# DuckDB
# A Modern Modular & Extensible Database System

Mark Raasveldt

DuckDB Labs
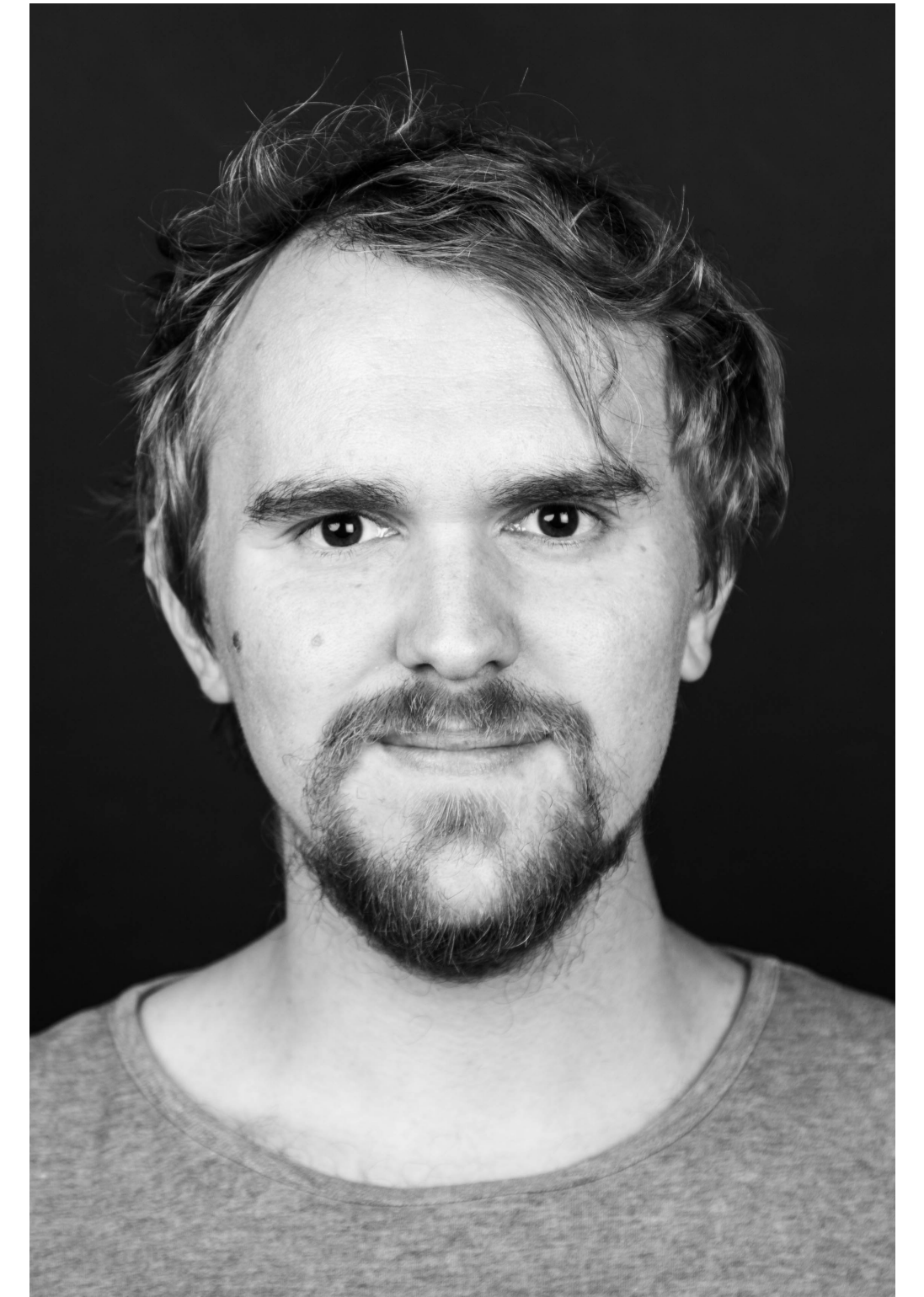
Mark Raasveldt

CTO of DuckDB Labs
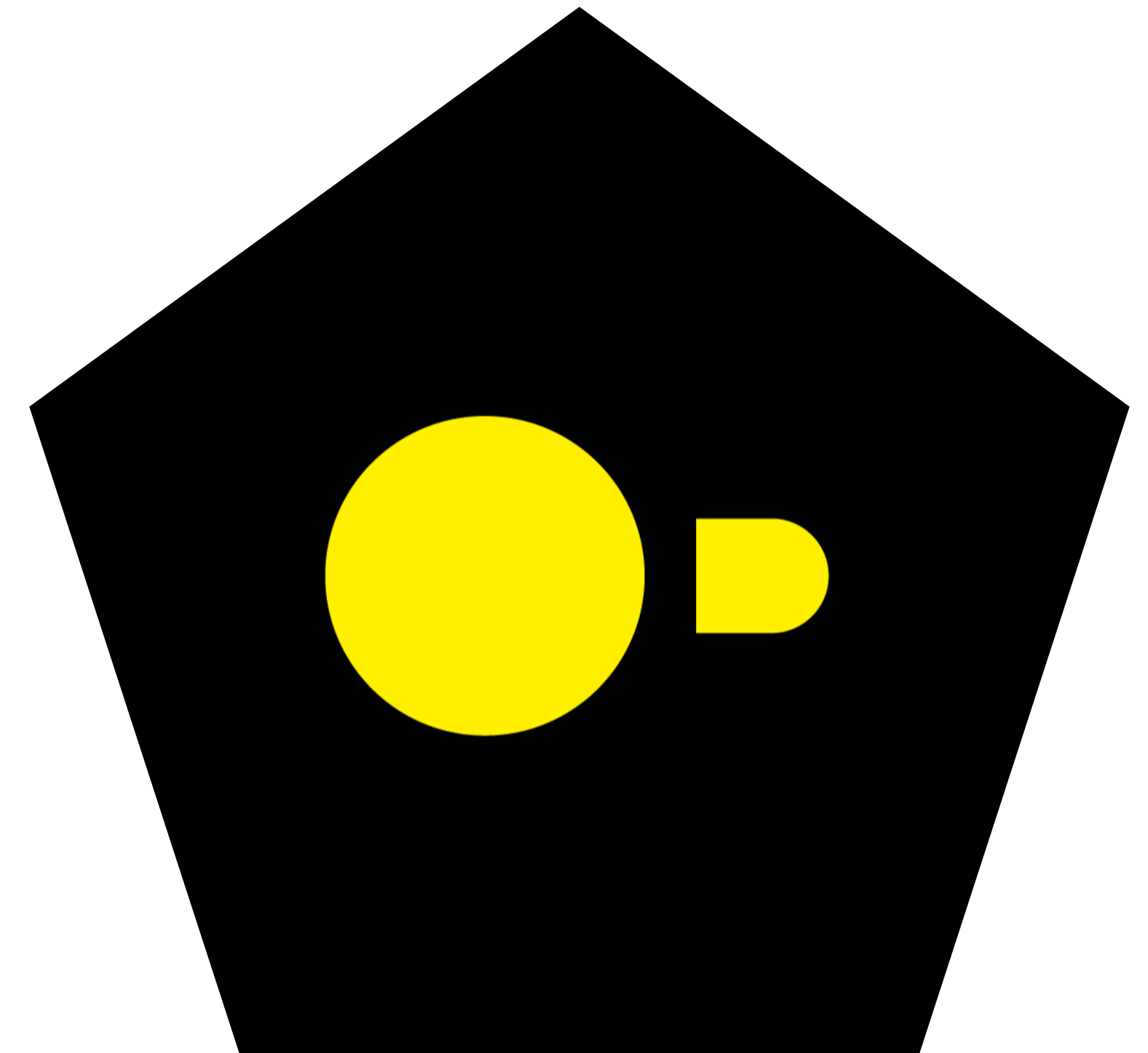
Postdoc at CWI

  Database Architectures Group

PhD at CWI

- DuckDB

- In-Process OLAP DBMS

  - "The SQLite for Analytics"

- Free and Open Source (MIT)
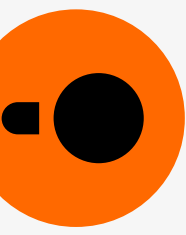
- [duckdb.org](duckdb.org)

# What is DuckDB?

- SQLite inspired us in many ways:

  - Easy installation

  - Ease of use

  - Robustness

- DuckDB aims to be the "SQLite for Analytics"

Analytics has **unique challenges**



Transactional workloads = **simple queries**

SQLite can be feature complete with a very small footprint

Analytical workloads = **complex queries**

DuckDB needs **many more** operations, functions, optimizers,….

