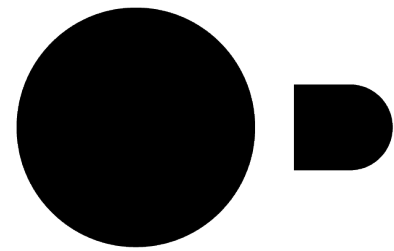


Mark Raasveldt & Pedro Holanda

DuckDB **an Embeddable Analytical RDBMS**

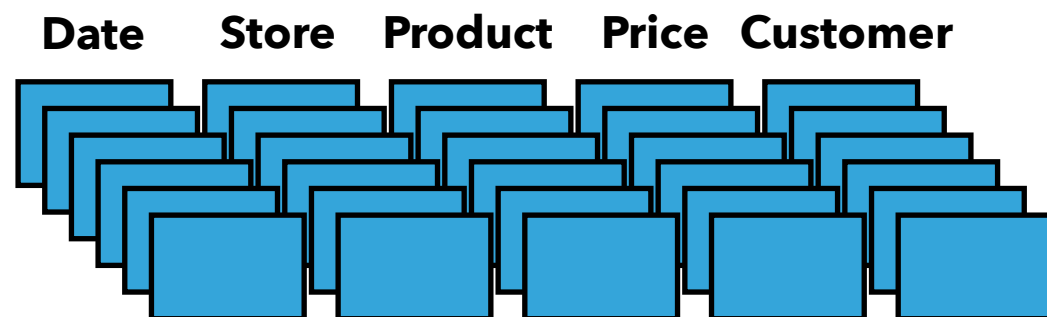
- ▶ **Internals at a Glance**
- ▶ Query processing pipeline
- ▶ Query execution
- ▶ Hands-On



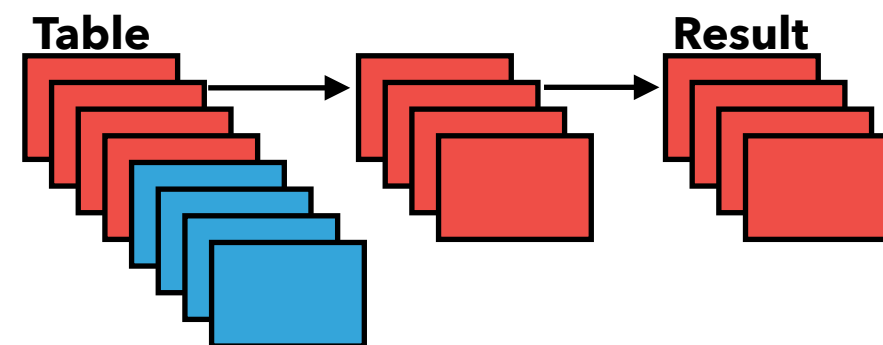
DuckDB

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

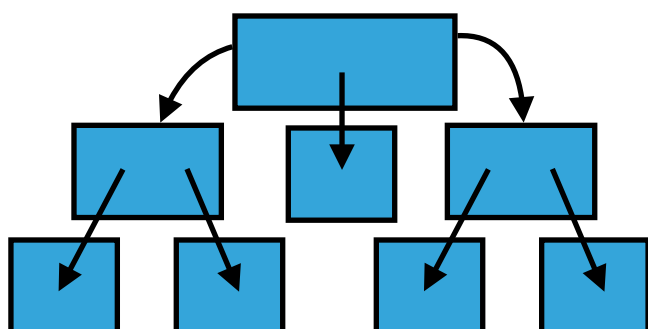
Column-Store



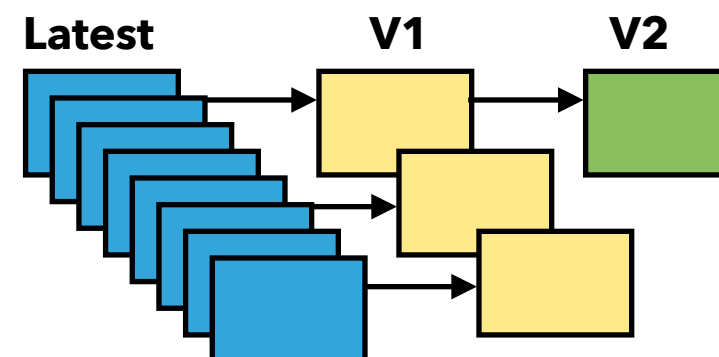
Vectorized Processing



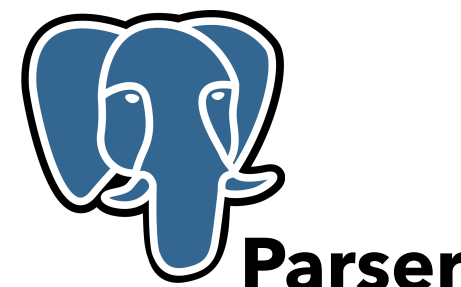
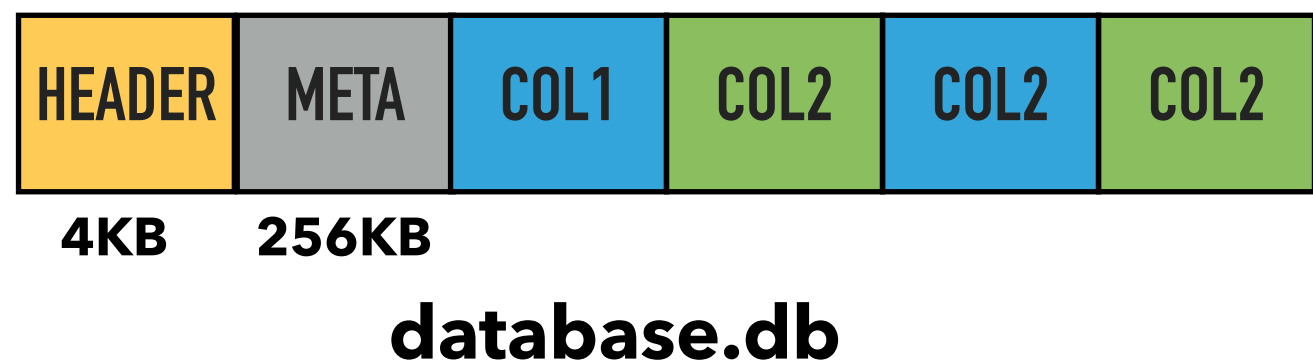
ART Index



