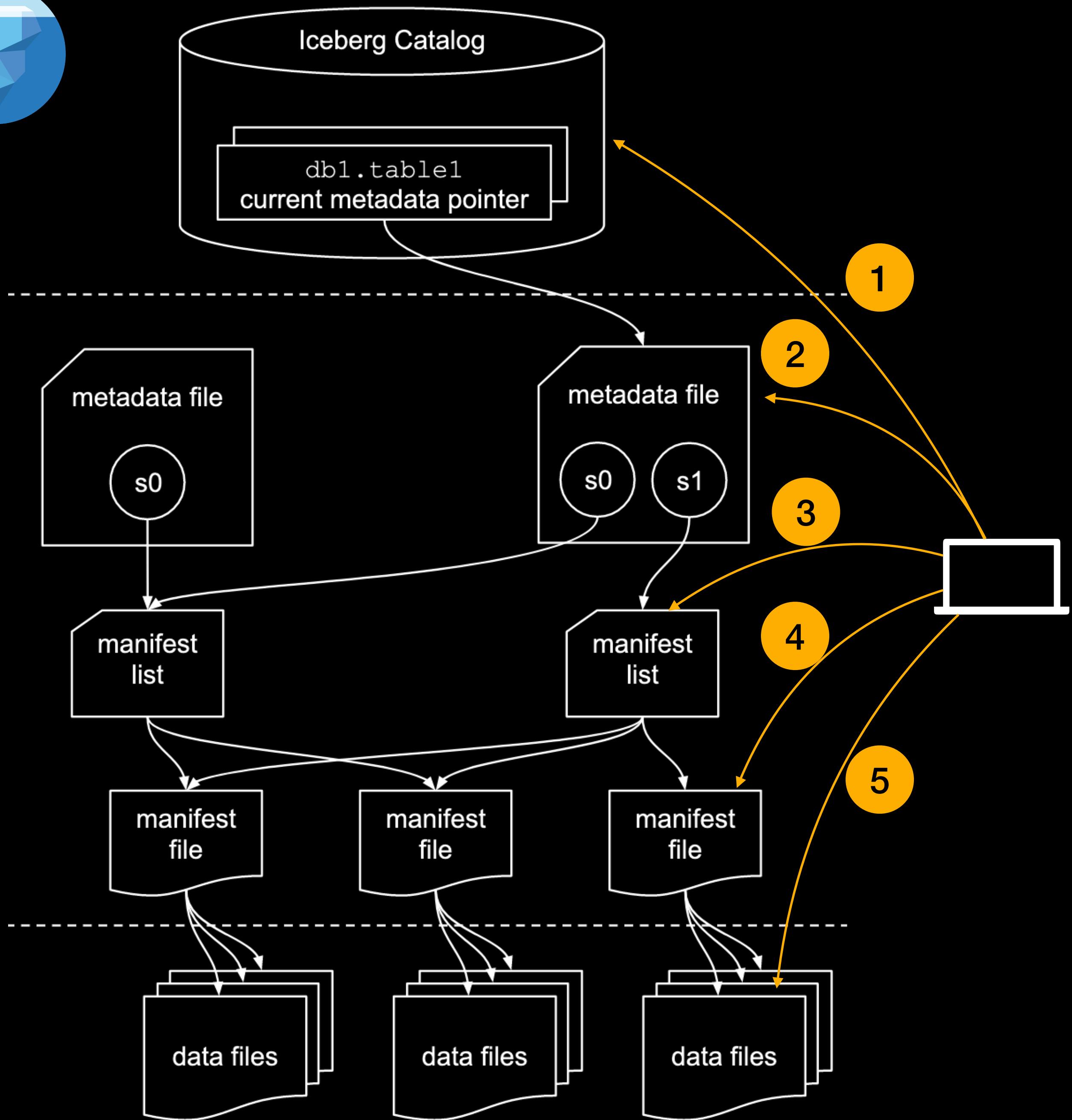
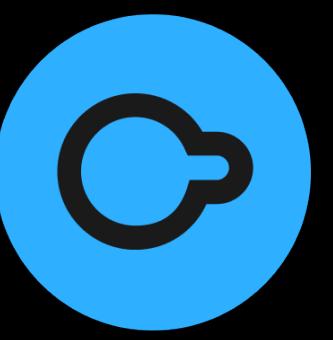
```
ATTACH 'ducklake:postgres:host=metadata'  
AS dl (DATA_PATH 's3://bucket/...');
```



- 1 Petabyte of data
- 10 Columns
- 100 M snapshots, 10 years
 - 270 GB/day

Planning Time	DuckDB	Postgres
Insert	~5ms	~3ms
Select All Files	~0.8s	~200s
Select 1 Month	~0.25s	~0.25s
Select 1 Day	~0.25s	~0.25s

Speed



Fair Benchmarking Considered Difficult: Common Pitfalls In Database Performance Testing

Mark Raasveldt, Pedro Holanda, Tim Gubner & Hannes Mühleisen

Centrum Wiskunde & Informatica (CWI)

Amsterdam, The Netherlands

[raasveld,holanda,tgubner,hannes]@cwi.nl

ABSTRACT

Performance benchmarking is one of the most commonly used methods for comparing different systems or algorithms, both in scientific literature and in industrial publications. While performance measurements might seem objective on the surface, there are many different ways to influence benchmark results to favor one system over the other, either by accident or on purpose. In this paper, we perform a study of the common pitfalls in DBMS performance comparisons, and give advice on how they can be spotted and avoided so a fair performance comparison between systems can be made. We illustrate the common pitfalls with a series of mock benchmarks, which show large differences in performance where none should be present.