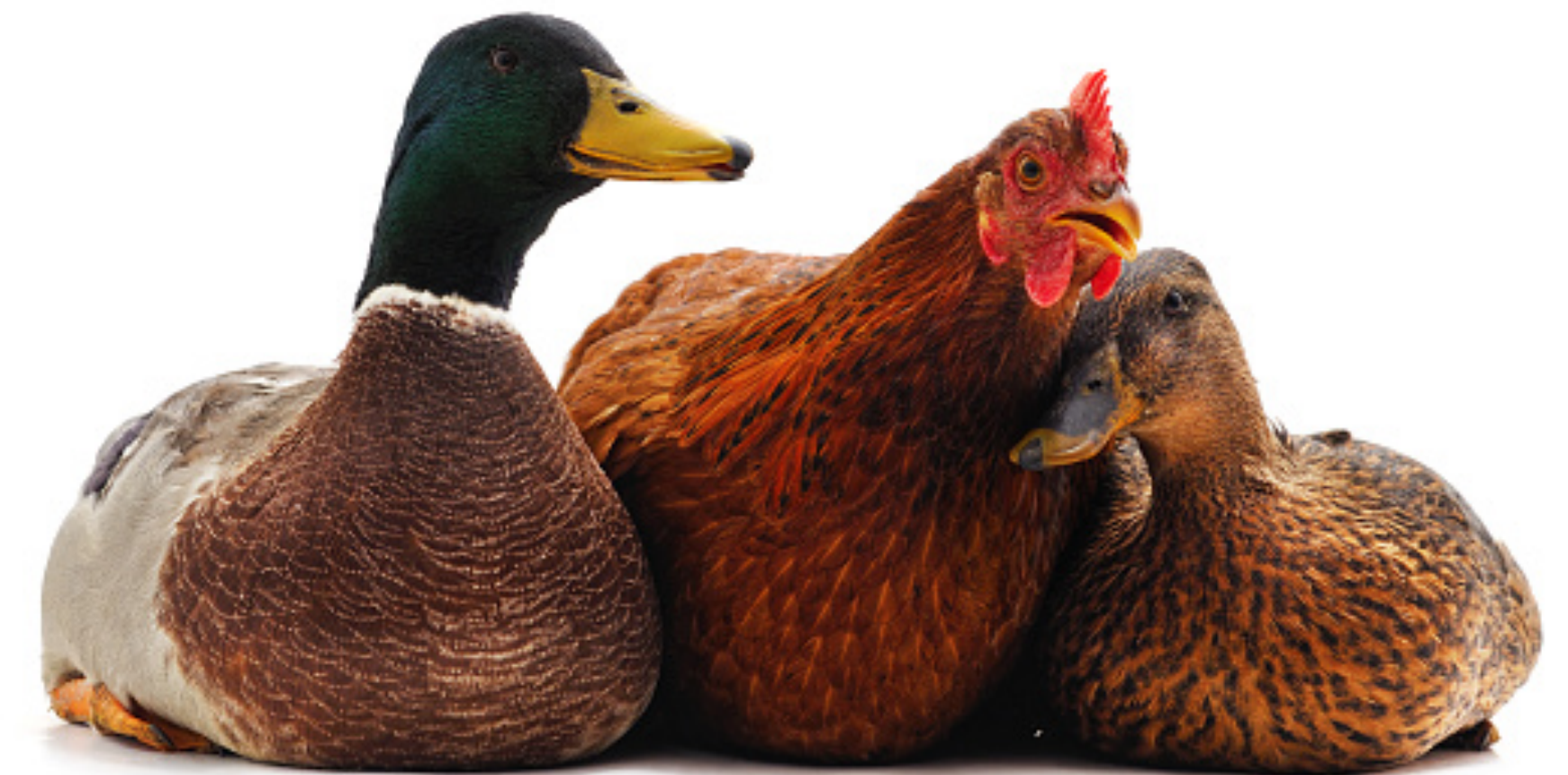Analytics requires a giant diversity of operations

**Collaboration is required!**

A single entity cannot hope to implement:

Data wrangling tools

Classification algorithms, ML toolkits

Data cleaning tools

etc...

**Research** other important collaboration area

DuckDB **originates from** research world

At CWI and elsewhere, people do their research in DuckDB

Important to allow **extending/modifying system**

e.g. add new join operator, new optimizer, ...

**System builders** want to use **components** of DuckDB

For example:

Use only the front-end

Use only the back-end

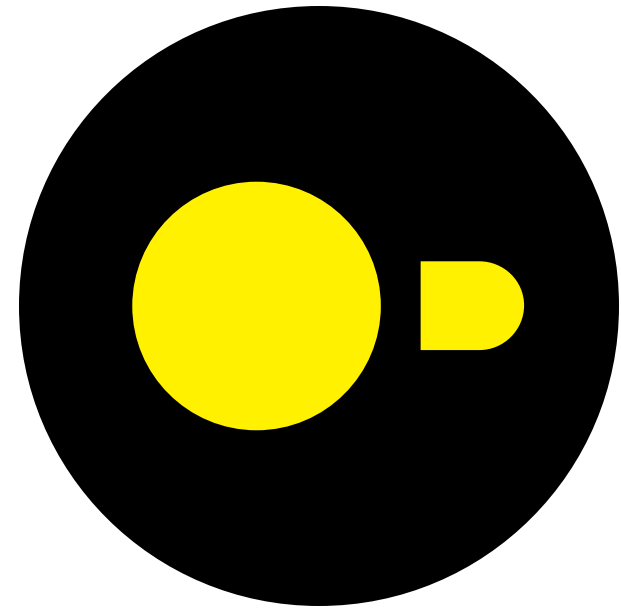How do we **enable this collaboration**?