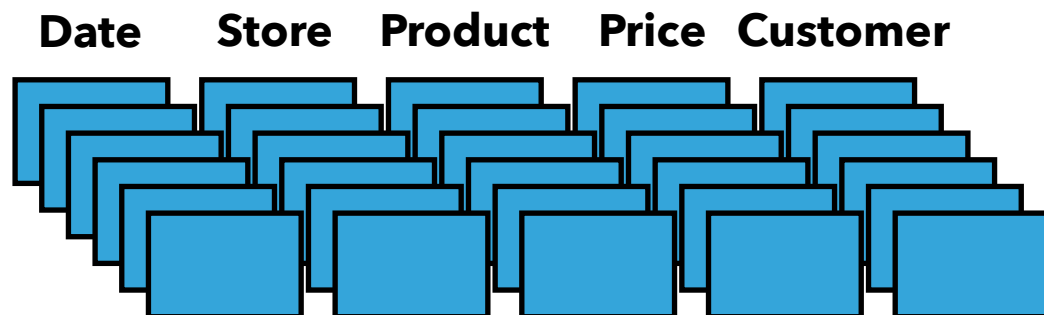
Multi-Version Concurrency Control



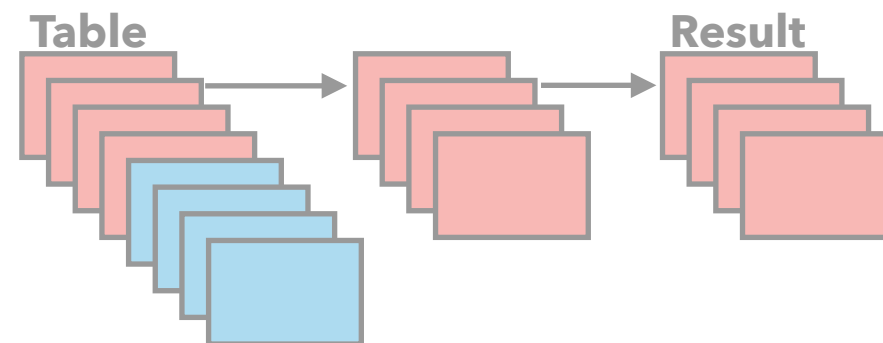
Single-File Storage



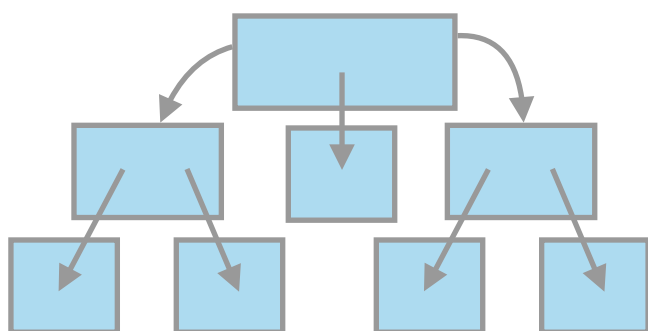
Column-Store



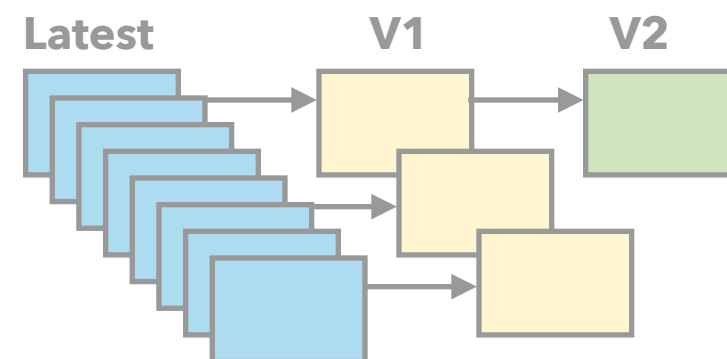
Vectorized Processing



ART Index



Multi-Version Concurrency Control



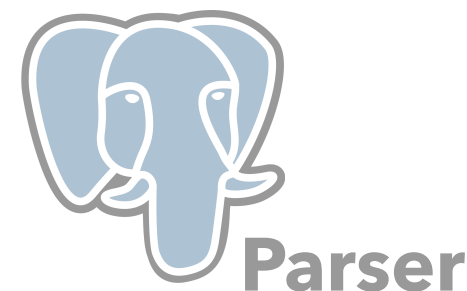
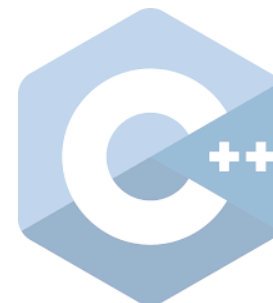
Single-File Storage



