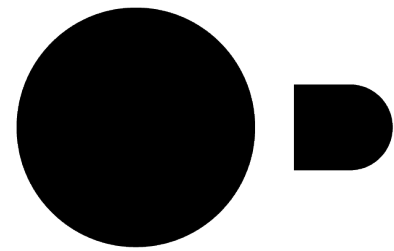


Mark Raasveldt & Pedro Holanda

DuckDB **an Embeddable Analytical RDBMS**

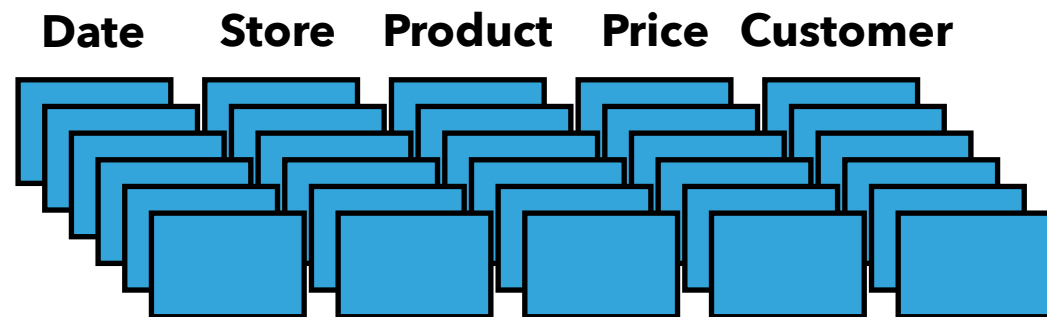
- ▶ **Internals at a Glance**
- ▶ Query processing pipeline
- ▶ Query execution
- ▶ Hands-On



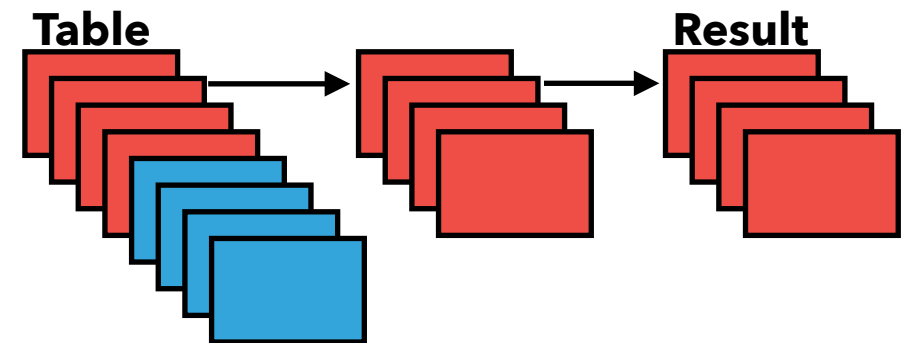
DuckDB

- ▶ Embedded analytical database
- ▶ Simple installation
 - ▶ `pip install duckdb`
- ▶ Fast and easy to use
- ▶ <https://www.duckdb.org>

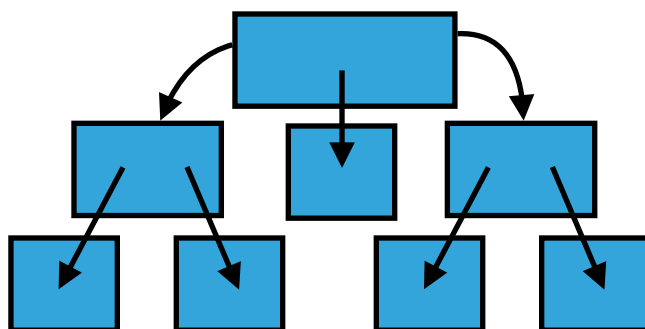
Column-Store



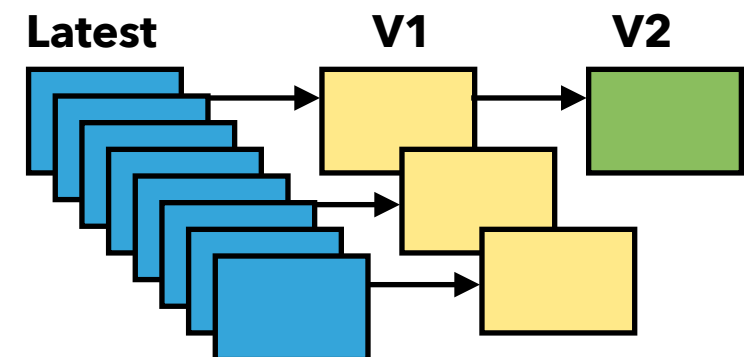
Vectorized Processing



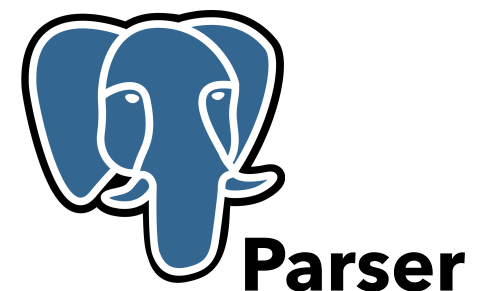
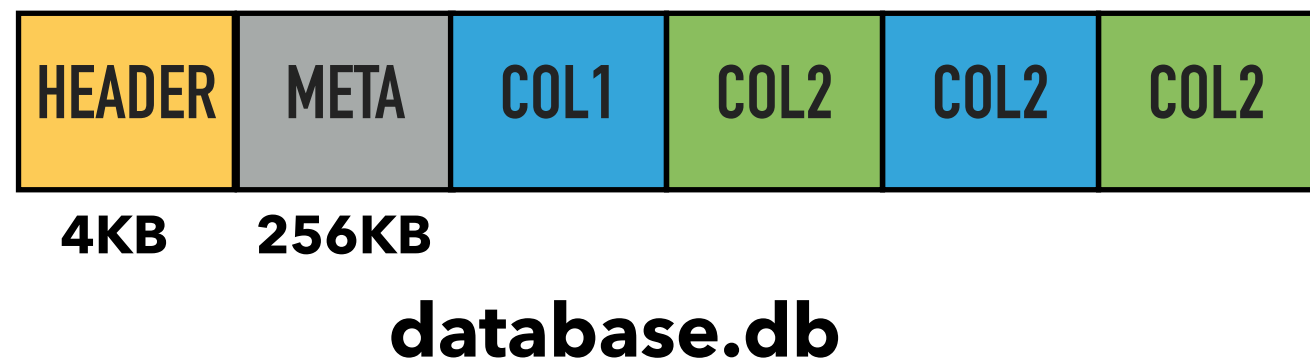
ART Index



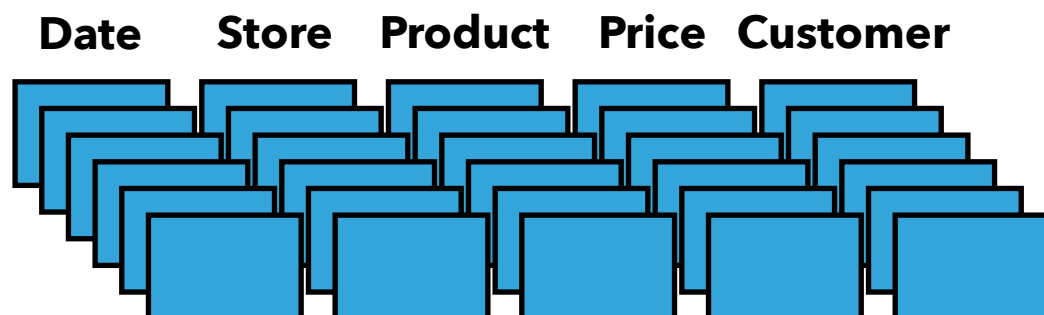
Multi-Version Concurrency Control



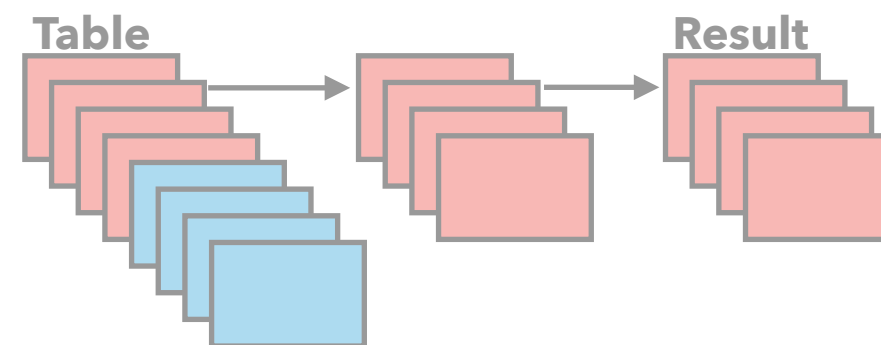
Single-File Storage



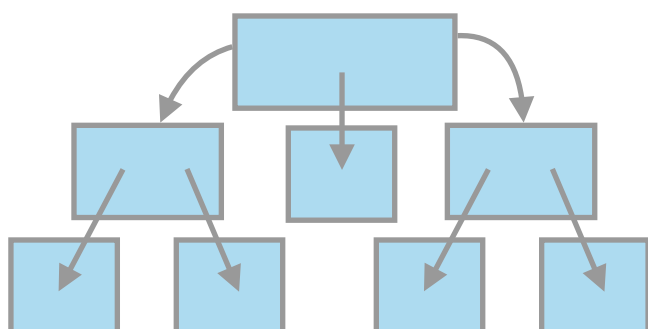
Column-Store



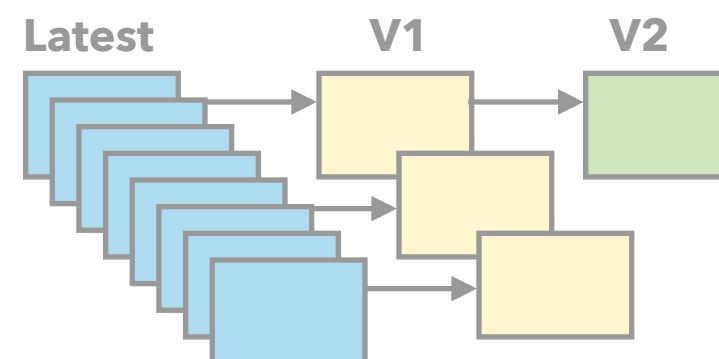
Vectorized Processing



ART Index



Multi-Version Concurrency Control



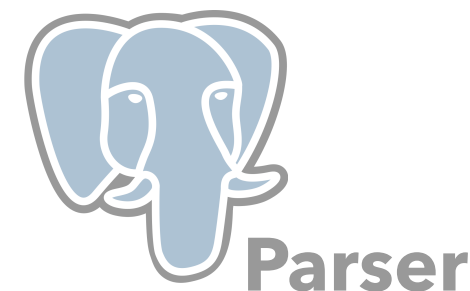
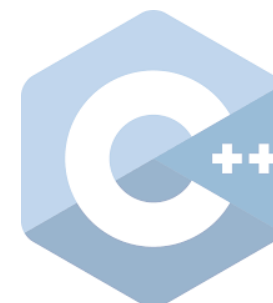
Single-File Storage