Three aspects:

Flexible **data import and export**

**Extensibility** of the system

Hooks in different locations in the system

# Flexible Data Import & Export

**Data import and export**
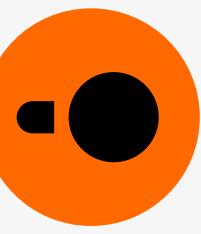
Crucial for communication **between libraries**

**Many use cases**

Load data exported from other systems

Pre-process in DuckDB→ use plotting/statistics/ML libraries

Mix usage DuckDB + other data wrangling libraries

Export data to persistent storage (e.g. Parquet files on S3)

What makes DuckDB different?

All database systems can import/export data

**… but very slowly!**

Don't Hold My Data Hostage -
A Case For Client Protocol Redesign

Ancient protocols, designed to transfer kilobytes of data

Unsuitable for modern workloads!

Import/export through text-based formats (CSV files)

DuckDB is designed for **bulk data export/import**

  **In-process = zero-copy data sharing**

Versatile **input & output APIs**



**Same API** used for both internal tables & external sources

  Allows for **streaming data, parallel scans, projection & filter pushdown, index usage, …**

# Data Import & Export

DuckDB can **efficiently** consume and output **many formats**

# Extensibility of the System

**Extensibility of the system**

DuckDB supports **extensions** that add new functionality

Allows users to **integrate new functionality in the system**

We create extensions ourselves liberally

**Eat your own duck food**

Allows us to add functionality…

…without **bloating the core system**

…that is **only included in certain distributions** (e.g. Python client)

…without introducing **external dependencies** to the core

## ICU Extension

Adds support for **collations (language-based sorting, comparisons)**

Adds support for **time-zone awareness**

Extension is **several megabytes**

**Includes ICU localization data**

Roughly same size as DuckDB core!

**As an extension we keep this optional**

## HTTP-FS Extension

Adds support for **reading/writing data over HTTP(S)**

Adds support for **reading/writing data with S3**

Adds **dependency on OpenSSL**

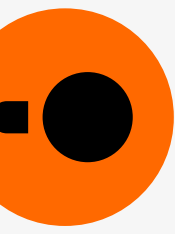As an extension, we **keep the core dependency-free**

Amazon S3

Extensions are **powerful**

**Goal:** Allow **every component of the system** to be extended

Currently extensions can:
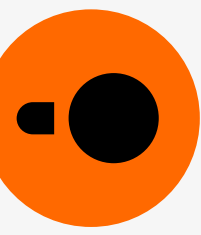
Add functions (scalar, aggregate, window)

Add new types

Add data sources and sinks

Add collations, time zones

Add parser functionality, optimizers

**Goal:** allow users to **create and maintain their own extensions**

   No need to **talk to us**

   No need for us **to maintain their code**

**Extension Repository**
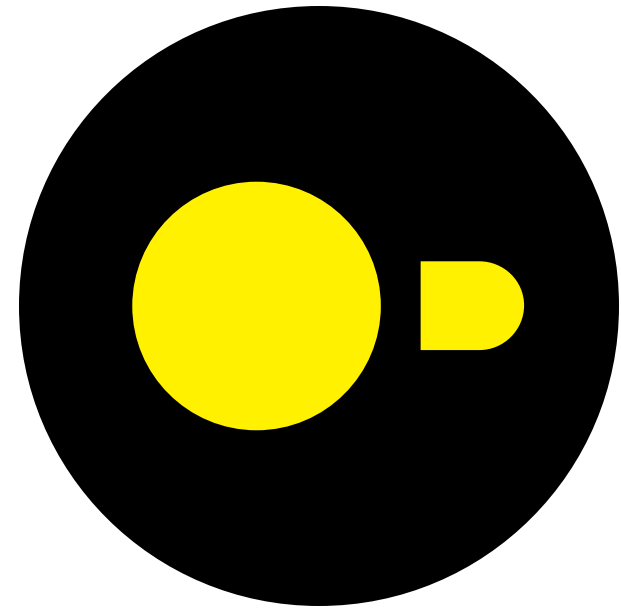
Extensions can be installed using SQL

```
INSTALL httpfs;
LOAD httpfs;
```
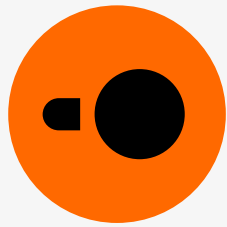
By default from our repository

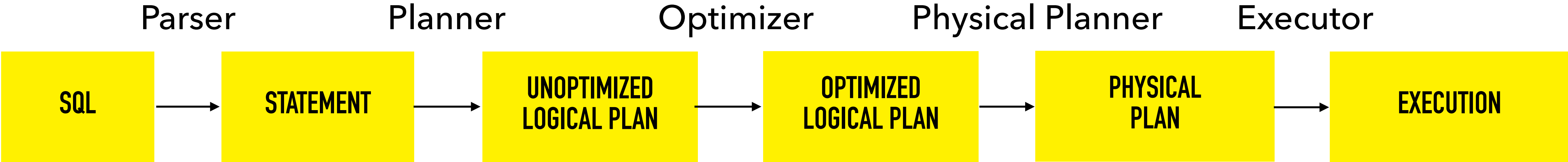   Custom repository can be used

# Different Hooks Into The System

DuckDB uses a typical pipeline for query processing

| Parser | | Planner | | Optimizer | | Physical Planner | | Executor |
|---|---|---|---|---|---|---|---|---|
| SQL | → | STATEMENT | → | UNOPTIMIZED LOGICAL PLAN | → | OPTIMIZED LOGICAL PLAN | → | PHYSICAL PLAN | → | EXECUTION |