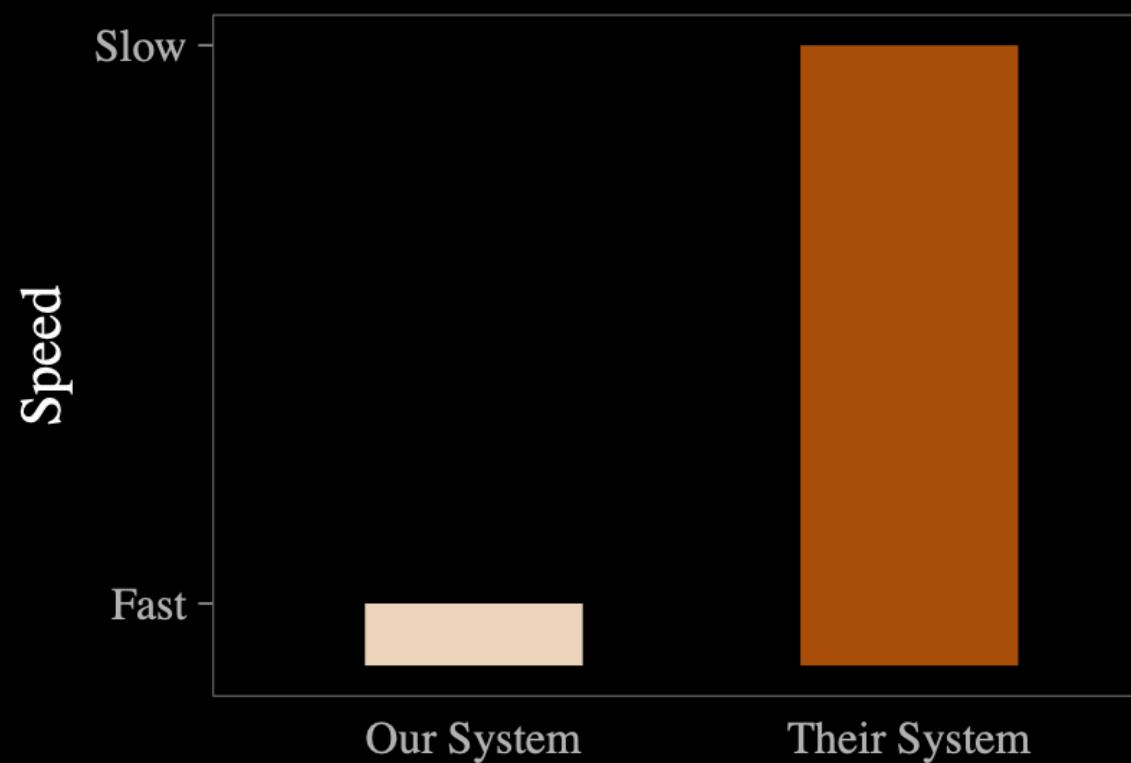
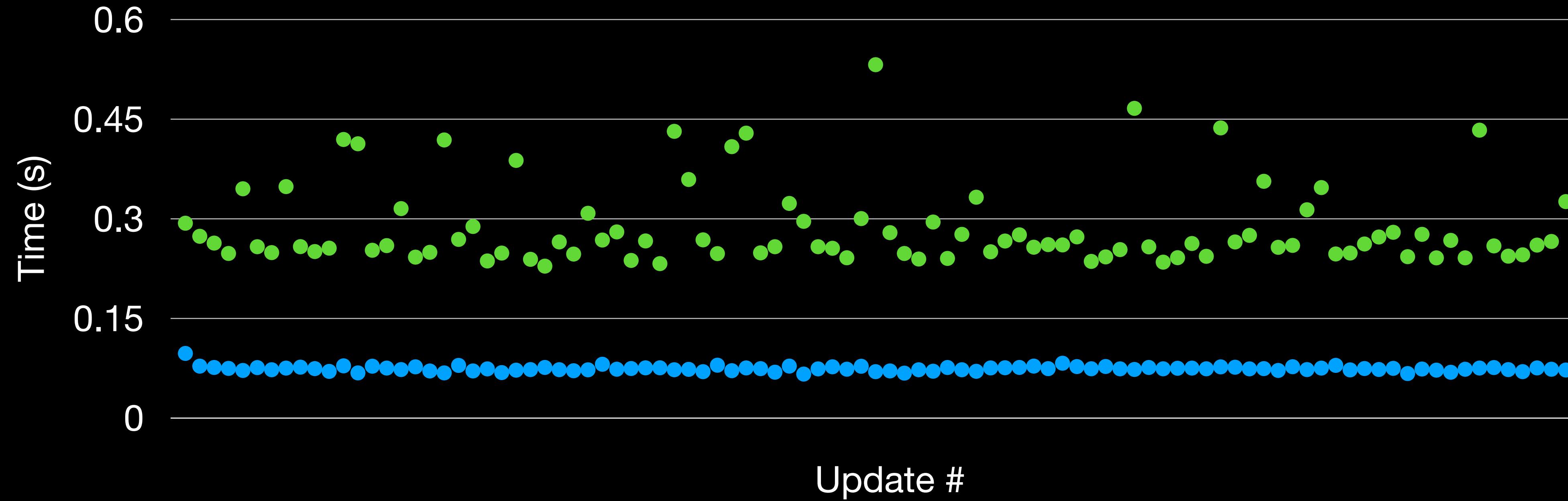