4KB

256KB

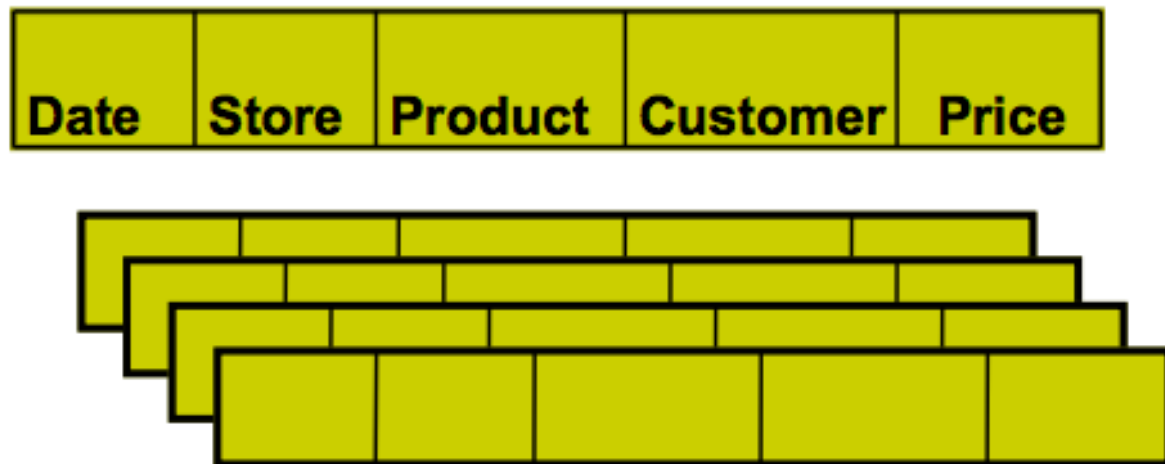
database.db



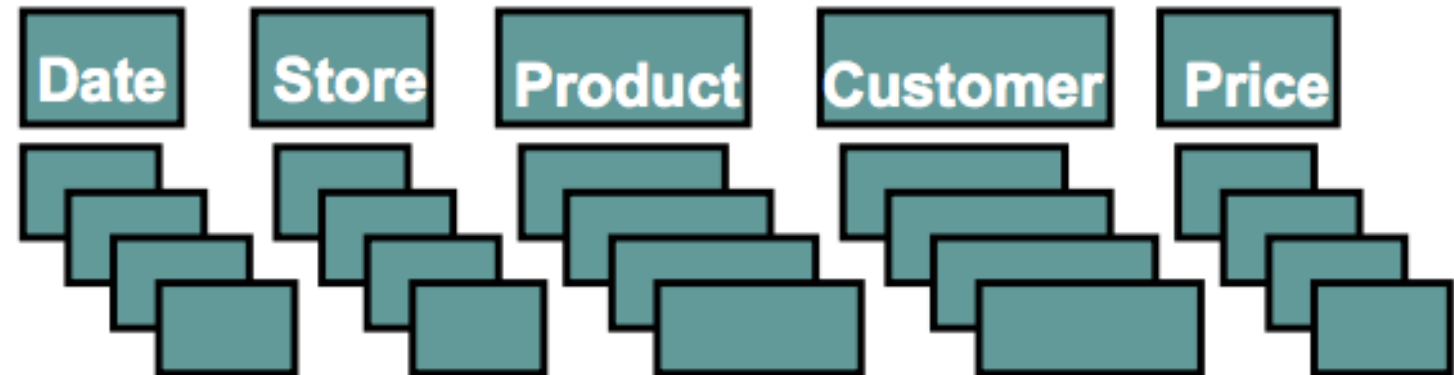
► Storage Model

- Traditional RDBMS use a row-storage model
- DuckDB uses a columnar storage model

row-store



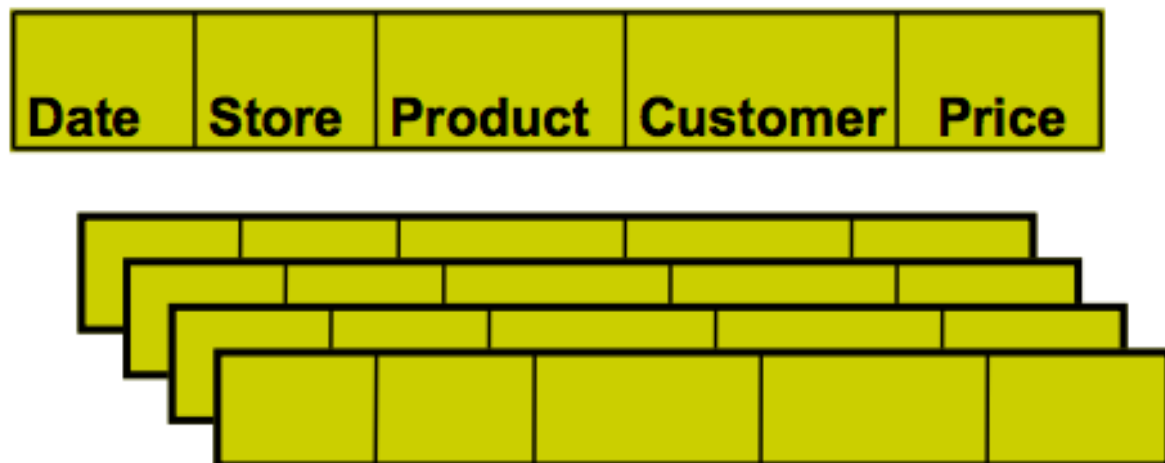
column-store



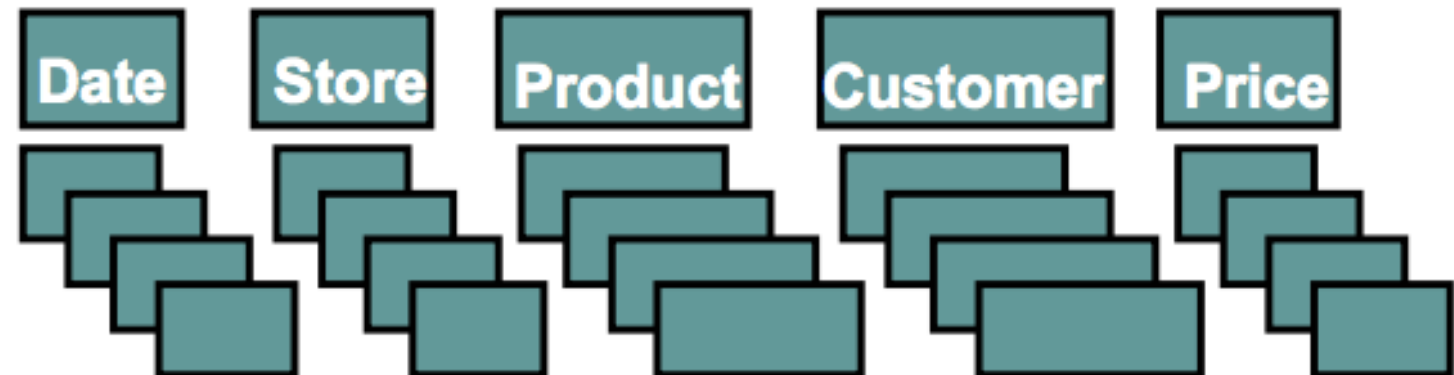
▶ Row-Storage:

- ▶ Individual rows can be fetched cheaply
- ▶ However, **all columns must always be fetched!**
- ▶ What if we only use a few columns?
- ▶ **e.g.:** What if we are only interested in the price of a product, not the stores in which it is sold?

row-store



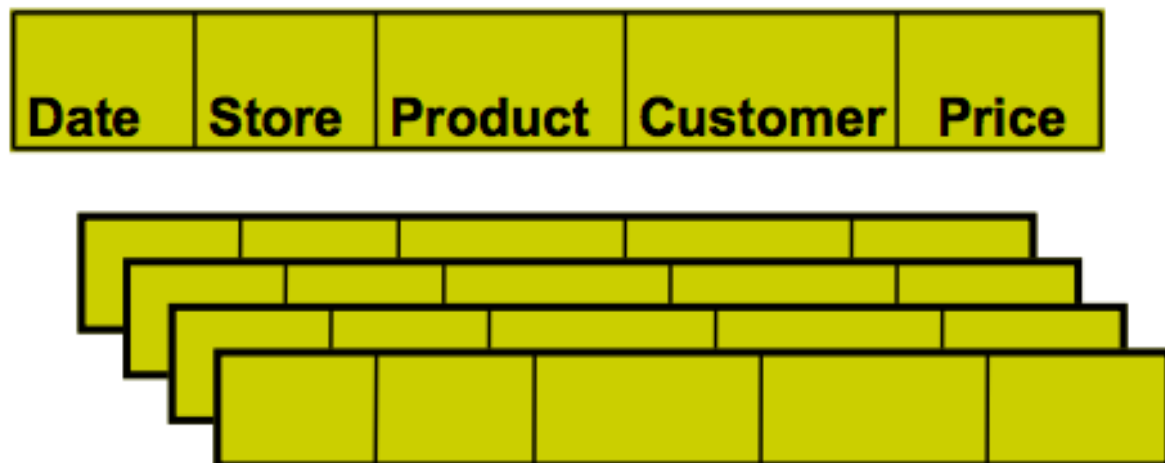
column-store



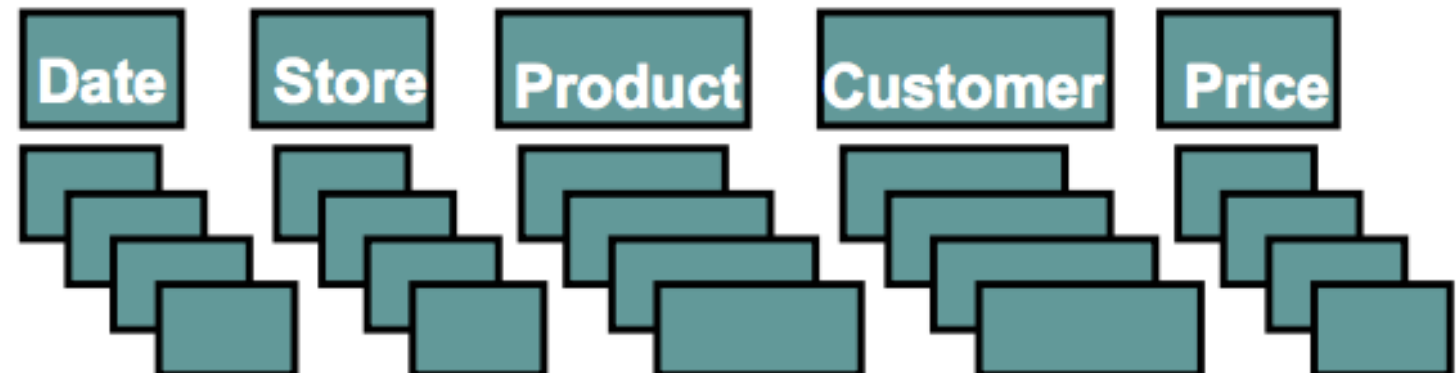
▶ Column-Storage:

- ▶ We can fetch individual columns
- ▶ Immense savings on disk IO/memory bw when only using few columns
- ▶ Queries that would take hours in a row-store can take seconds in a column-store

row-store



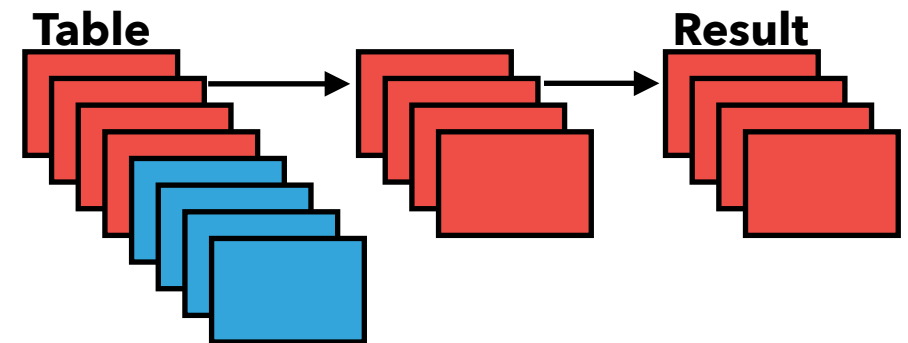
column-store



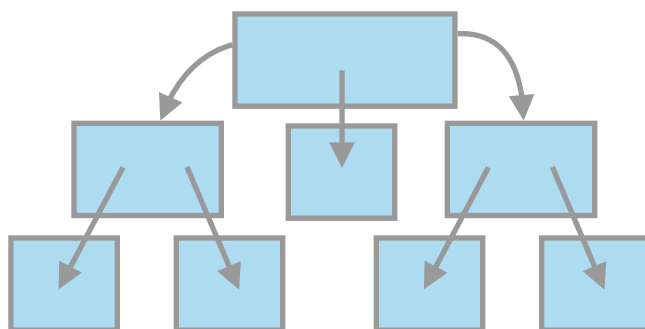
Column-Store



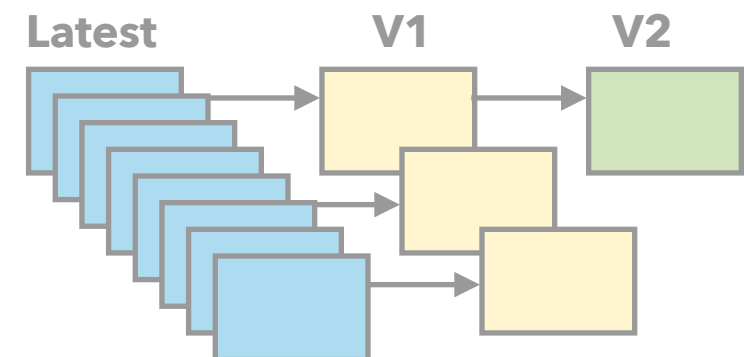
Vectorized Processing



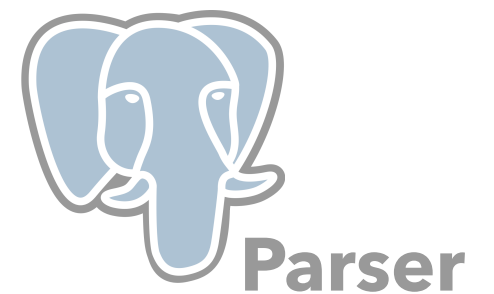
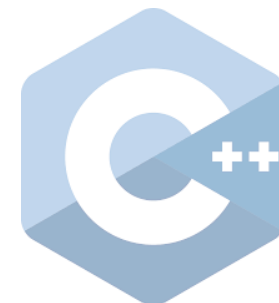
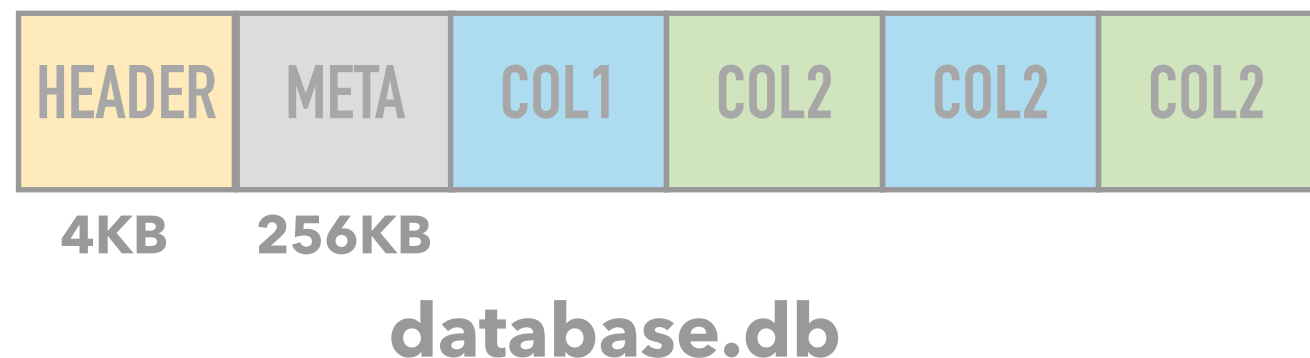
ART Index



Multi-Version Concurrency Control

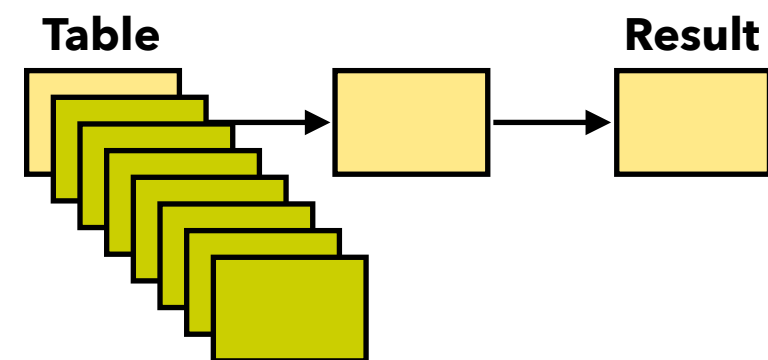


Single-File Storage

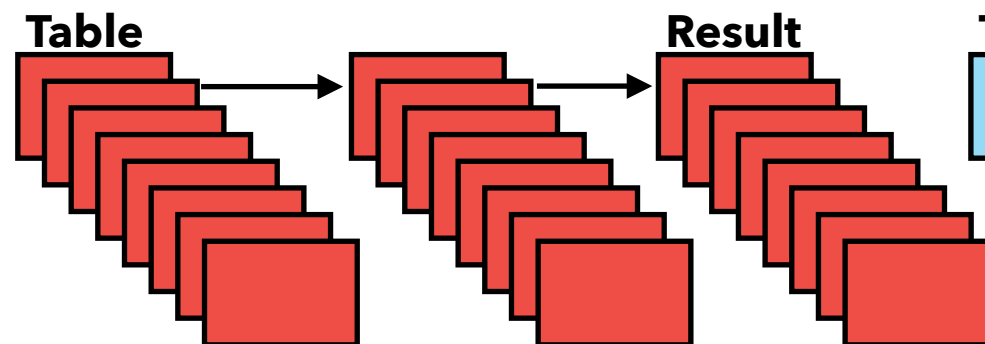


- ▶ **Query Execution**
- ▶ **Traditional RDBMS** use tuple-at-a-time processing
 - ▶ Process **one row** at a time
- ▶ **NumPy/R** use column-at-a-time processing
 - ▶ Process entire columns at once
- ▶ **DuckDB** uses **vectorized** processing
 - ▶ Process **batches** of columns at once