**Standard workflow: use all components**
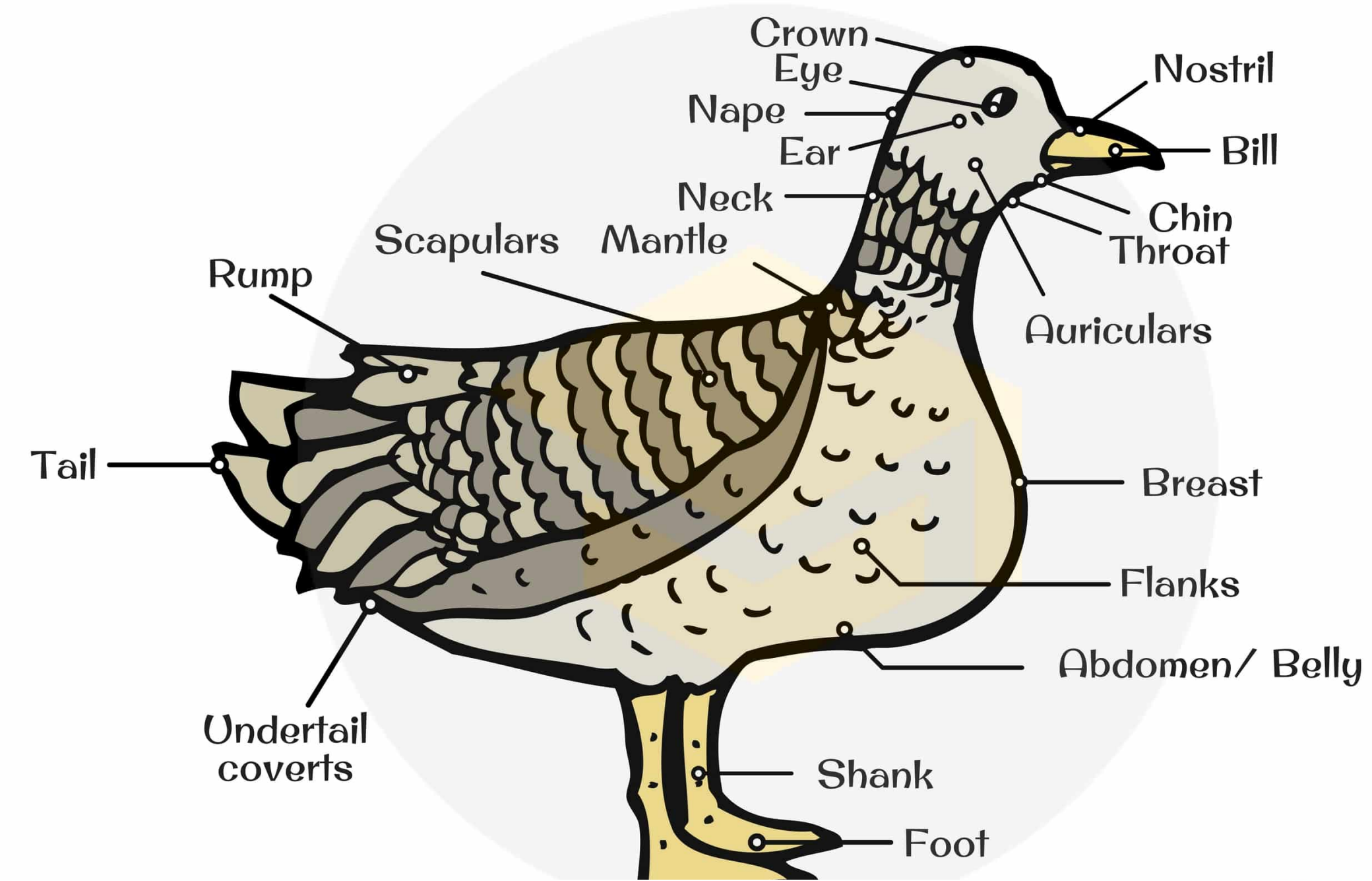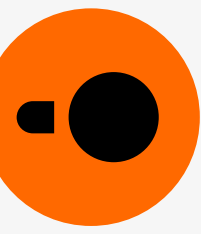
**Standard pipeline: SQL input → Data output**