4KB

256KB

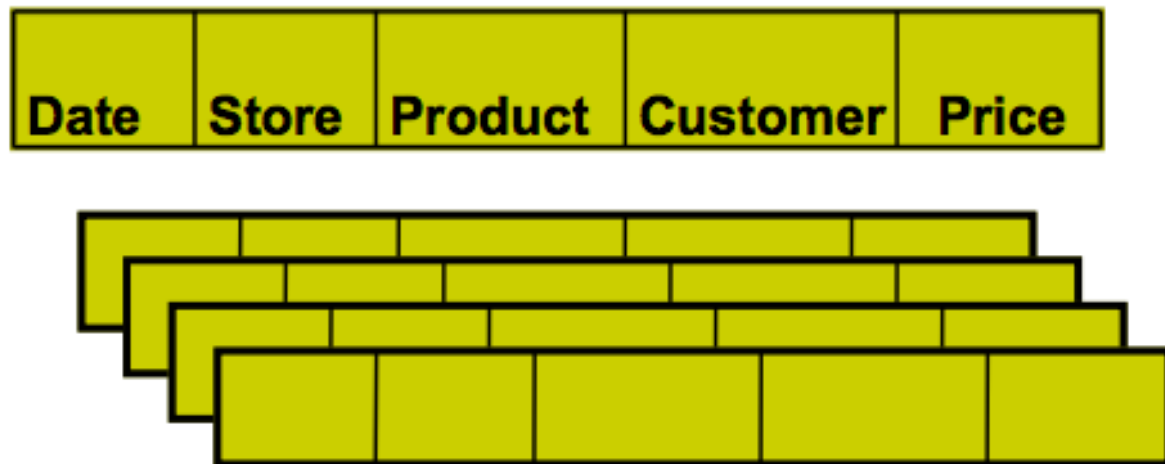
database.db



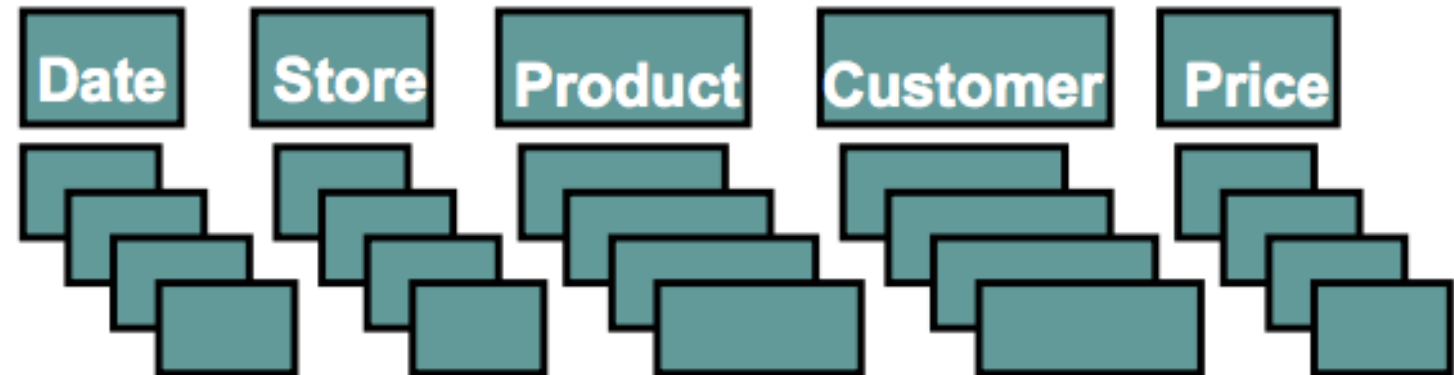
► Storage Model

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

row-store



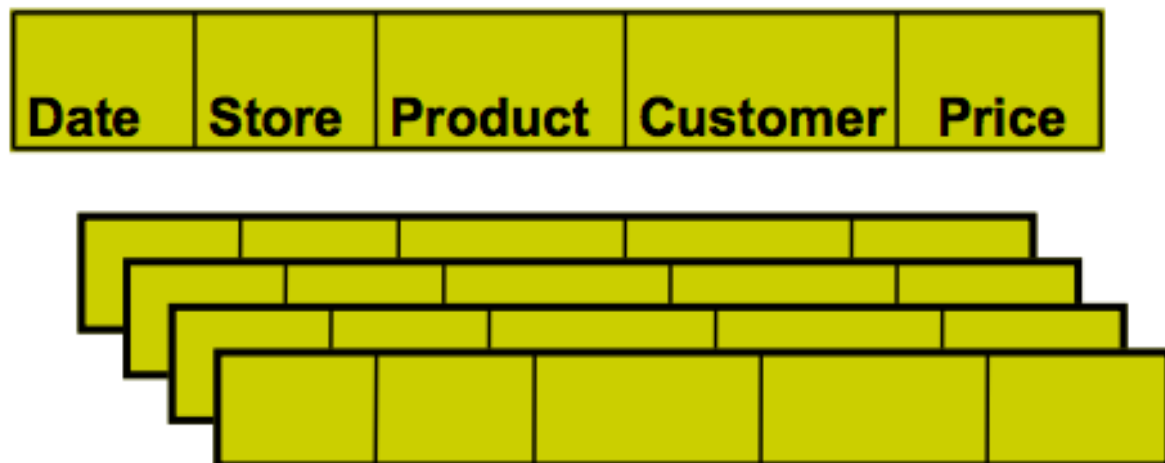
column-store



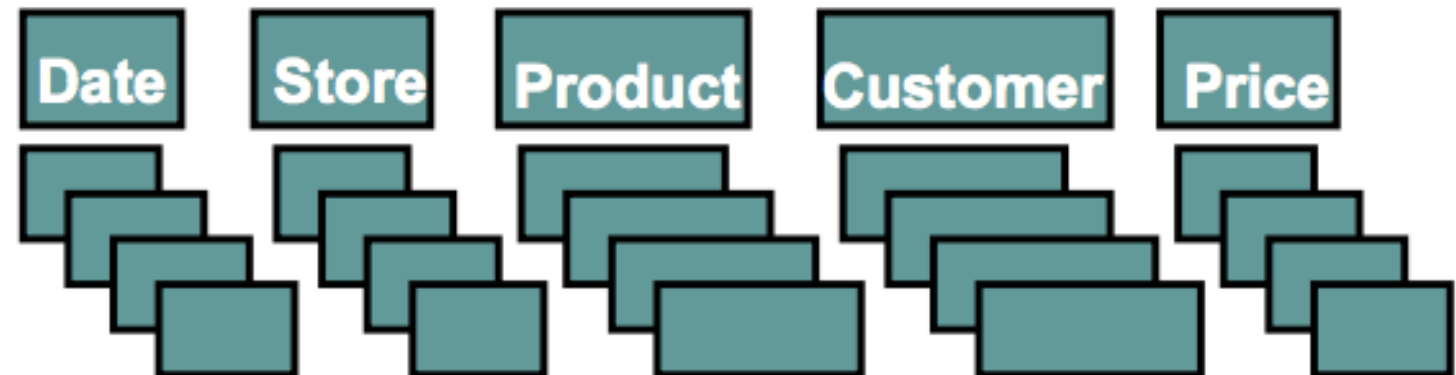
▶ Row-Storage:

- ▶ 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



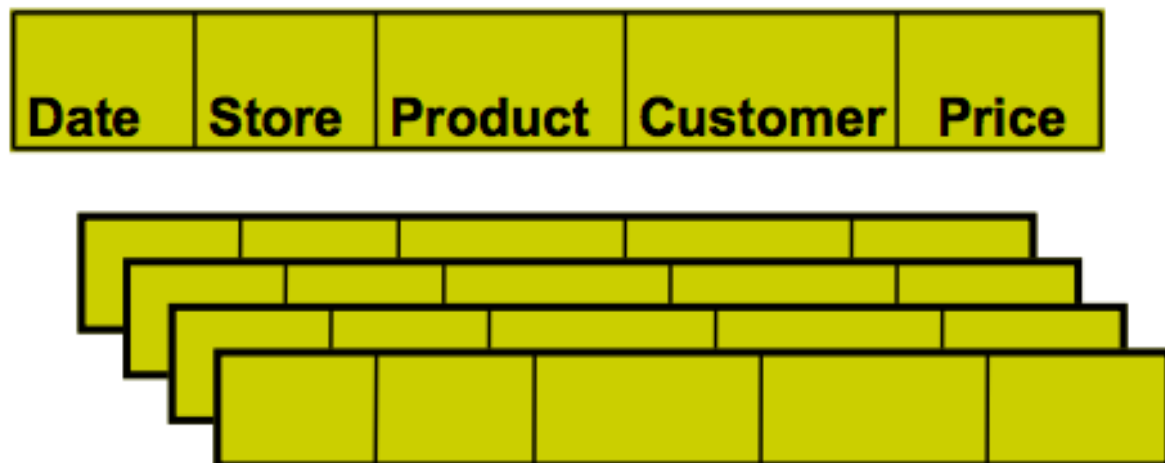
column-store



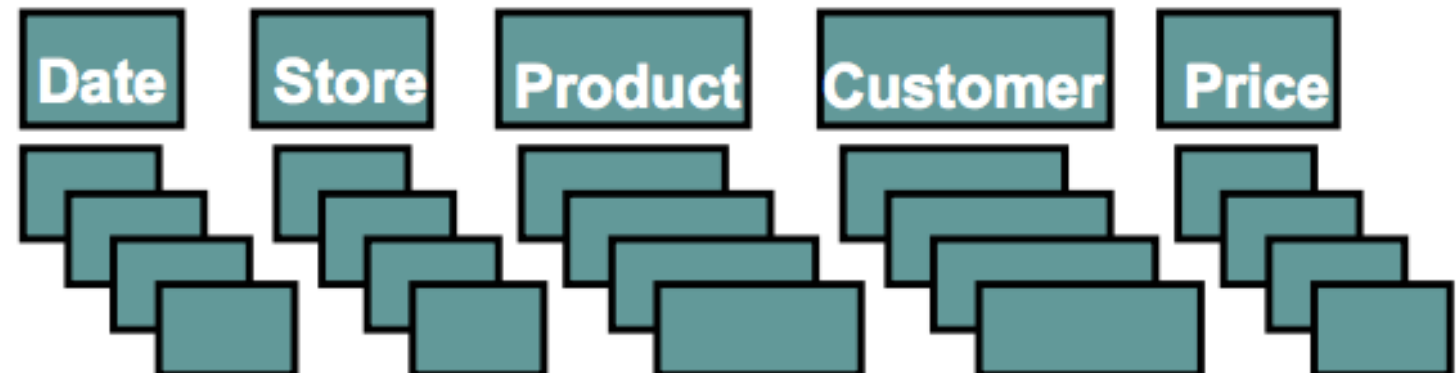
▶ Column-Storage:

- ▶ 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



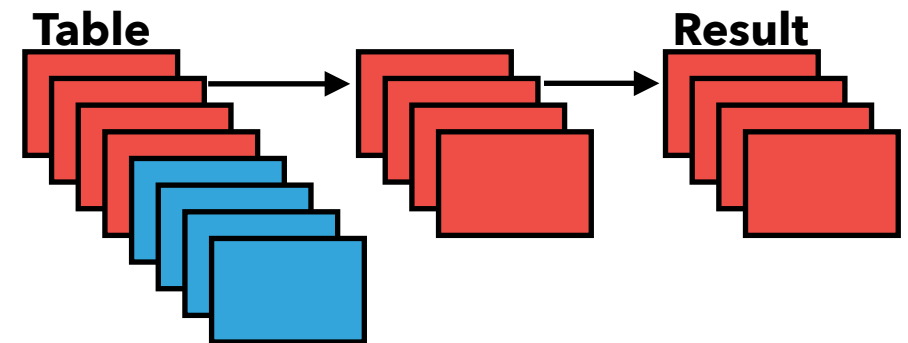
column-store



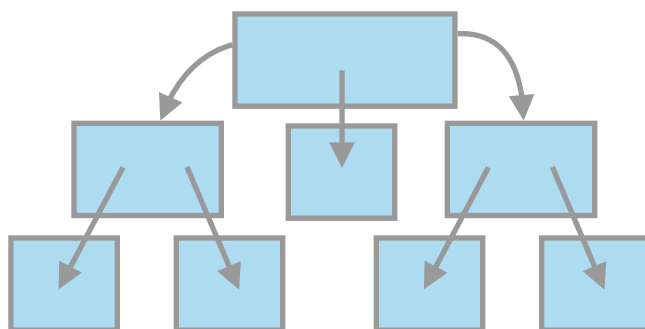
Column-Store



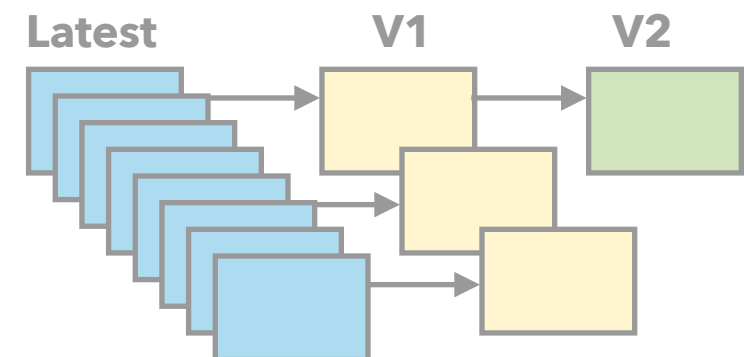
Vectorized Processing



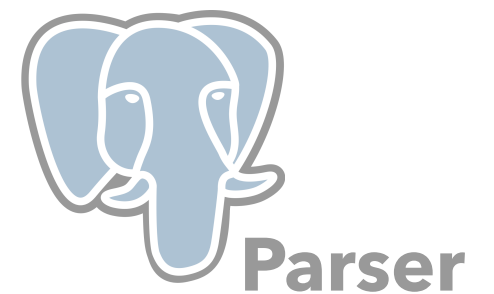
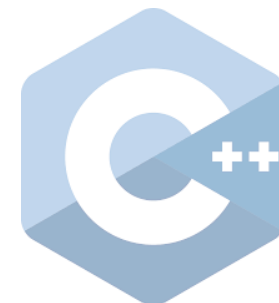
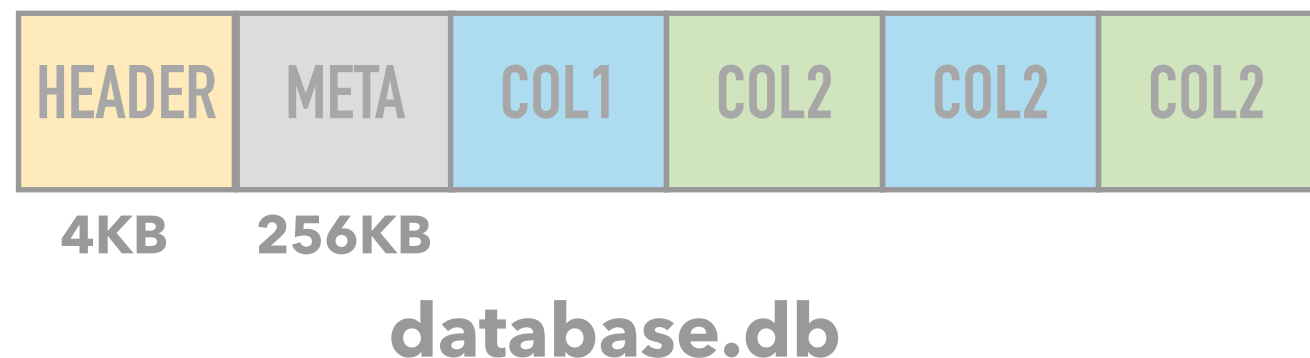
ART Index



Multi-Version Concurrency Control

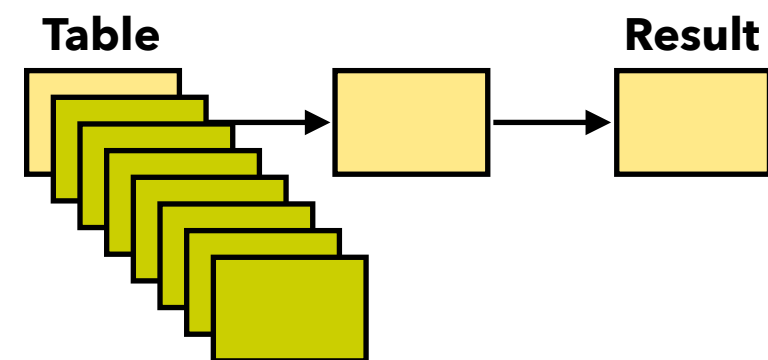


Single-File Storage

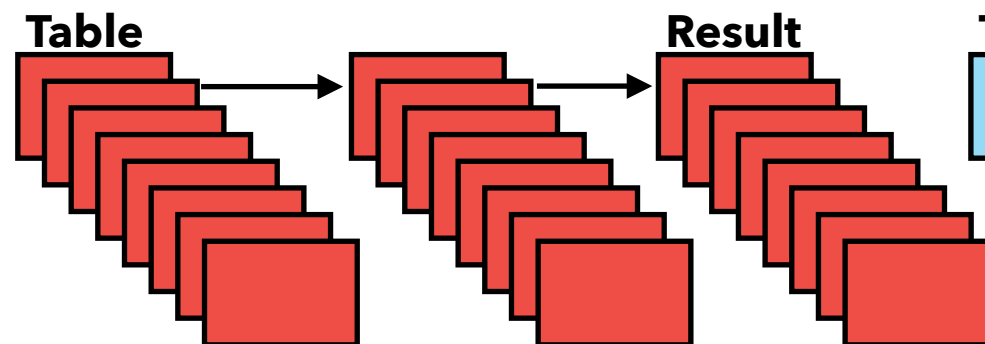


- ▶ **Query Execution**
- ▶ **Traditional RDBMS** use tuple-at-a-time processing
 - ▶ Process **one row** at a time
- ▶ **NumPy/R** use column-at-a-time processing
 - ▶ Process entire columns at once
- ▶ **DuckDB** uses **vectorized** processing
 - ▶ Process **batches** of columns at once

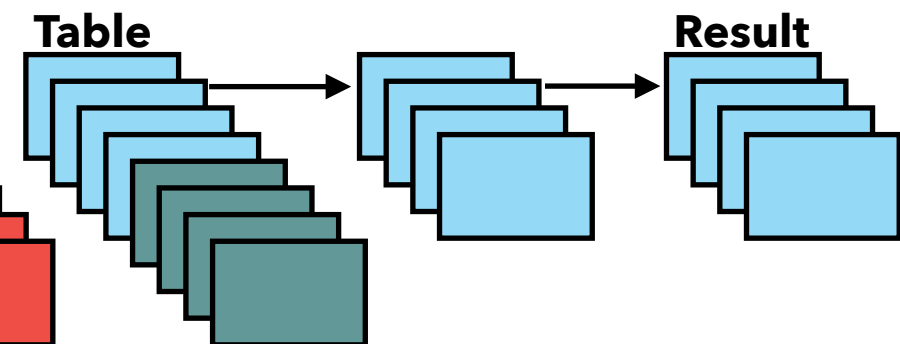
Tuple-at-a-Time



Column-at-a-Time

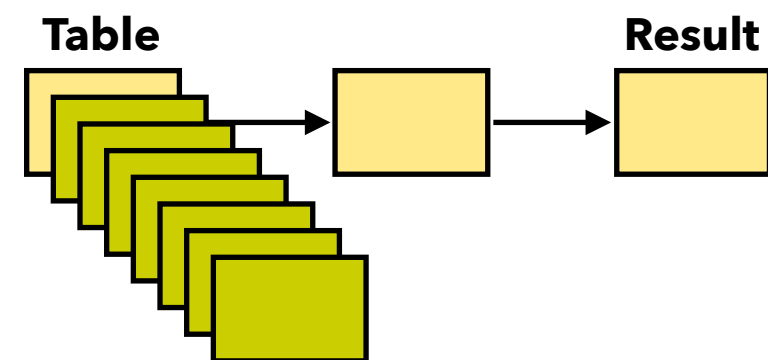


Vectorized Processing

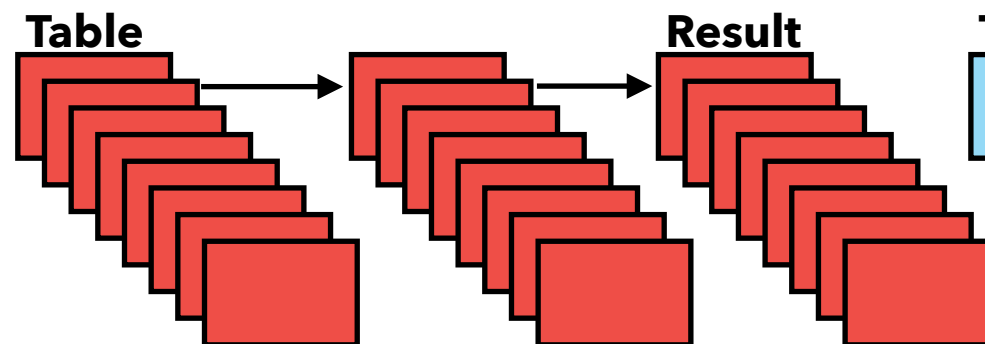


- ▶ **Tuple-at-a-Time (Traditional RDBMS)**
 - ▶ Optimize for low memory footprint
 - ▶ Only need to keep **single row** in memory
- ▶ Comes from a time when **memory was expensive**
- ▶ **High CPU overhead per tuple!**