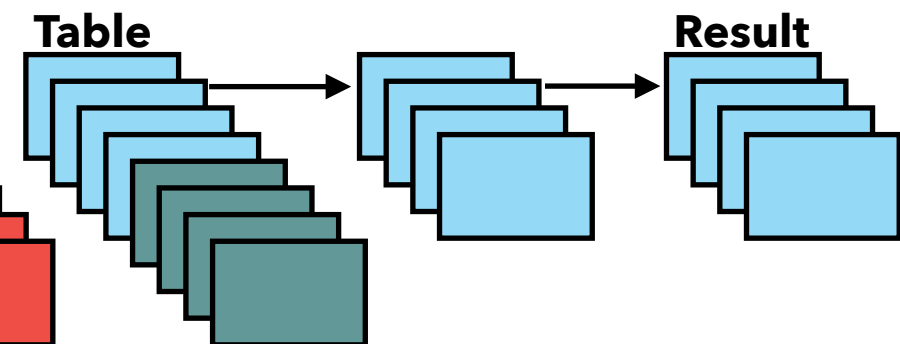
Tuple-at-a-Time



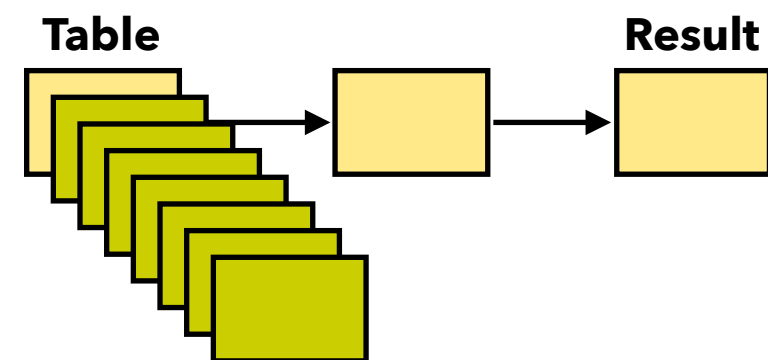
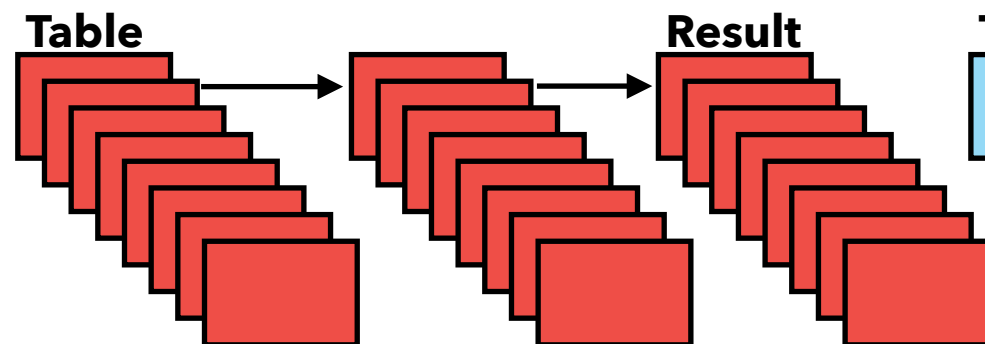
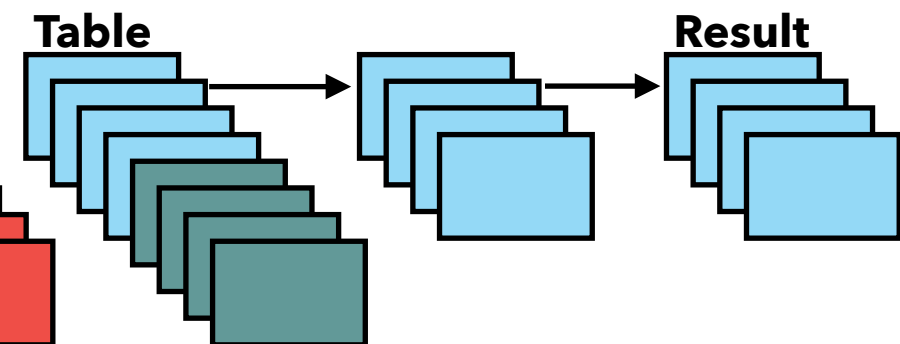
Column-at-a-Time



Vectorized Processing

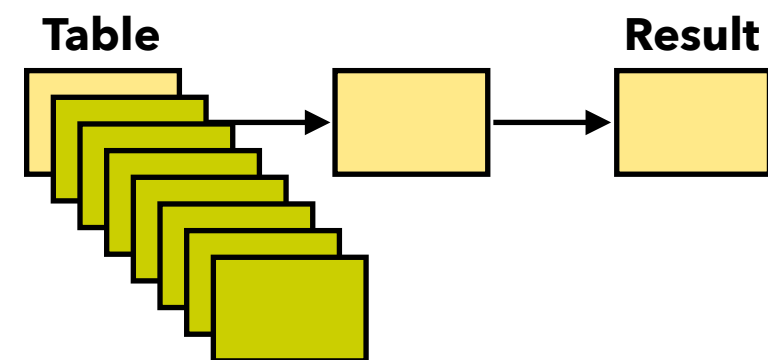


- ▶ **Tuple-at-a-Time (Traditional RDBMS)**
 - ▶ Optimize for low memory footprint
 - ▶ Only need to keep **single row** in memory
- ▶ Comes from a time when **memory was expensive**
- ▶ **High CPU overhead per tuple!**

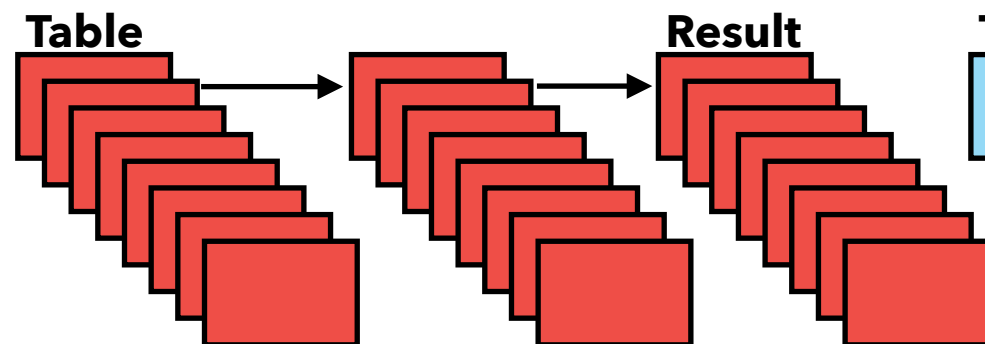
Tuple-at-a-Time**Column-at-a-Time****Vectorized Processing**

- ▶ **Column-at-a-Time (NumPy/R)**
 - ▶ Better CPU utilization, allows for SIMD
 - ▶ Materialize **large intermediates** in memory!
- ▶ Intermediates can be gigabytes each...
- ▶ **Problematic** when data sizes are large

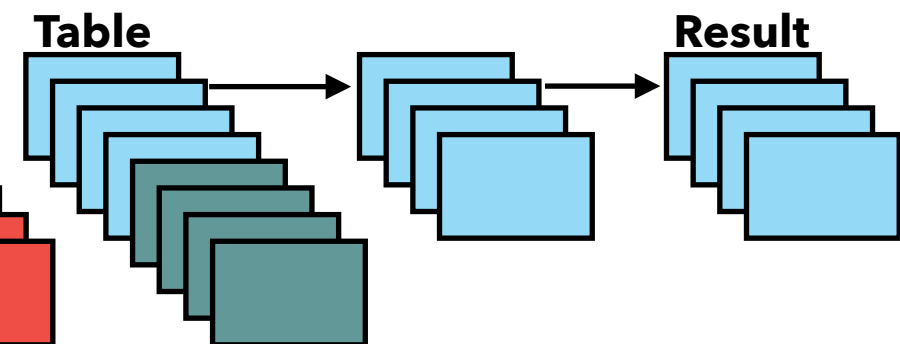
Tuple-at-a-Time



Column-at-a-Time

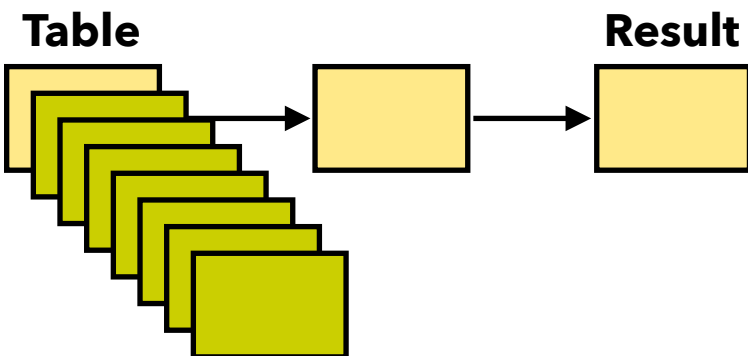


Vectorized Processing

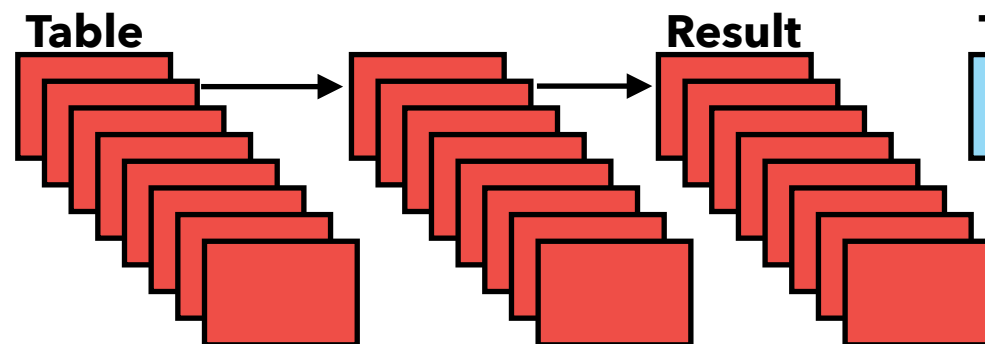


- ▶ **Vectorized Processing (DuckDB)**
 - ▶ Optimized for CPU Cache locality
 - ▶ SIMD instructions, Pipelining
 - ▶ **Small intermediates (fit in L3 cache)**

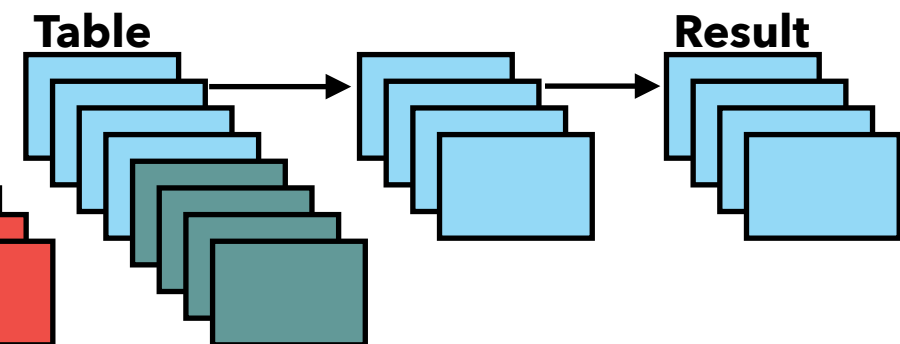
Tuple-at-a-Time



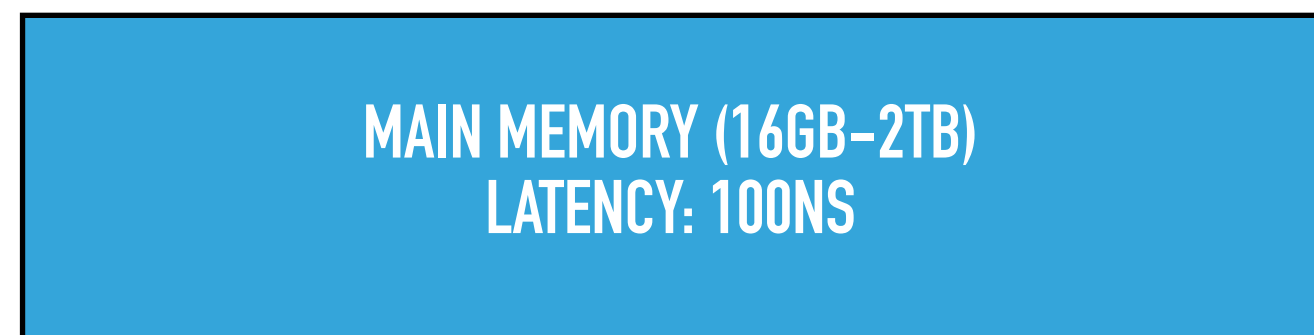
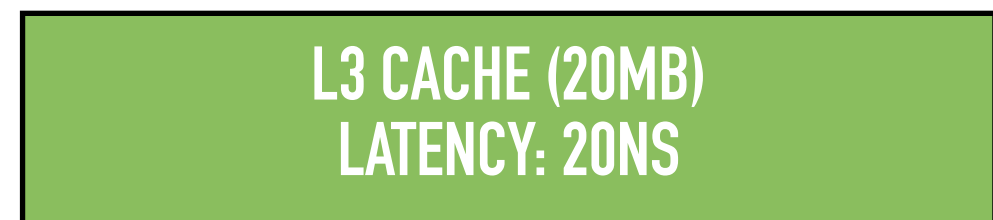
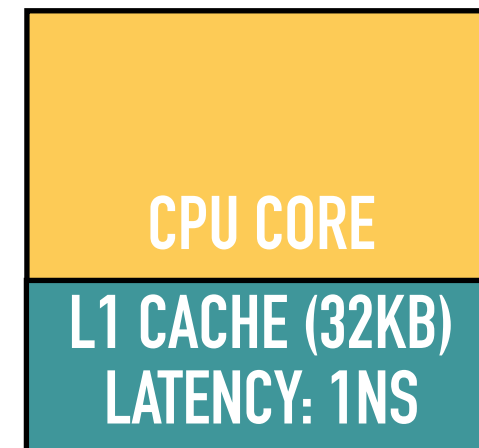
Column-at-a-Time