**System builders/researchers** often want to use **parts of DuckDB**

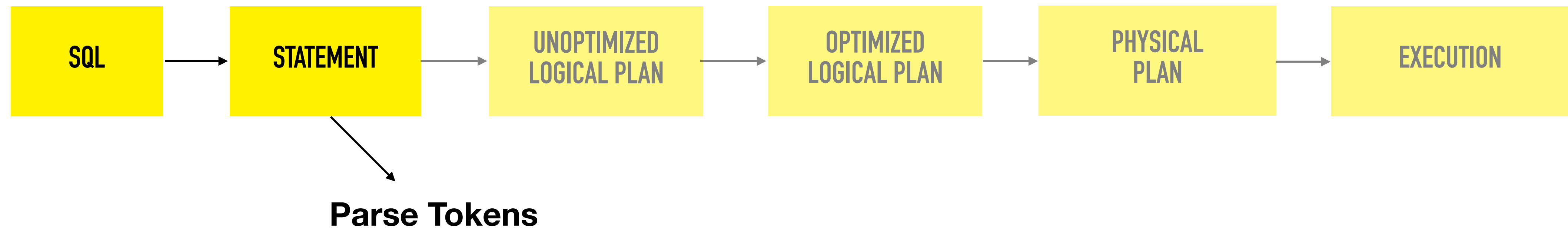**Different hooks in/out of the system**

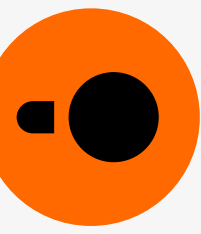# Different Hooks Into The System

Velox: to use the parser
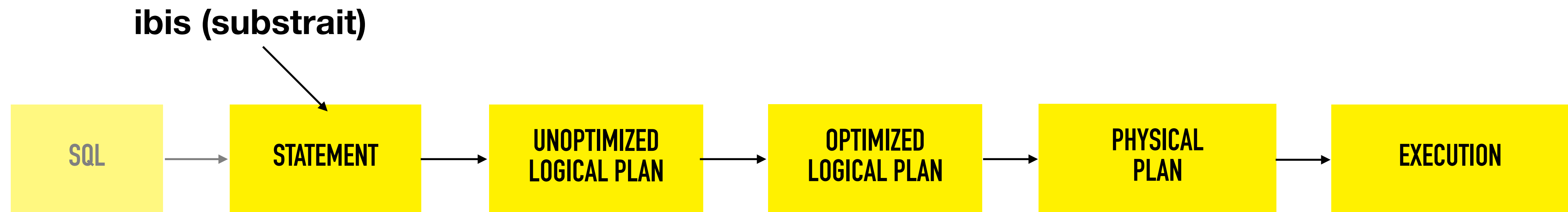
**Parser pipeline: SQL Input → Parse Tokens Output**



| SQL | → | STATEMENT | → | UNOPTIMIZED LOGICAL PLAN | → | OPTIMIZED LOGICAL PLAN | → | PHYSICAL PLAN | → | EXECUTION |

**Parse Tokens**

Ibis: provide a **new front-end to DuckDB**

**Ibis front-end pipeline: Substrait plan input → Data Output**

**ibis (substrait)**

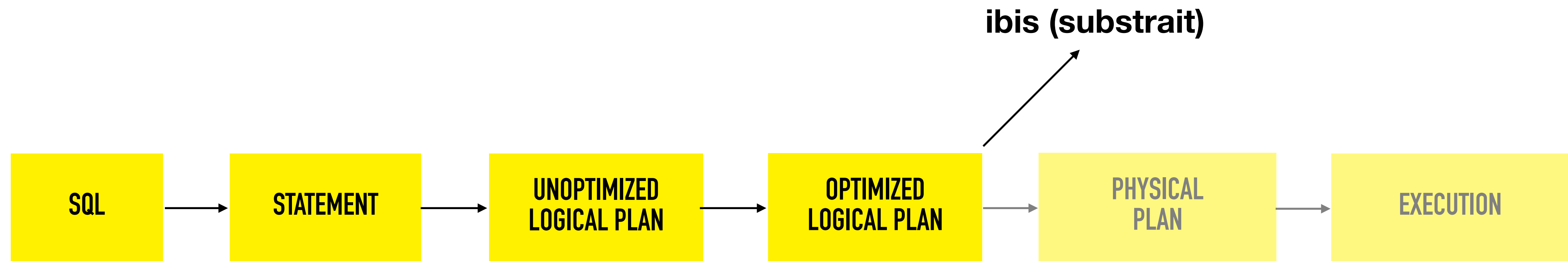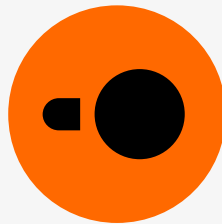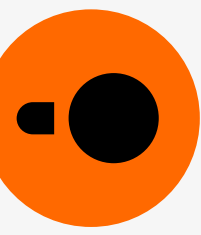| SQL | → | STATEMENT | → | UNOPTIMIZED LOGICAL PLAN | → | OPTIMIZED LOGICAL PLAN | → | PHYSICAL PLAN | → | EXECUTION |