Figure 1: Generic benchmark results.



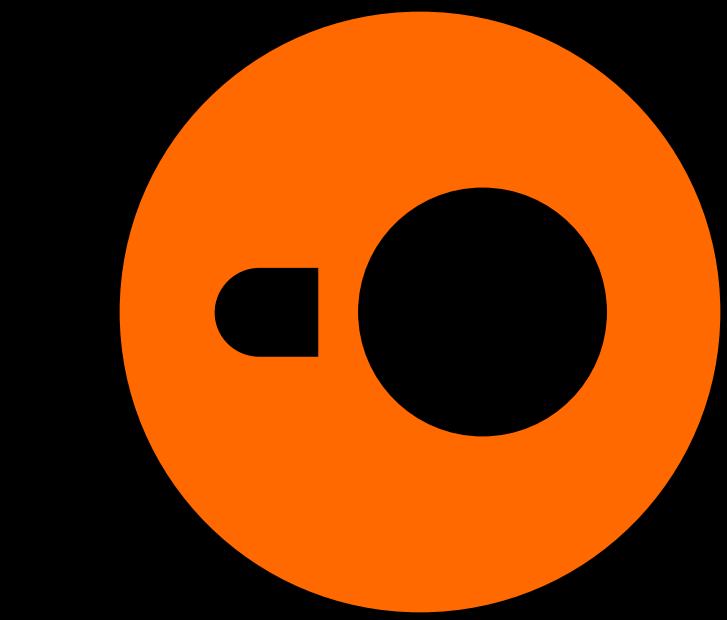
```
-----  
:: retrieving :: org.apache.spark#spark-submit-parent-1223a723-f206-4458-857d-d14611b6a7a9  
  confs: [default]  
    0 artifacts copied, 2 already retrieved (0kB/2ms)  
25/06/18 13:51:11 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-jav  
a  
Setting default log level to "WARN".  
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).  
[Stage 1:>                                         25/06/18 13:51:24 WARN GarbageCollectionMetrics: To enable  
n-built-in garbage collector(s) List(G1 Concurrent GC), users should configure it(them) to spark.eventLog.gcMetrics.yc  
GenerationGarbageCollectors or spark.eventLog.gcMetrics.oldGenerationGarbageCollectors  
[Stage 1:====>                               [Stage 1:=====>  
Stage 1:=====>                               [Stage 1:=====>  
0.2873690128326416
```

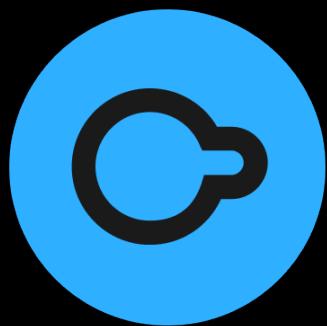
⚠ Program crashed: Bad pointer dereference at 0x0000db8a0000db95

Thread 0:

```
0          0x0000000180764c34 mach_msg2_trap + 8  in libsystem_kernel.dylib  
1 [ra]      0x000000018076d764 mach_msg_overwrite + 484 in libsystem_kernel.dylib  
2 [ra]      0x0000000180764fa8 mach_msg + 24  in libsystem_kernel.dylib  
3 [ra]      0x0000000180891e7c __CFRunLoopServiceMachPort + 160 in CoreFoundation  
4 [ra]      0x0000000180890798 __CFRunLoopRun + 1208 in CoreFoundation  
5 [ra]      0x000000018088fc58 CFRunLoopRunSpecific + 572 in CoreFoundation  
6 [ra]      0x0000000104ac326c CreateExecutionEnvironment + 400 in libjli.dylib  
7 [ra]      0x0000000104abf344 JLI_Launch + 1128 in libjli.dylib  
8 [ra]      0x0000000104a83bac main + 392 in java  
9 [ra] [system] 0x0000000180406b98 start + 6076 in dyld
```