Vectorized Processing



- ▶ **Vectorized Processing**
- ▶ Intermediates fit in L3 cache
- ▶ **Column-at-a-Time**
- ▶ Intermediates go to memory



- ▶ Internals at a Glance
- ▶ **Query processing pipeline**
- ▶ Query execution
- ▶ Hands-On

▶ **Life of a query**

▶ How does the system go from query to result?

▶ We will focus on the following query:

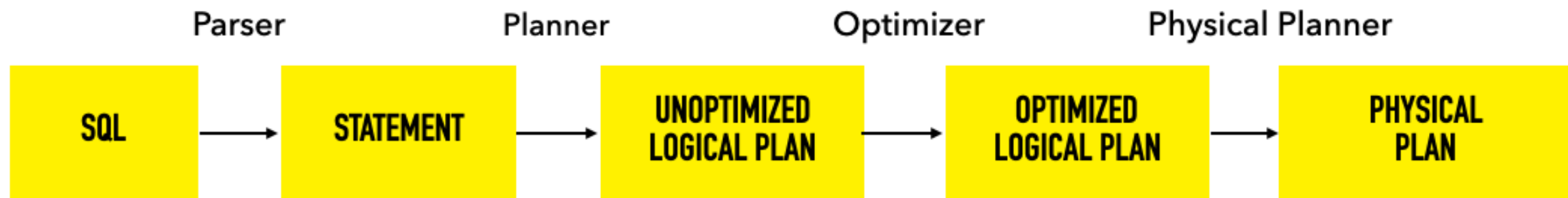
```
SELECT COUNT(*)  
FROM lineitem, orders  
WHERE l_orderkey=o_orderkey  
      AND o_orderstatus='F'  
      AND l_tax > 0.04;
```

▶ **Aggregate:** COUNT (*)

▶ **Implicit join:** lineitem, orders **on** orderkey

▶ **Filters:** o_orderstatus='F' **and** l_tax>0.04

DuckDB uses a typical pipeline for query processing

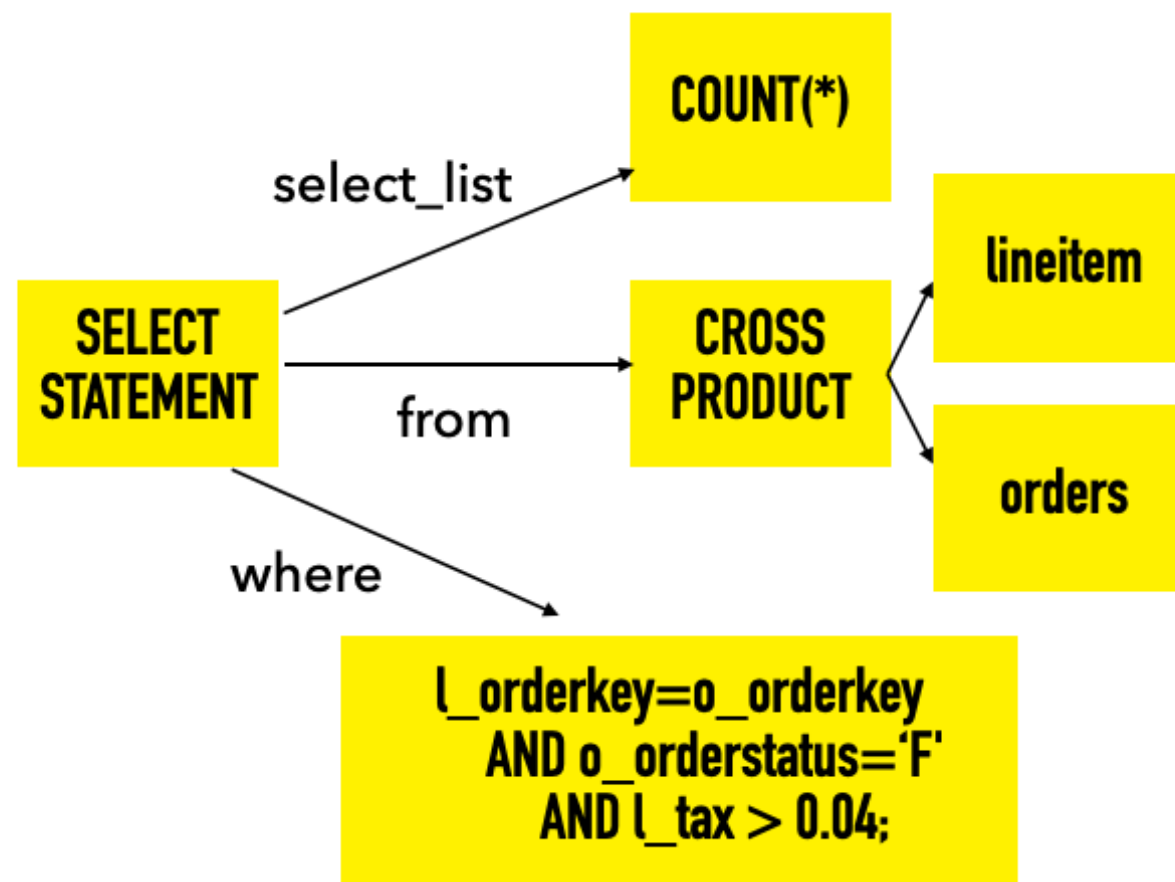




- Query is input into the system as a string
- The lexer and parser take the input string and convert it into a set of **statements, parsed expressions** and **table references**
- Note that this is not yet a query tree!
- We utilize the **Postgres parser + Transformer** for this part

```
SELECT COUNT(*)
FROM lineitem, orders
WHERE l_orderkey=o_orderkey
      AND o_orderstatus='F'
      AND l_tax > 0.04;
```

- The result of the **parsing** stage is the following:



```
SELECT COUNT(*)
FROM lineitem, orders
WHERE l_orderkey=o_orderkey
      AND o_orderstatus='F'
      AND l_tax > 0.04;
```



- Look up tables (**lineitem** and **orders**)
- Look up **columns** within these tables
- Assign **table** and **column** indexes

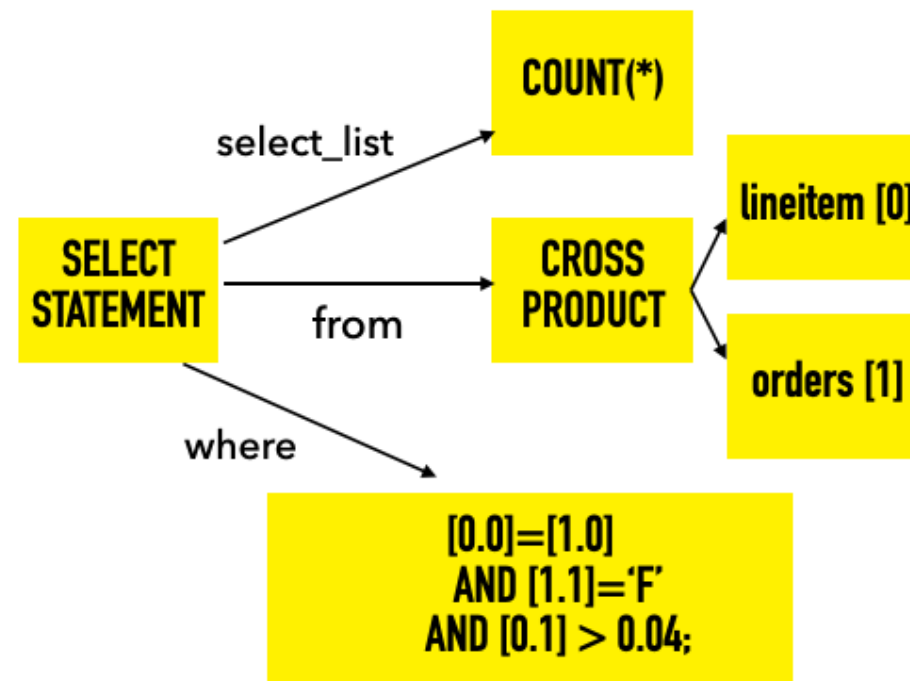


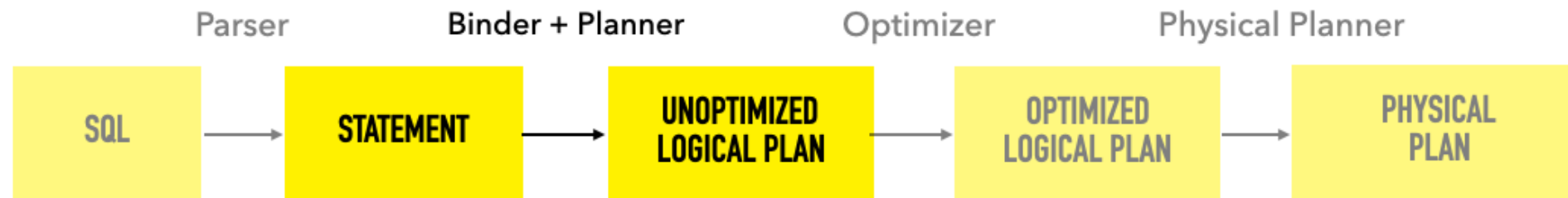
table index
lineitem: 0
orders: 1

column index
l_orderkey: 0
l_tax: 1

column index
o_orderkey: 0
o_orderstatus: 1

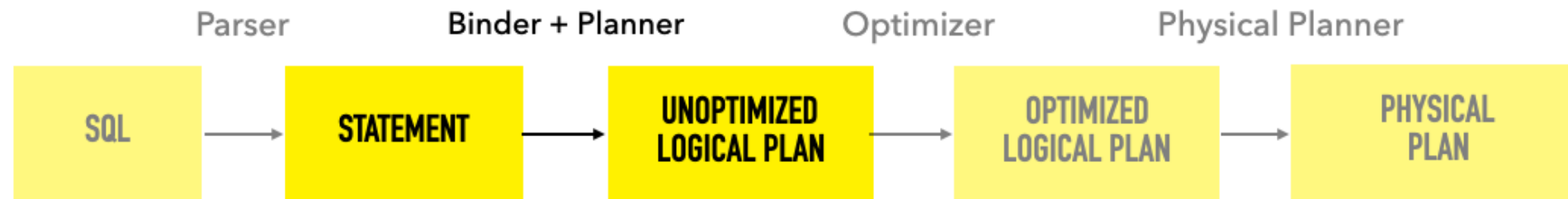
```

SELECT COUNT(*)
FROM lineitem, orders
WHERE l_orderkey=o_orderkey
AND o_orderstatus='F'
AND l_tax > 0.04;
  