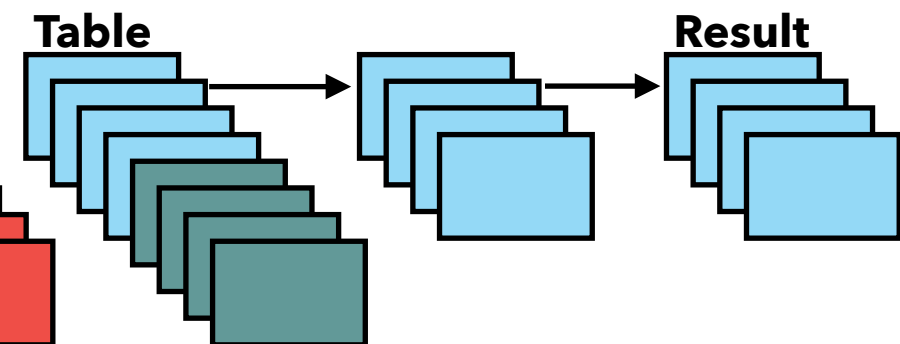
Tuple-at-a-Time



Column-at-a-Time

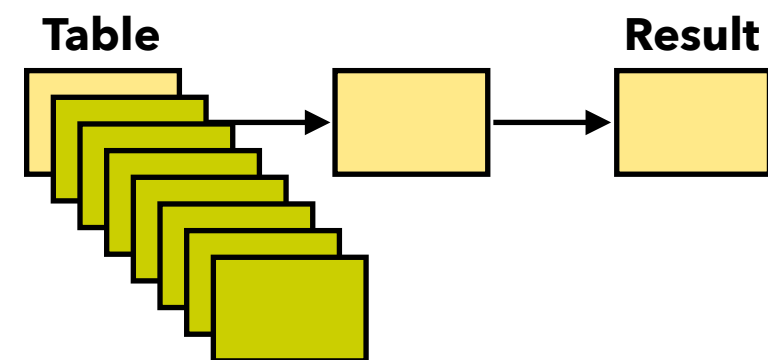


Vectorized Processing

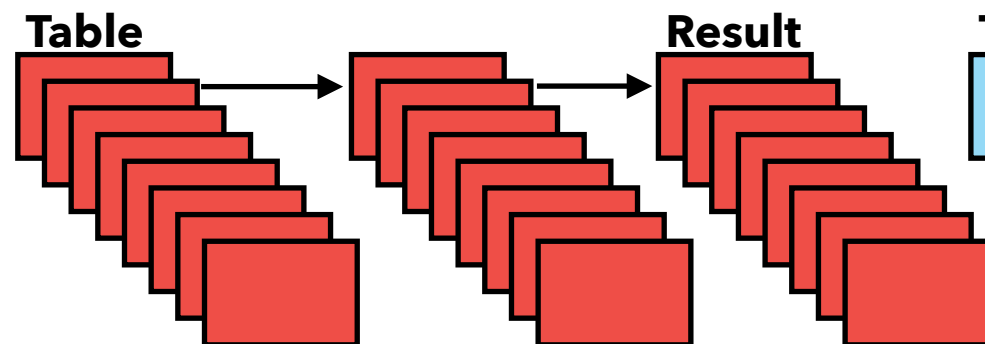


- ▶ **Column-at-a-Time (NumPy/R)**
 - ▶ Better CPU utilization, allows for SIMD
 - ▶ Materialize **large intermediates** in memory!
- ▶ Intermediates can be gigabytes each...
- ▶ **Problematic** when data sizes are large

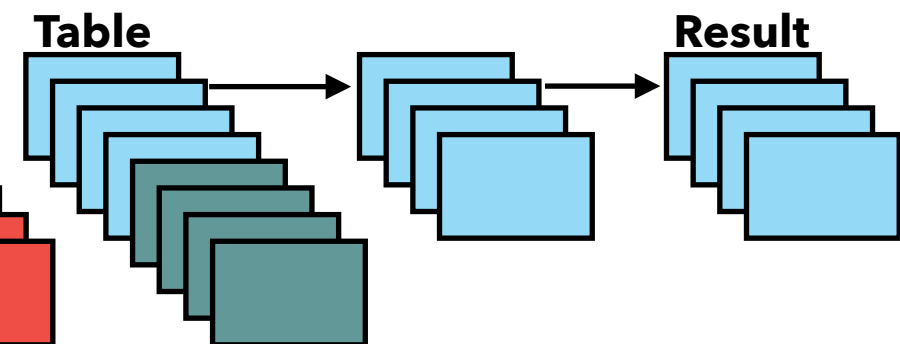
Tuple-at-a-Time



Column-at-a-Time

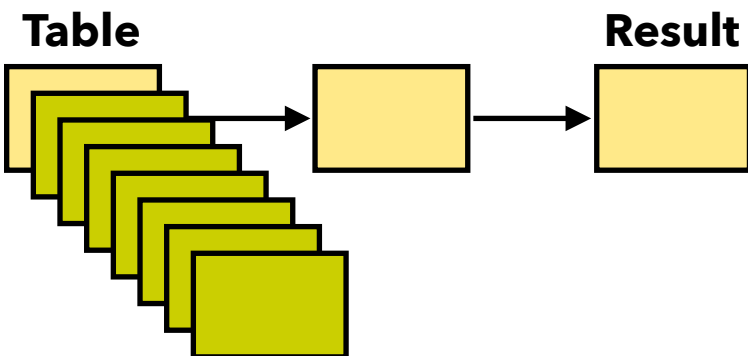


Vectorized Processing

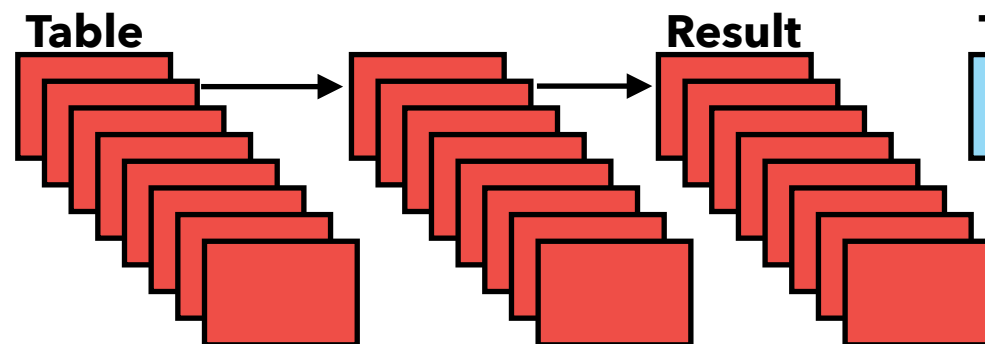


- ▶ **Vectorized Processing (DuckDB)**
 - ▶ Optimized for CPU Cache locality
 - ▶ SIMD instructions, Pipelining
 - ▶ **Small** intermediates (**fit in L3 cache**)