# Different Hooks Into The System

Ibis: **consume SQL** and use a different back-end

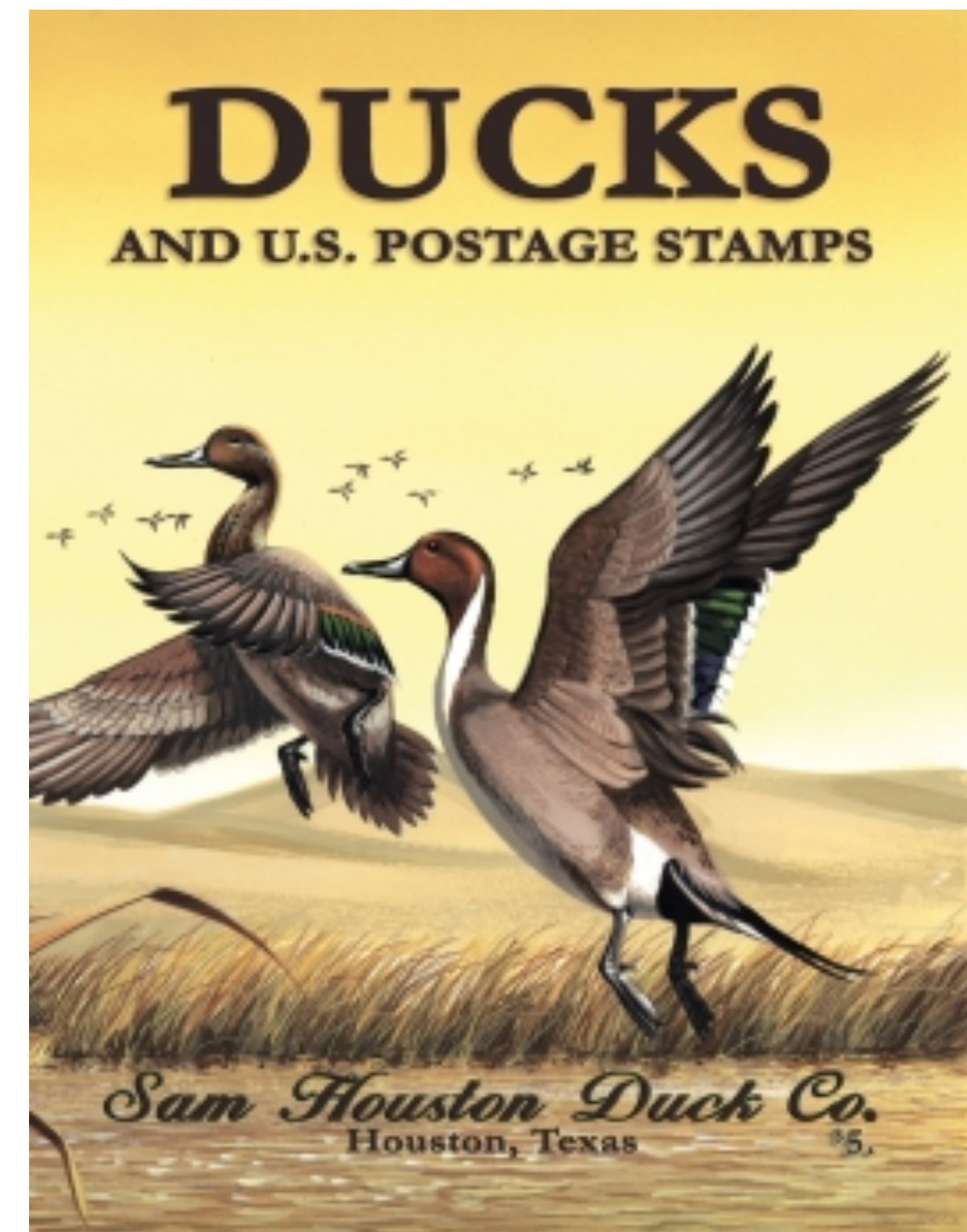**Ibis back-end pipeline: SQL input → Substrait plan output**

ibis (substrait)

| SQL | → | STATEMENT | → | UNOPTIMIZED LOGICAL PLAN | → | OPTIMIZED LOGICAL PLAN | → | PHYSICAL PLAN | → | EXECUTION |

**Goal:** more **extensibility** options

New index types

Alternate versions of existing operators (e.g. new joins)
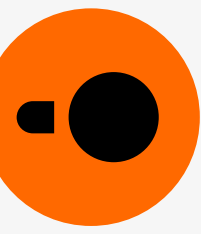
Define new casts between types