```



- **Type resolution** happens in this stage
 - `l_orderkey` : INTEGER
 - `o_orderkey` : INTEGER
- `l_orderkey = o_orderkey` : BOOLEAN
- Types are propagated through every expression/operator
- Types are used for function resolution
 - Potential implicit casts are added

```
SELECT COUNT(*)  
FROM lineitem, orders  
WHERE l_orderkey=o_orderkey  
      AND o_orderstatus='F'  
      AND l_tax > 0.04;
```



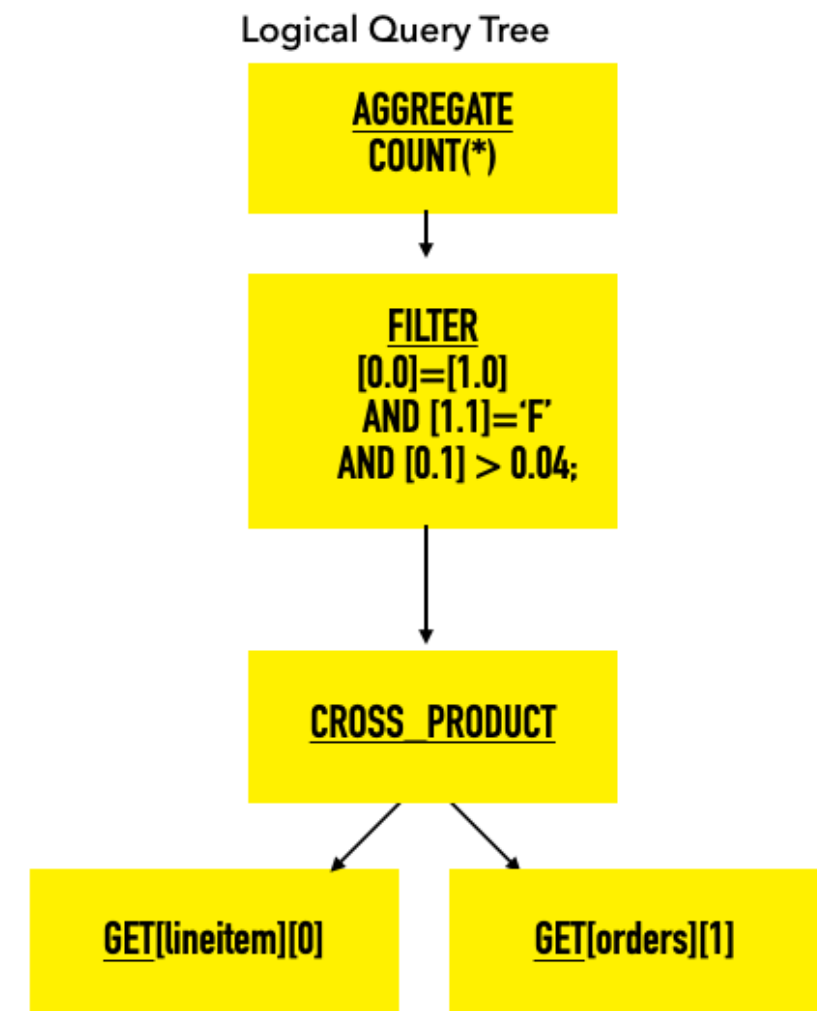
- **Planner:** Create **logical** query tree
- The logical query tree contains **logical operations**
 - Describes what to do, not how to do it
 - e.g. Join, not HashJoin or MergeJoin

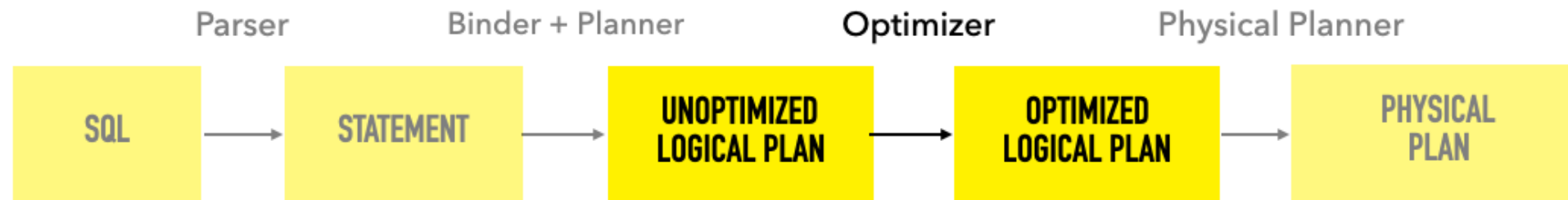
```
SELECT COUNT(*)  
FROM lineitem, orders  
WHERE l_orderkey=o_orderkey  
      AND o_orderstatus='F'  
      AND l_tax > 0.04;
```



- Query tree starts with **FROM** clause
- Followed by **WHERE**
- Followed by **SELECT**

```
SELECT COUNT(*)  
FROM lineitem, orders  
WHERE l_orderkey=o_orderkey  
      AND o_orderstatus='F'  
      AND l_tax > 0.04;
```





- **Optimizer: rewrite logical query tree**

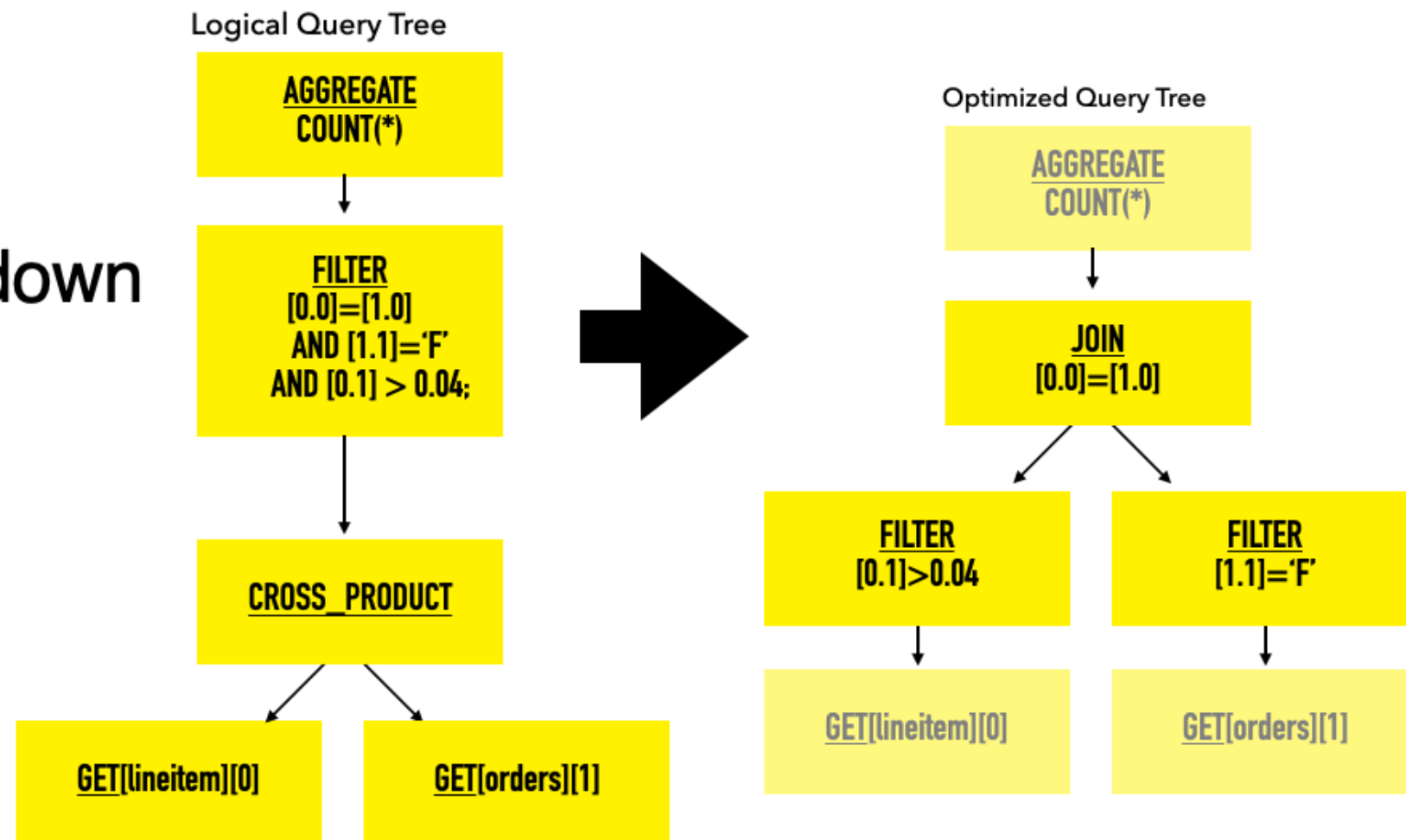
- **Mix of rule-based and cost-based optimizers**