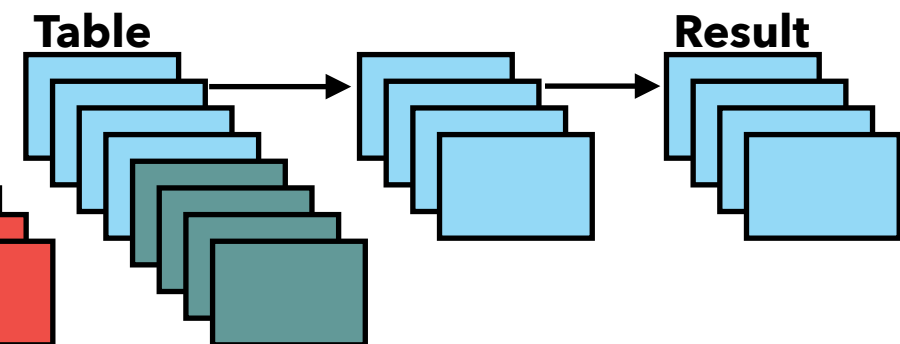
Tuple-at-a-Time



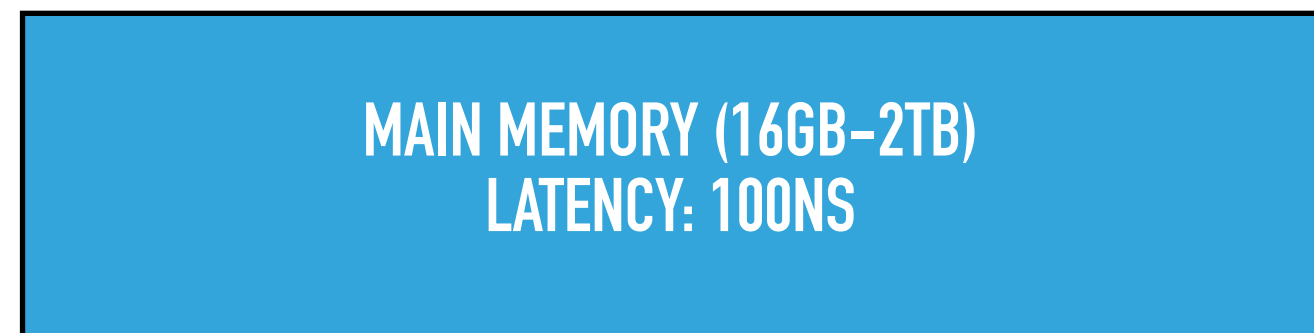
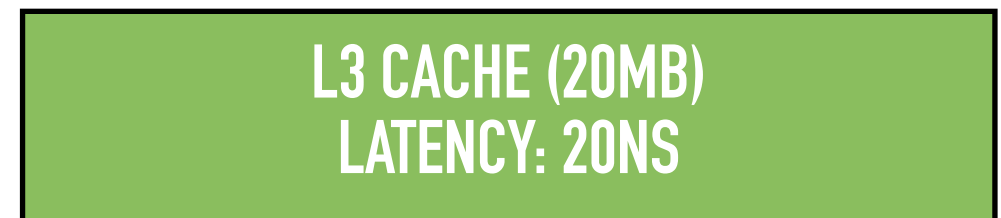
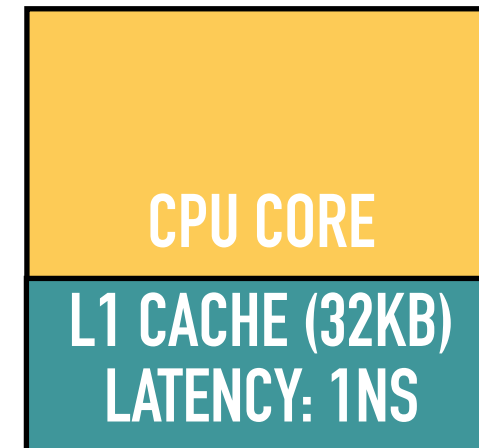
Column-at-a-Time



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:

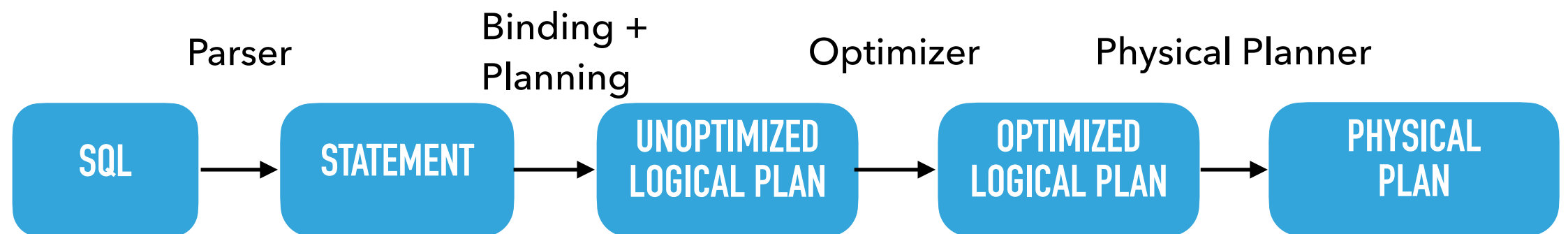
```
SELECT COUNT(*)  
FROM lineitem, orders  
WHERE l_orderkey=o_orderkey  
      AND o_orderstatus='X'  
      AND l_tax > 50;
```

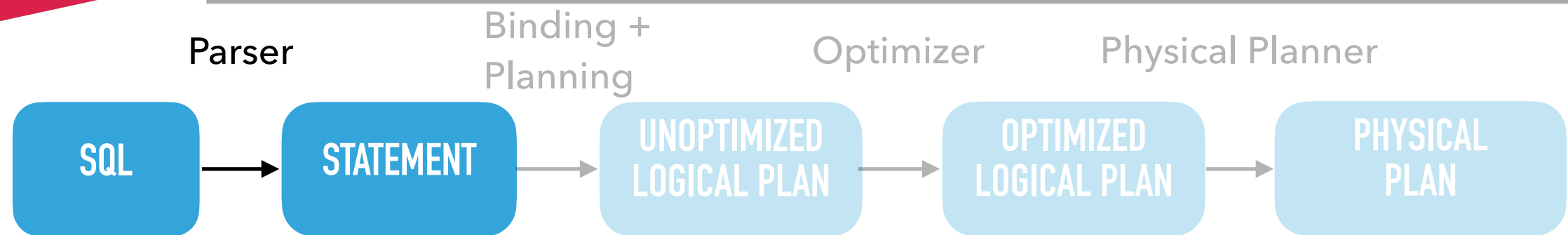
▶ **Aggregate:** COUNT (*)