Multi-Table

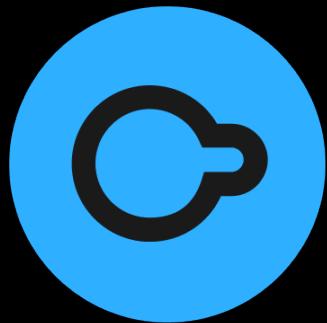
```
ATTACH 'ducklake:metadata.ducklake' AS dl;
```

```
CREATE SCHEMA dl.s1;
```

```
CREATE TABLE dl.s1.t1 ...;
```

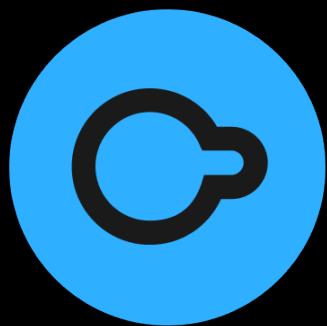
```
CREATE TABLE dl.s1.t2 ...;
```

```
...
```



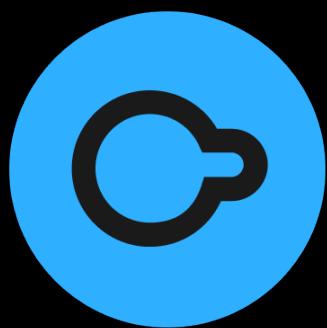
Holistic Transactions

```
BEGIN  
CREATE TABLE ...  
ALTER TABLE ...  
INSERT ... INTO t1  
UPDATE t2 ...  
DELETE FROM t3  
COMMIT
```



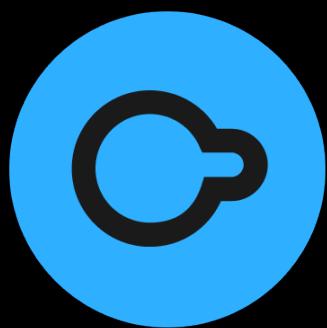
Time Travel

```
FROM tbl AT (VERSION => 3);  
FROM tbl AT (TIMESTAMP => ...);
```



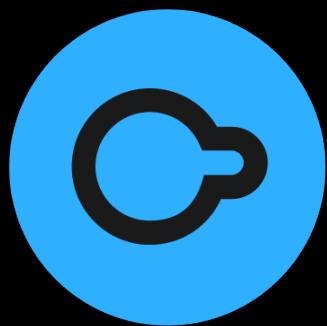
Change Feed

```
ATTACH 'ducklake:metadata.ducklake' AS dl;  
FROM dl.table_changes('tbl', 3, 4);
```



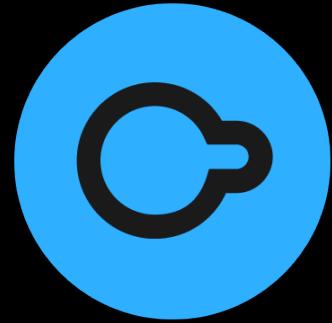
Inlining

```
ATTACH 'ducklake:metadata.ducklake' AS dl  
  (DATA_INLINING_ROW_LIMIT 100);
```



Encryption

```
ATTACH 'ducklake:metadata.ducklake' AS dl  
  (DATA_PATH '/some/untrusted/location' ENCRYPTED);
```