- Join ordering/filter pushdown

- Expression rewriter

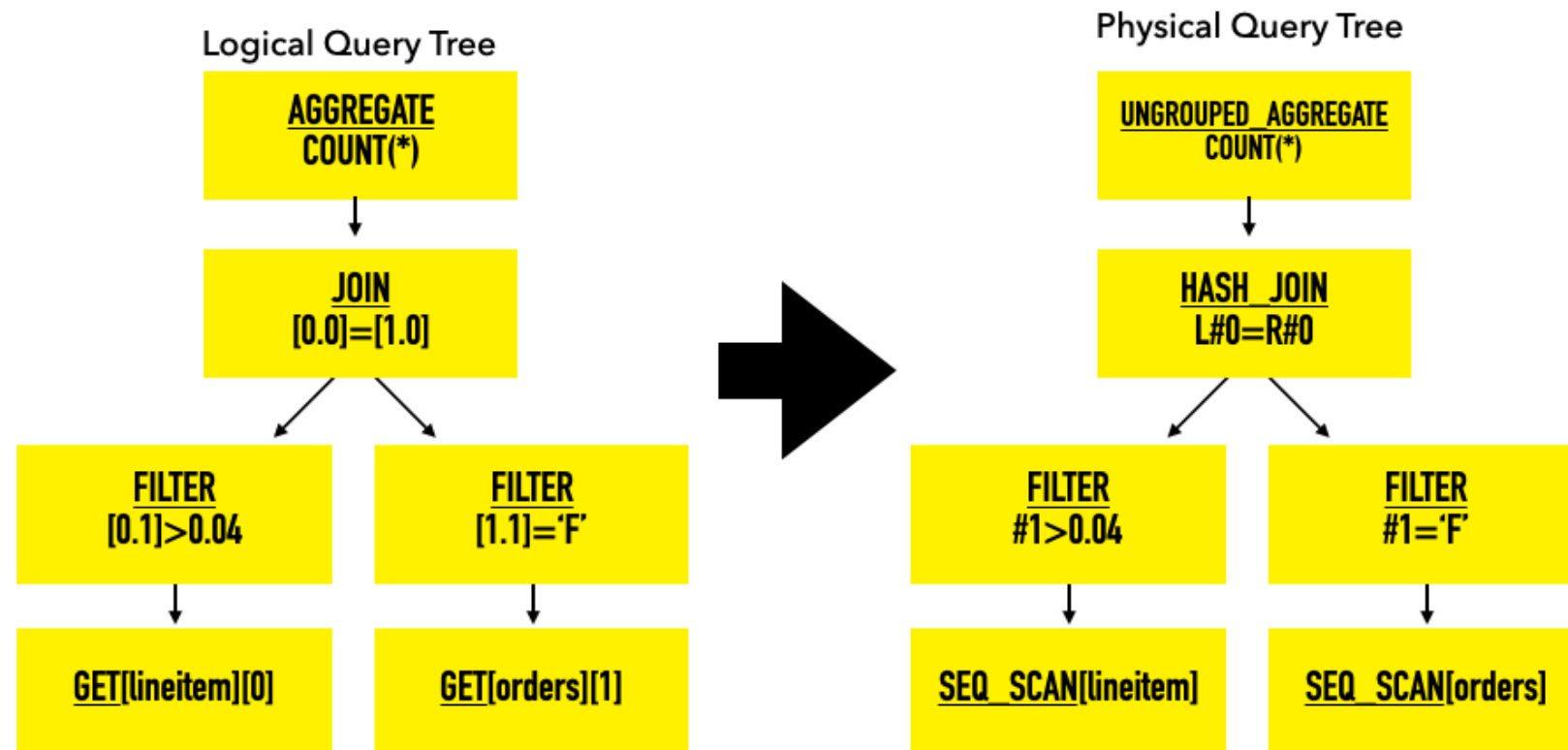
- Statistics propagation

- ...



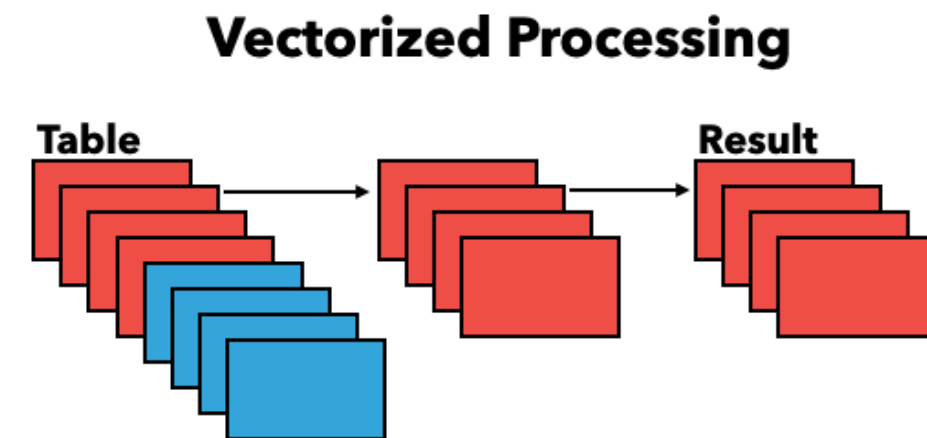


- **Physical planner:** convert logical into physical plan
- Make decisions on implementations of operators
- Convert **column bindings** into **references**

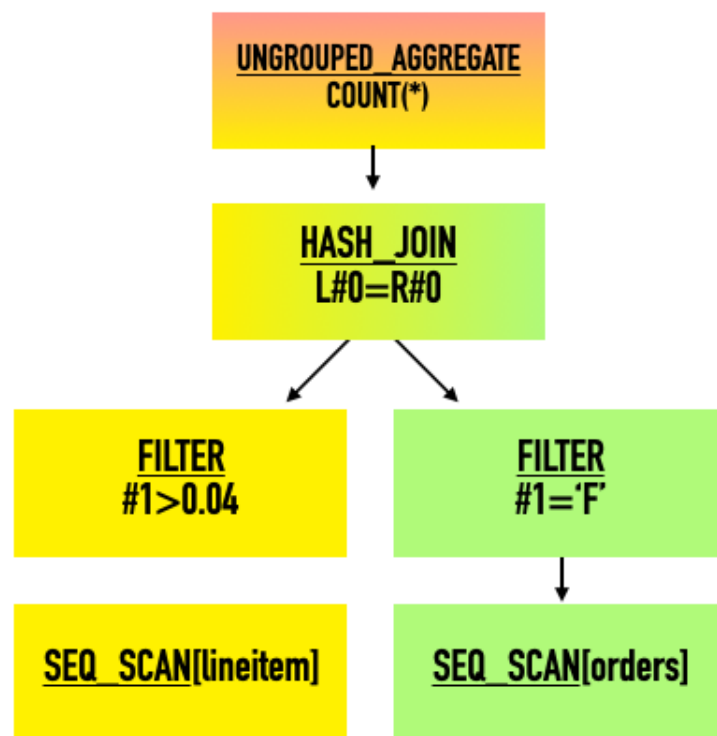


- ▶ Internals at a Glance
- ▶ Query processing pipeline
- ▶ **Query execution**
- ▶ Hands-On

- DuckDB uses a **vectorized push-based model**
- Query tree is divided into **pipelines**
 - Pipelines have a **source**, **operators** and **sink**
- **Data flows from source to sink one chunk at a time**



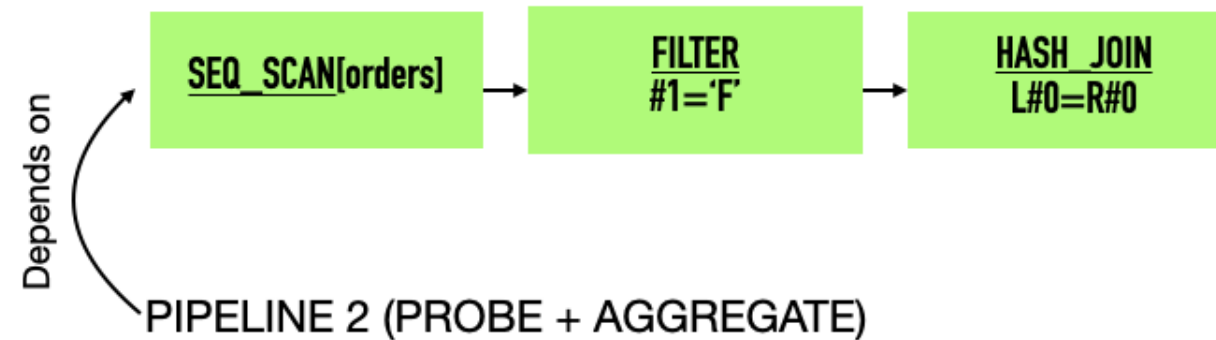
Physical Query Tree



```

SELECT COUNT(*)
FROM lineitem, orders
WHERE l_orderkey=o_orderkey
      AND o_orderstatus='F'
      AND l_tax > 0.04;
  
```

PIPELINE 1 (HT BUILD)



PIPELINE 2 (PROBE + AGGREGATE)



PIPELINE 3 (RESULT)



PIPELINE 1 (HT BUILD)



```
SELECT COUNT(*)  
FROM lineitem, orders  
WHERE l_orderkey=o_orderkey  
      AND o_orderstatus='F'  
      AND l_tax > 0.04;
```

SEQ_SCAN[orders]

Output

o_orderkey	1	o_orderstatus	0
	2		0
	3		F

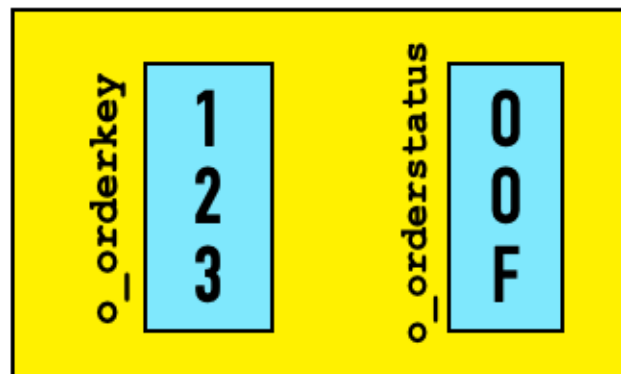
PIPELINE 1 (HT BUILD)



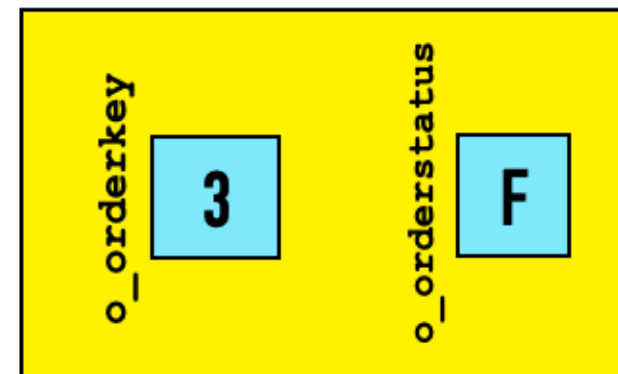
```
SELECT COUNT(*)  
FROM lineitem, orders  
WHERE l_orderkey=o_orderkey  
      AND o_orderstatus='F'  
      AND l_tax > 0.04;
```