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

▶ **Filters:** o_orderstatus='X' **and** l_tax>50

- DuckDB uses a typical pipeline for query processing

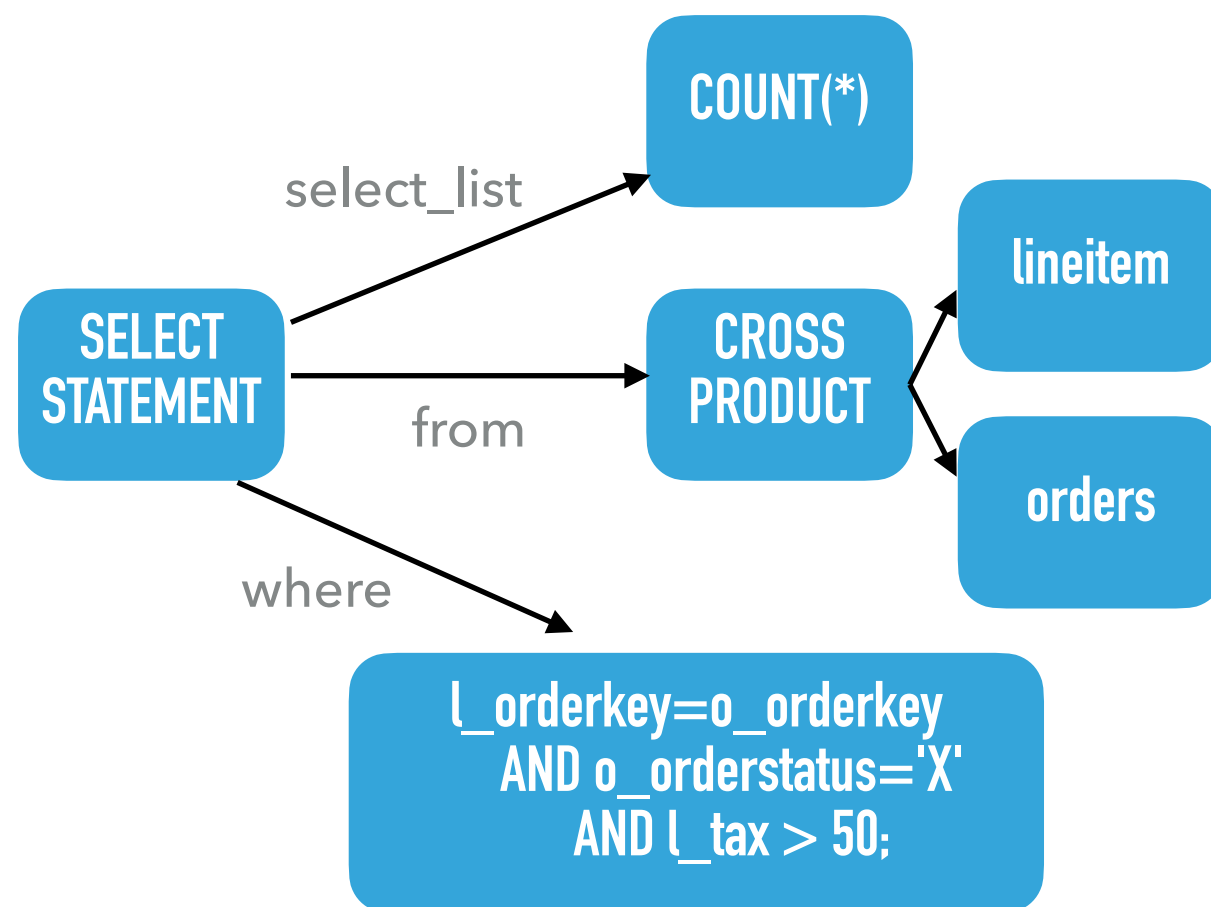




- ▶ Query is input into the system as a string
- ▶ The **lexer and parser** take the input string and convert it into a set of **statements, expressions** and **table references**
 - ▶ Note that this is not yet a query tree!
- ▶ We utilize the **Postgres parser** for this part

```
SELECT COUNT(*)  
FROM lineitem, orders  
WHERE l_orderkey=o_orderkey  
      AND o_orderstatus='X'  
      AND l_tax > 50;
```

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



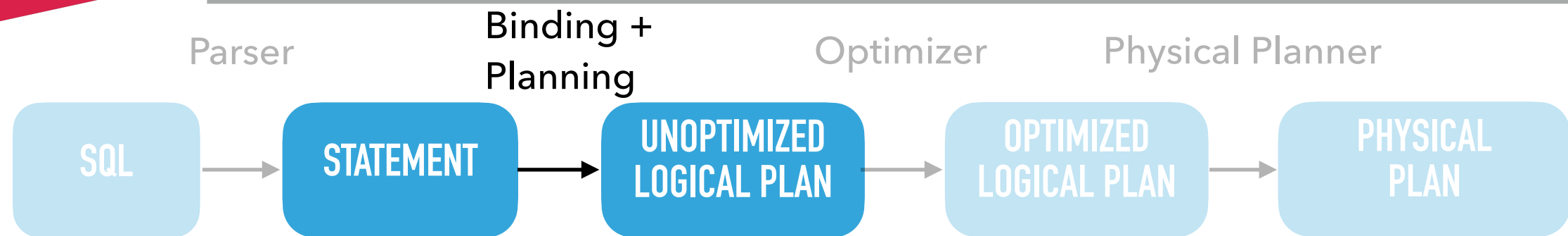
```
SELECT COUNT(*)
FROM lineitem, orders
WHERE l_orderkey=o_orderkey
      AND o_orderstatus='X'
      AND l_tax > 50;
```

- ▶ In (pseudo) code, this is as follows:

```
class SelectStatement : Statement {  
    vector<Expression*> select_list = { COUNT(*) }  
    TableRef from_clause = CROSS_PRODUCT("lineitem", "orders");  
    Expression *where_clause = "l_orderkey"="o_orderkey"  
                                AND "o_orderstatus"="X"  
                                AND "l_tax" > 50;  
}
```

```
SELECT COUNT(*)  
FROM lineitem, orders  
WHERE l_orderkey=o_orderkey  
      AND o_orderstatus='X'  
      AND l_tax > 50;
```

- ▶ Note: table/column names *are not resolved yet*
 - ▶ e.g. if **lineitem** table does not exist, no error will be thrown yet

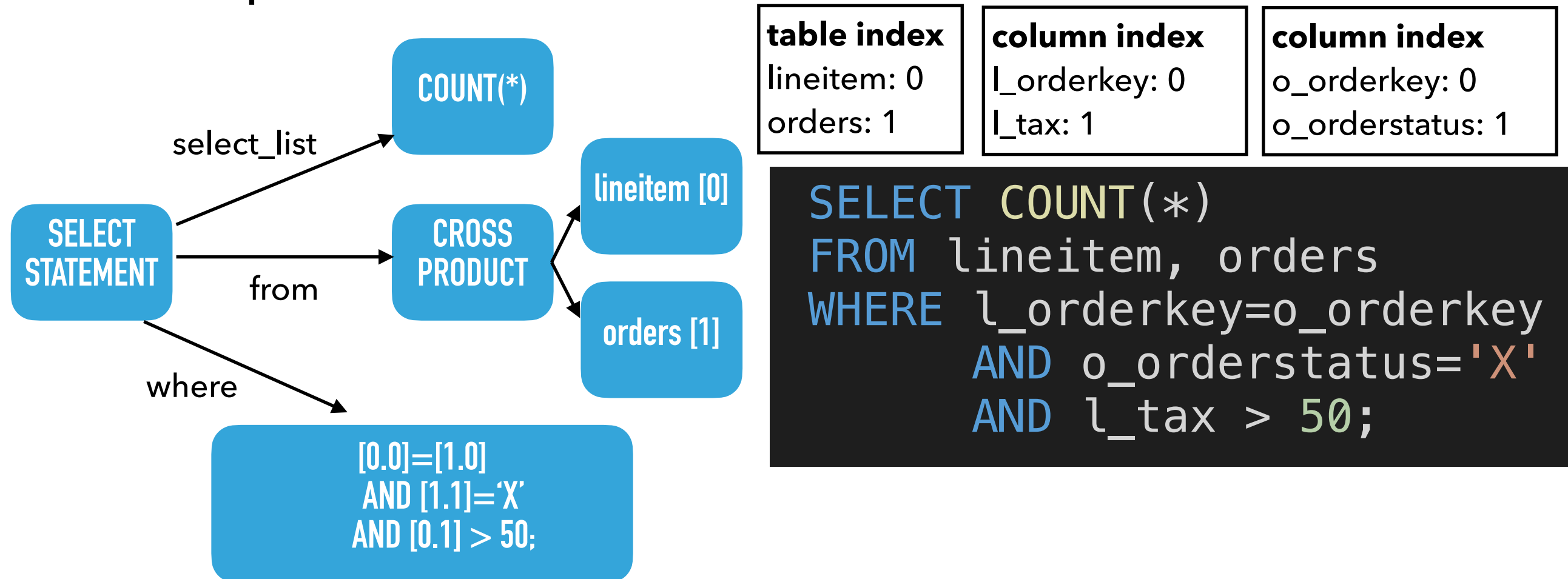


► Binding phase

- Catalog lookup of table/column names
- Type resolution of expressions

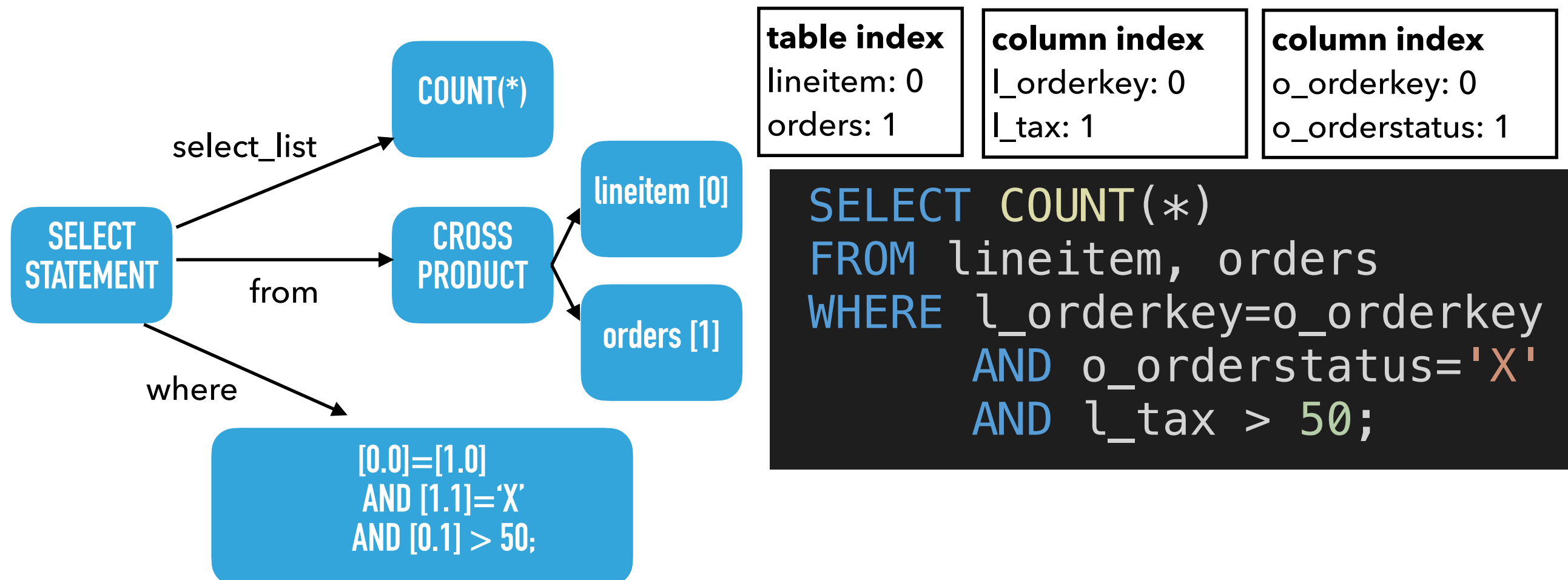
```
SELECT COUNT(*)  
FROM lineitem, orders  
WHERE l_orderkey=o_orderkey  
      AND o_orderstatus='X'  
      AND l_tax > 50;
```