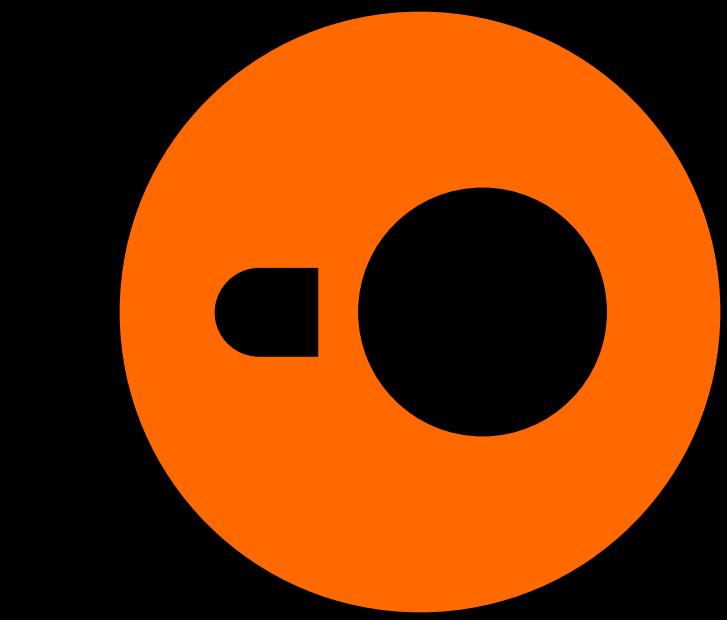
Interoperability

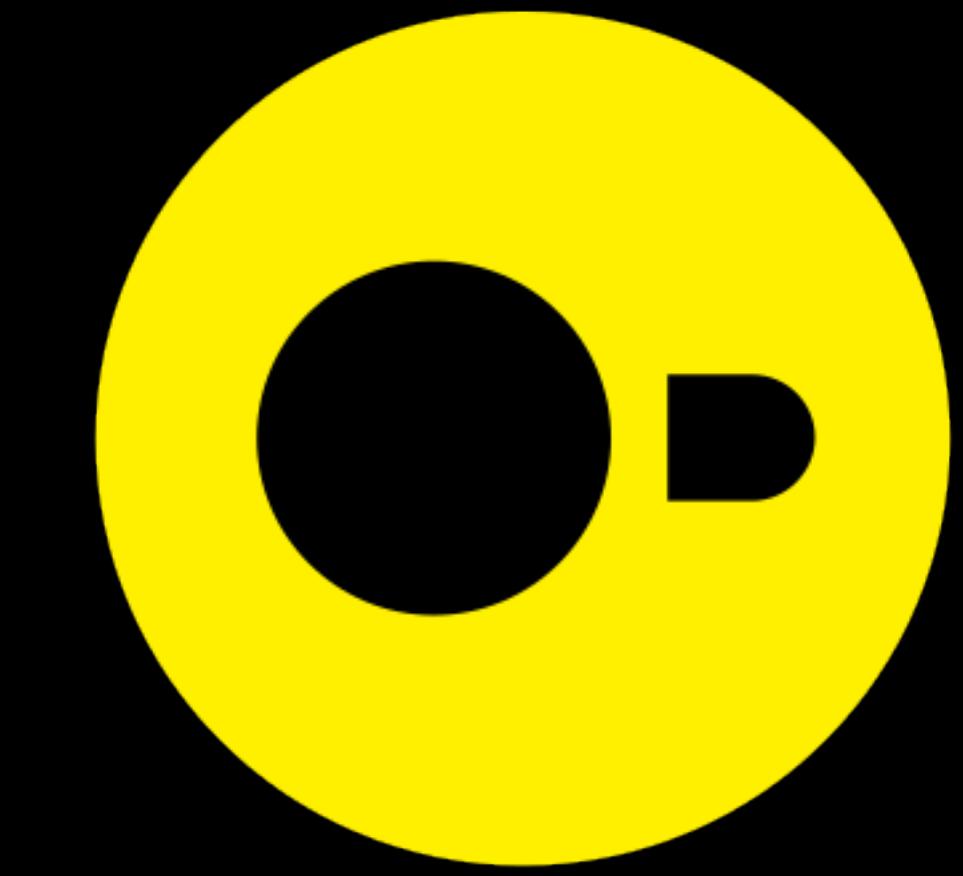




Free

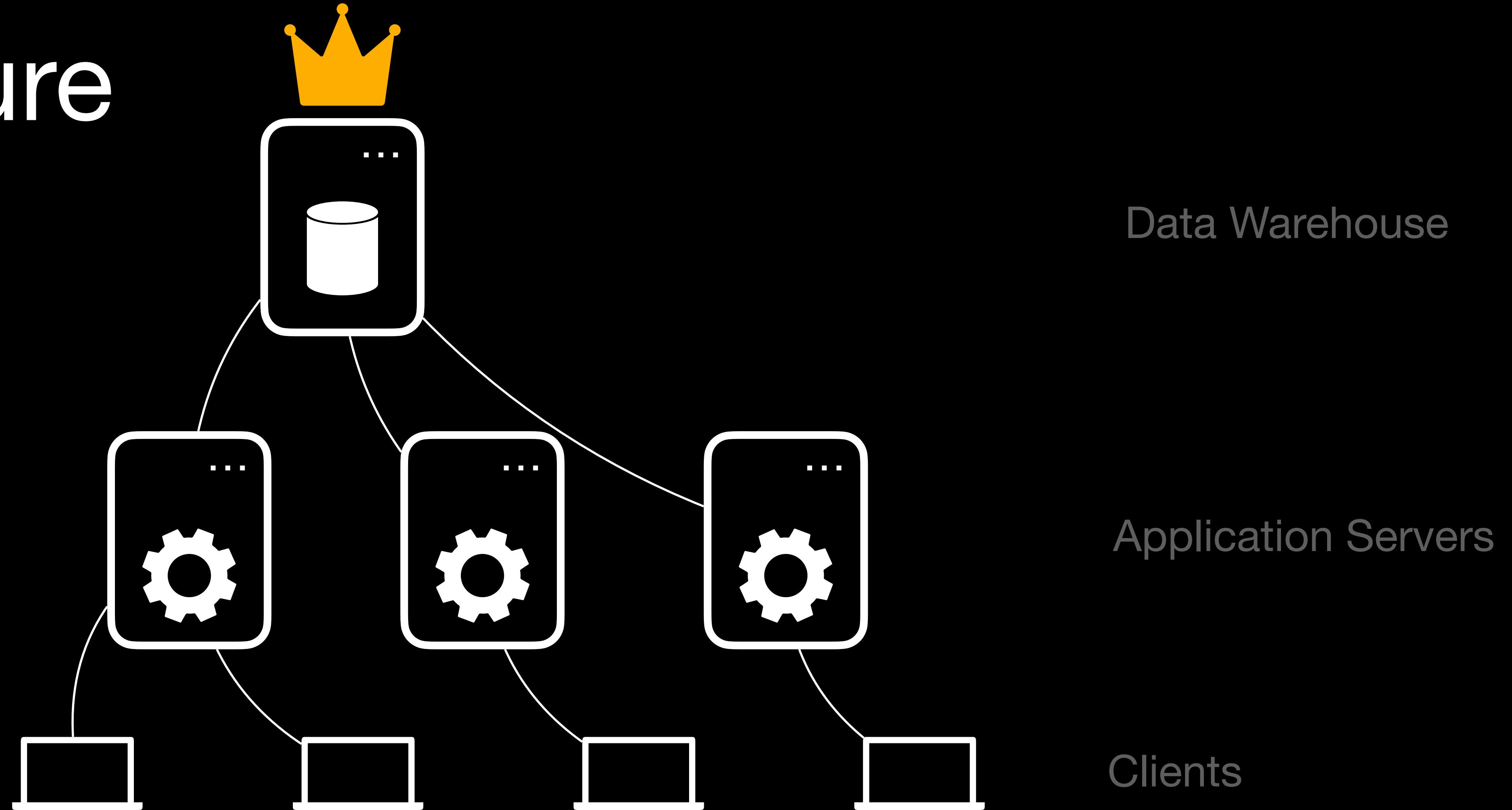