**Goal:** support fully custom storage back-end

Currently support **data-input** and **copy output**

**Missing support** for routing **insert, update, delete** to custom storage
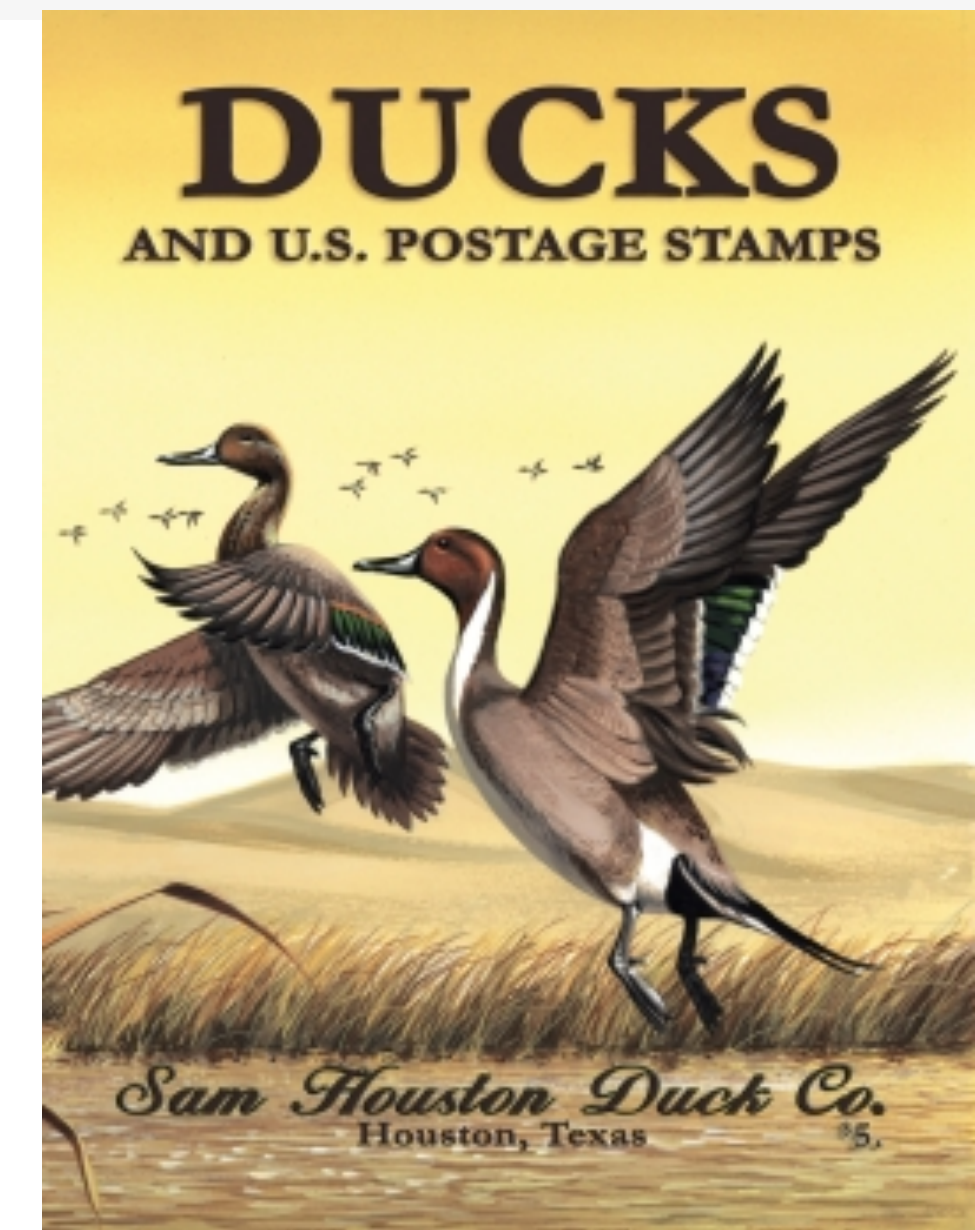
**Goal:** extended custom **catalog** support

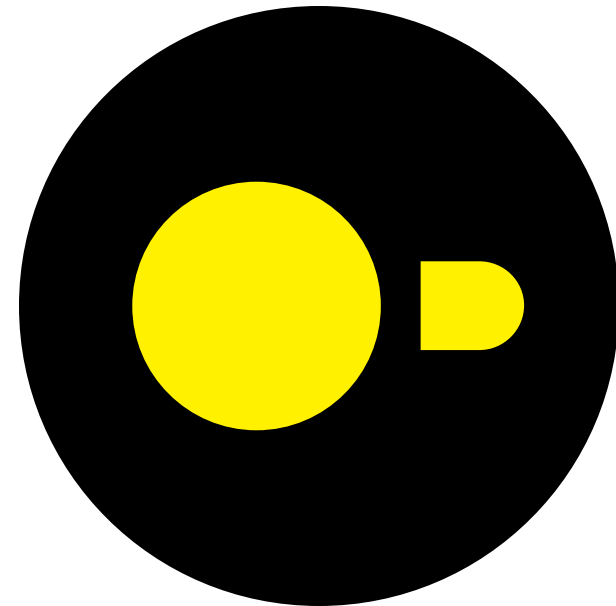Currently support **replacement scans**

This allows for **views** to be stored in a different catalog

Needs to be extended to fully support custom catalogs

Including for all **different catalog types**

Thanks for having me!

Any questions?