Input



Output



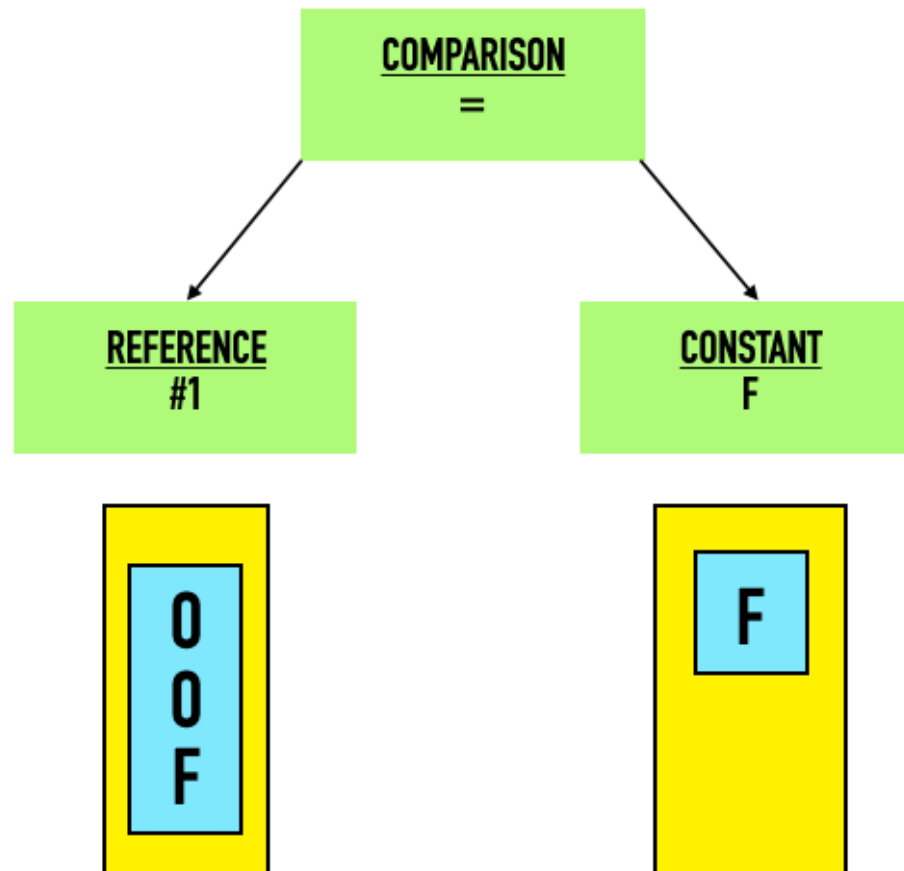
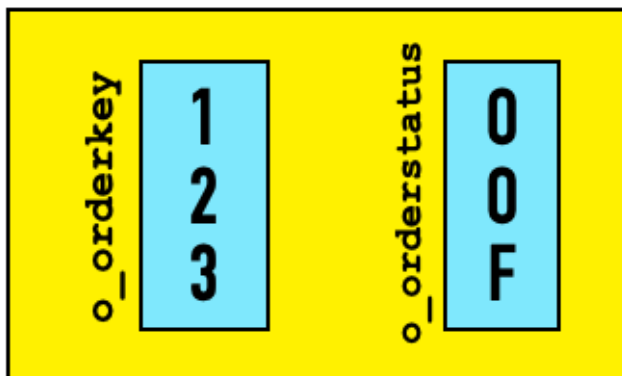
PIPELINE 1 (HT BUILD)



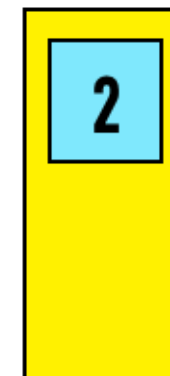
```
SELECT COUNT(*)  
FROM lineitem, orders  
WHERE l_orderkey=o_orderkey  
AND o_orderstatus='F'  
AND l_tax > 0.04;
```



Input



Matches



SelectionVector

PIPELINE 1 (HT BUILD)

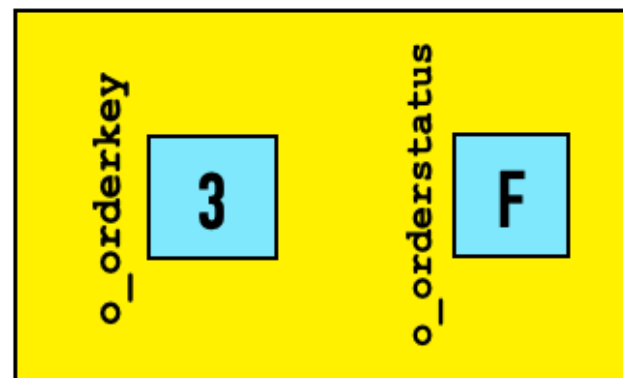
SEQ_SCAN[orders]

FILTER
#1='F'HASH_JOIN
L#0=R#0

```
SELECT COUNT(*)  
FROM lineitem, orders  
WHERE l_orderkey=o_orderkey  
      AND o_orderstatus='F'  
      AND l_tax > 0.04;
```

HASH_JOIN
L#0=R#0

Sink



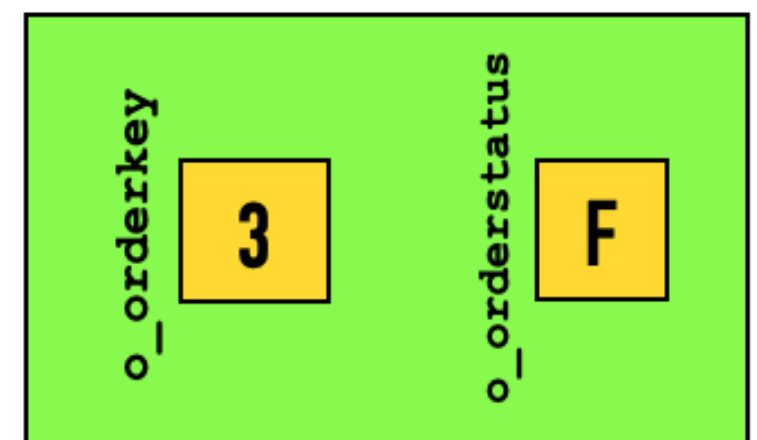
SEQ SCAN[lineitem]

FILTER
#1>0.04

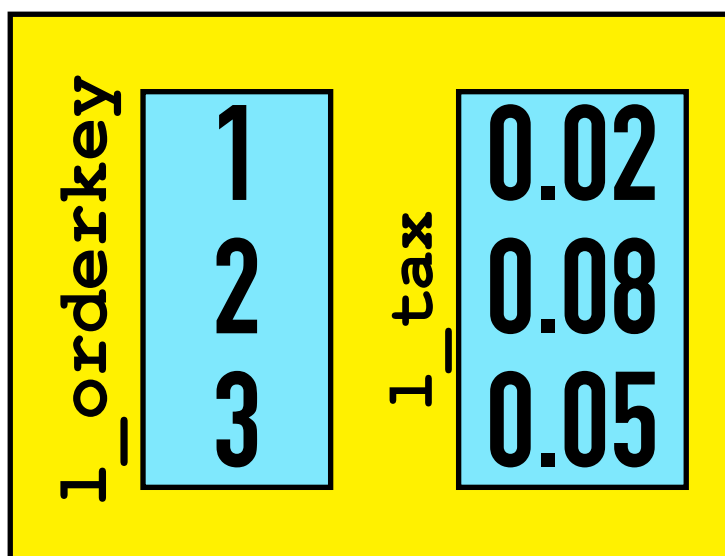
HASH JOIN
L#0=R#0

UNGROUPED AGGREGATE
COUNT(*)

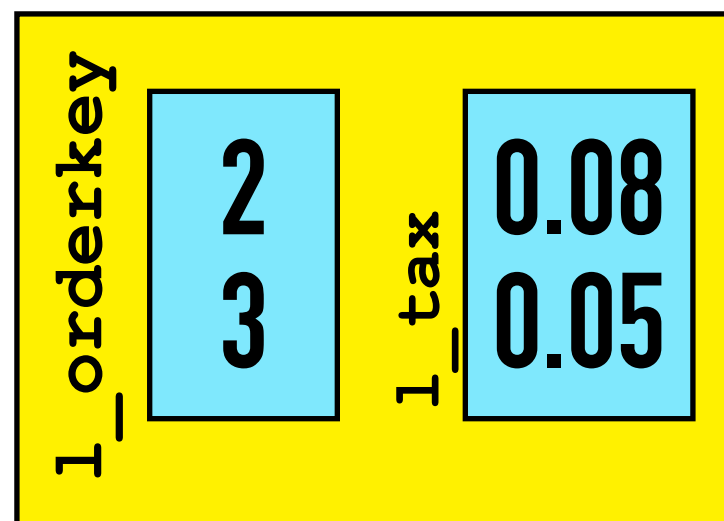
HashTable



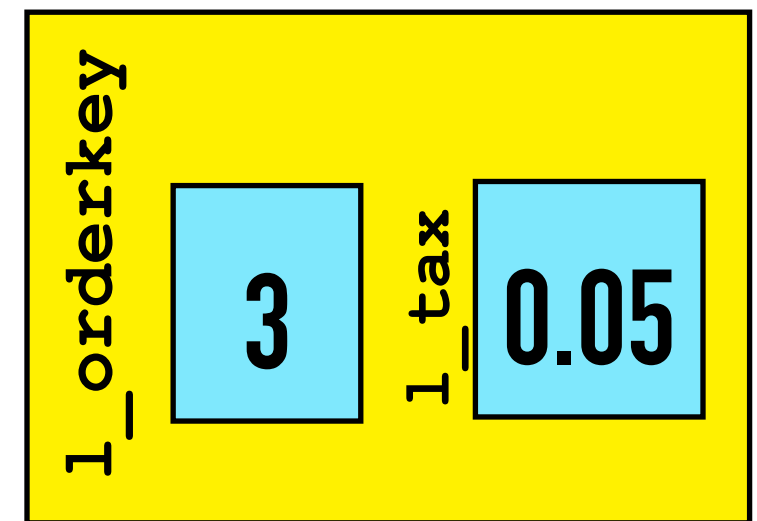
Scan



Filter



Join



- ▶ Internals at a Glance
- ▶ Query processing pipeline
- ▶ Query execution
- ▶ **Hands-On**

- ▶ **Slides are online**
- ▶ <https://github.com/pdet/duckdb-tutorial>
- ▶ Feel free to ask any questions!