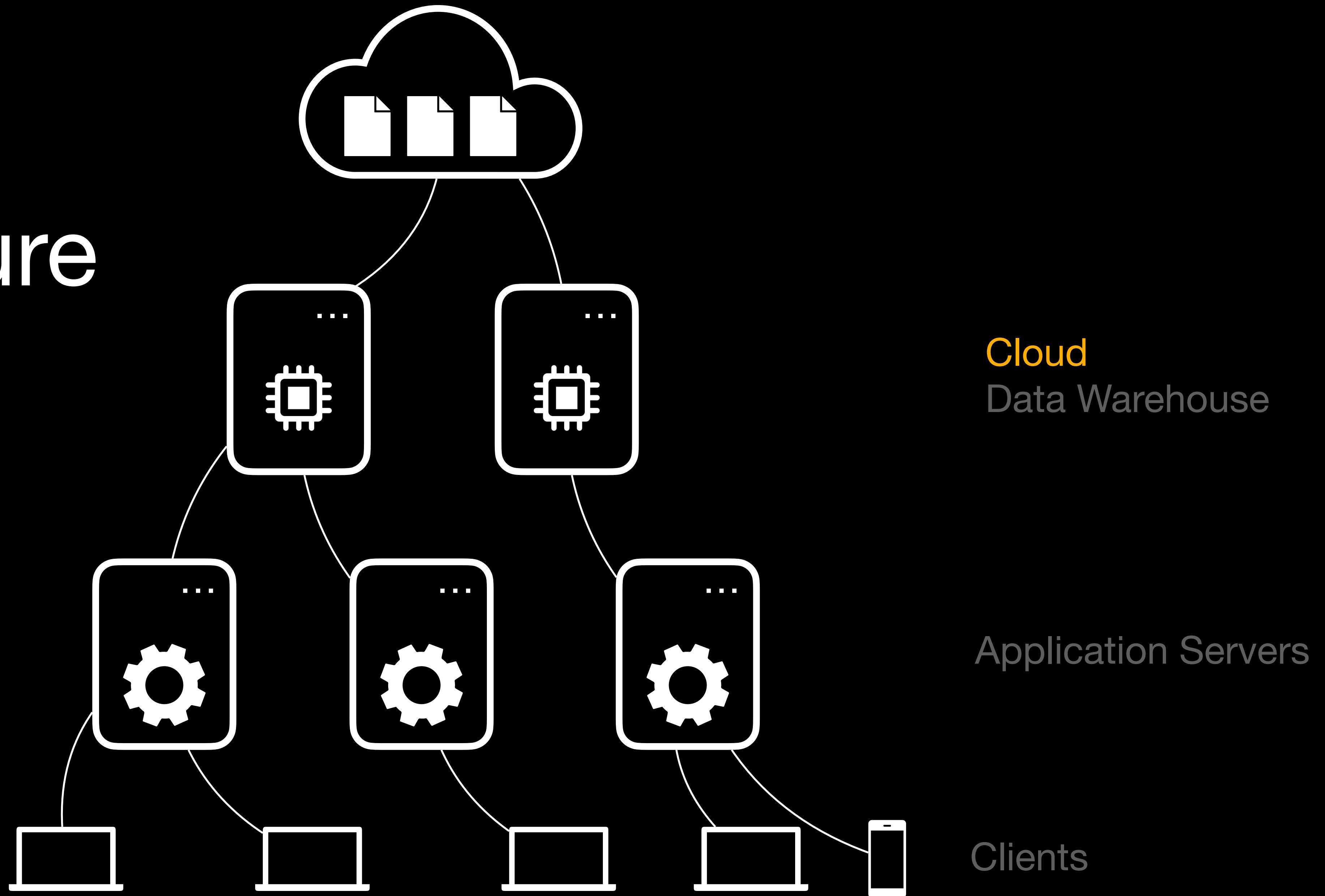


DuckDB

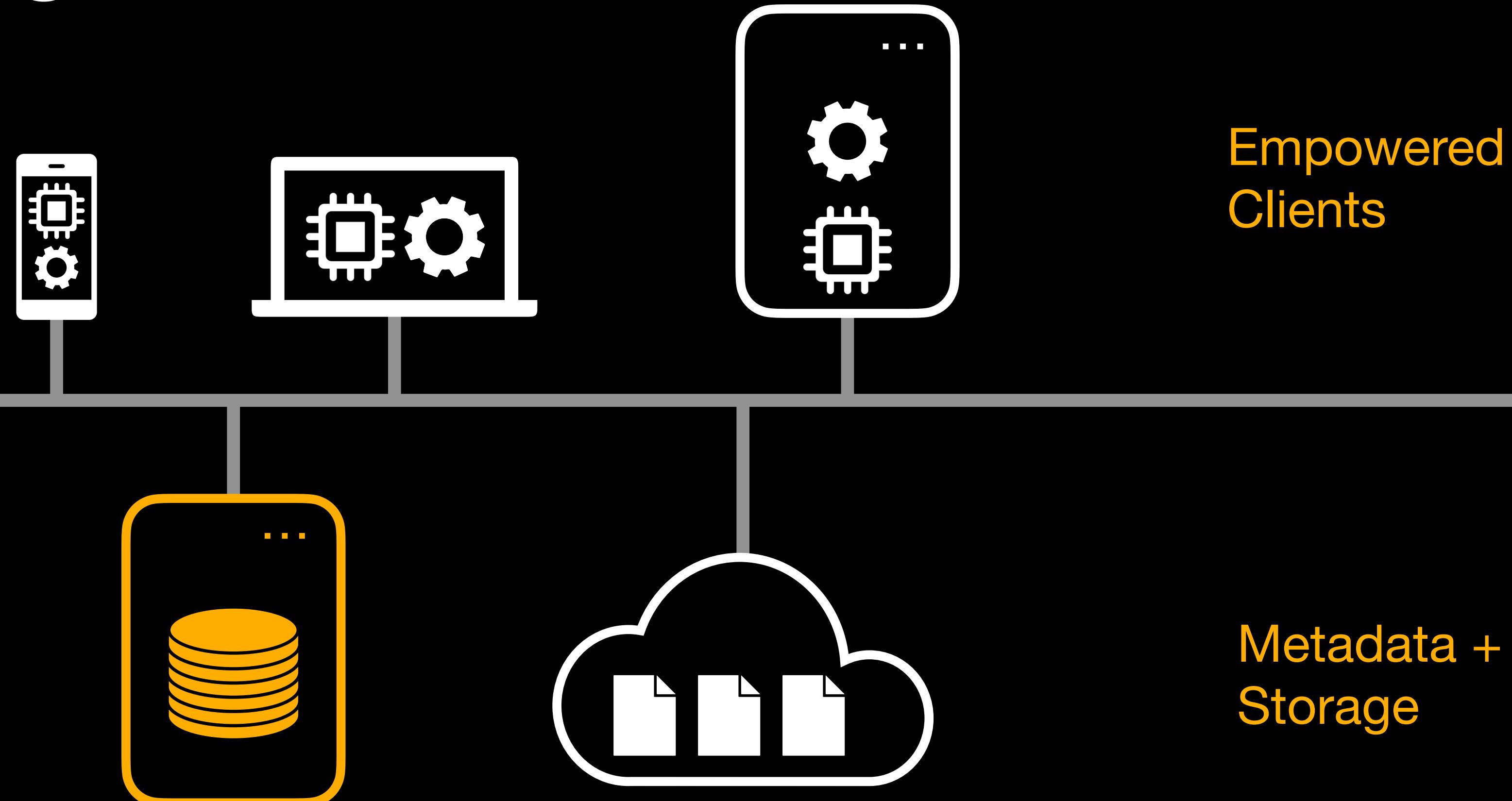
1985 Data Architecture



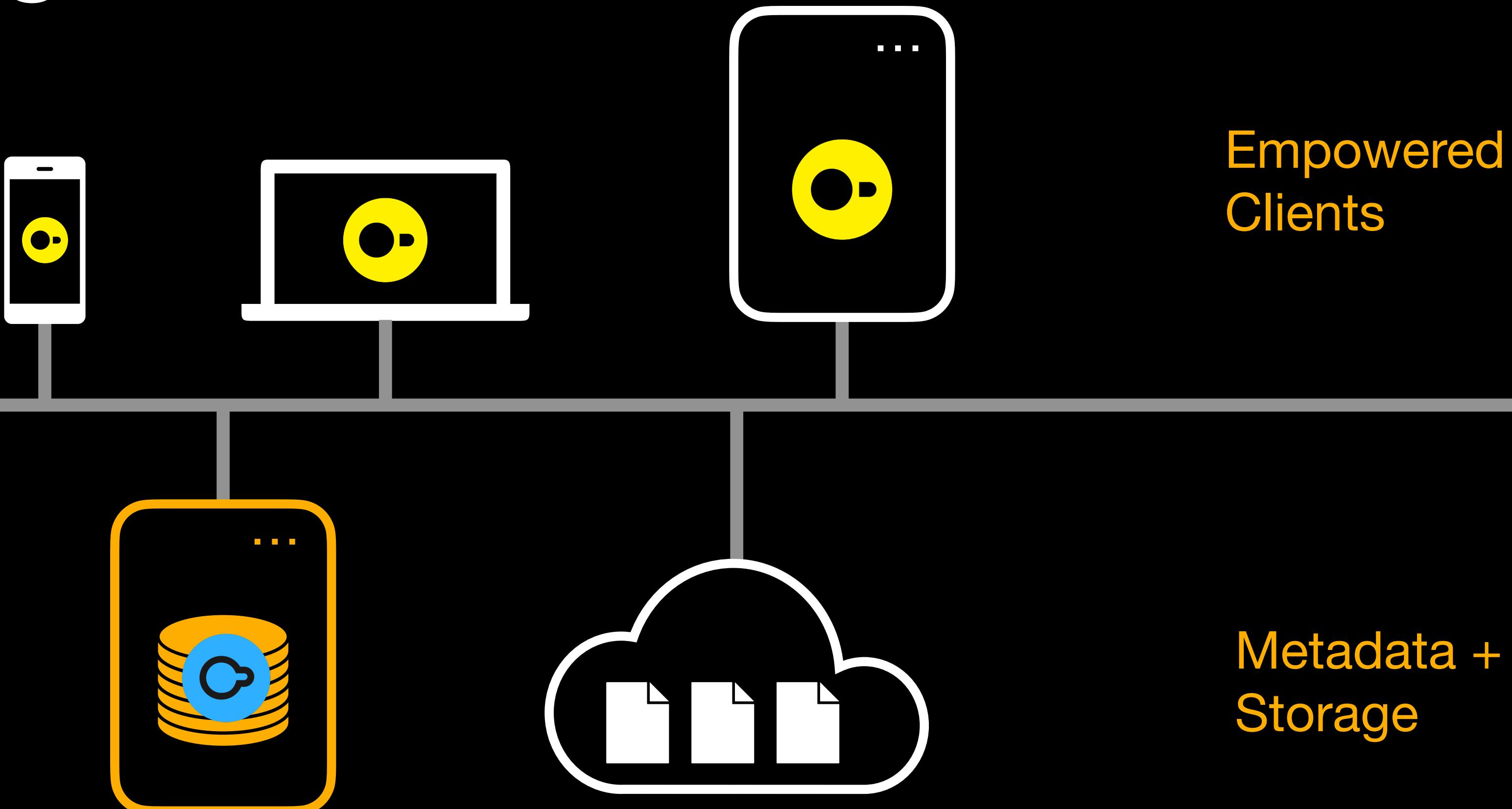
2015 Data Architecture

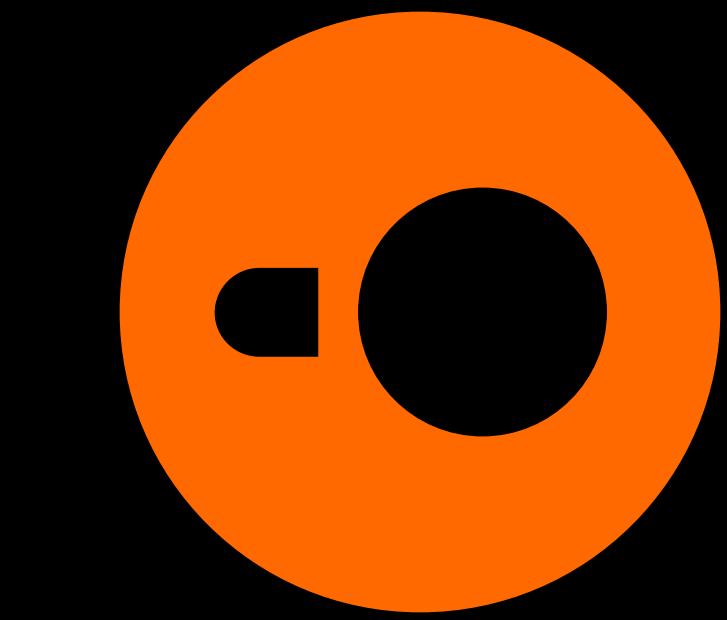


2025 Data Architecture

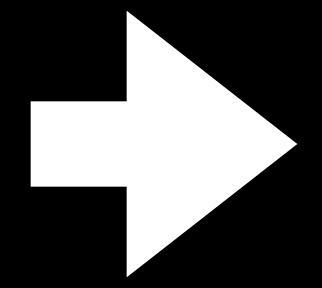


2025 Data Architecture





Data Fear



Data Confidence



DuckLake

ducklake.select
hannes.muehleisen.org