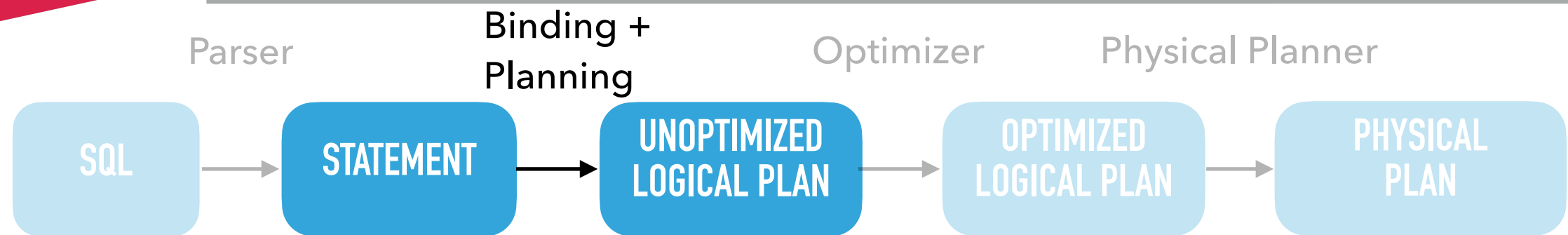
- ▶ Look up tables `lineitem` and `orders` tables
- ▶ Look up **columns** within these tables



- ▶ Replace table names with **table indexes**
- ▶ Replace column names with **table+column indexes**

- ▶ **Type resolution:** look up the types from the tables
 - ▶ `l_orderkey` : INTEGER
 - ▶ `o_orderkey` : INTEGER
 - ▶ `l_orderkey = o_orderkey` : BOOLEAN





- ▶ **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='X'  
      AND l_tax > 50;
```

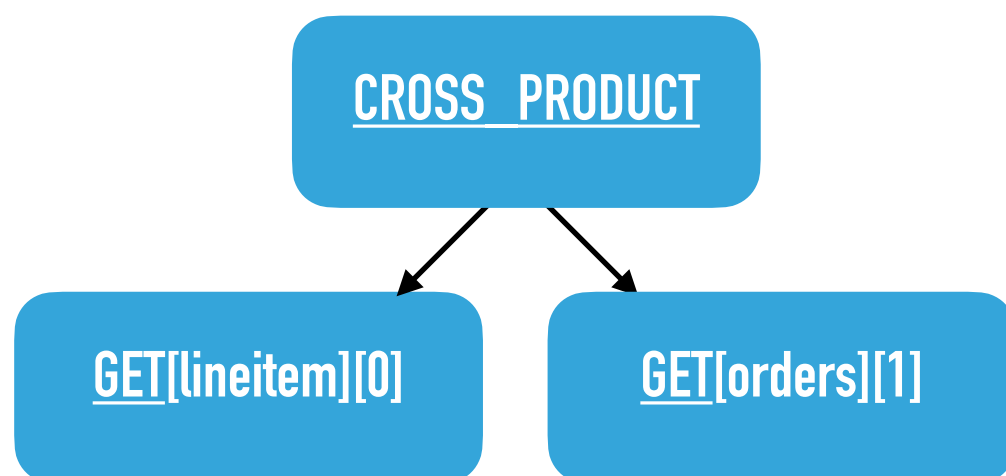

- ▶ Query tree starts with **tables**
- ▶ We have two tables: **lineitem** and **orders**
- ▶ These will result in two **LogicalGet** operations

GET[lineitem][0]

GET[orders][1]

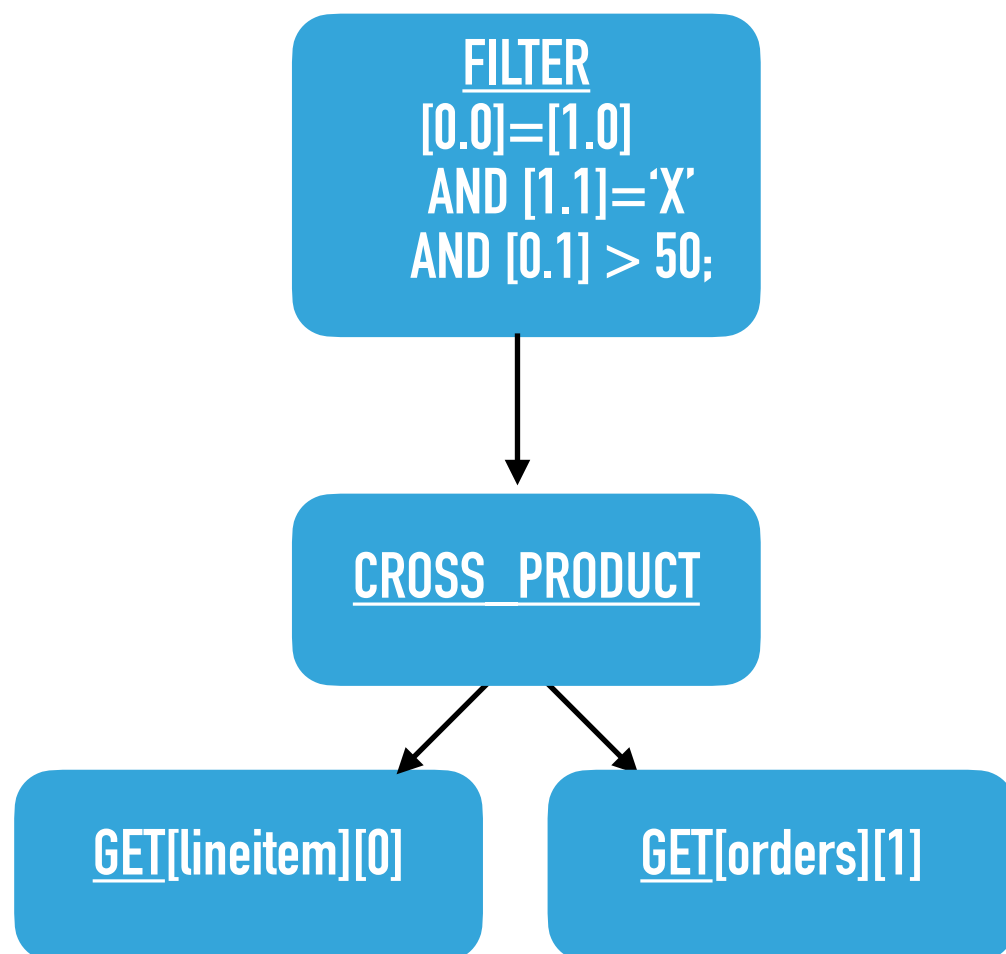
```
SELECT COUNT(*)  
FROM lineitem, orders  
WHERE l_orderkey=o_orderkey  
      AND o_orderstatus='X'  
      AND l_tax > 50;
```

- ▶ The tables are combined with a **cross product**
 - ▶ There is no explicit join here
- ▶ The optimizer will later convert this into a join



```
SELECT COUNT(*)  
FROM lineitem, orders  
WHERE l_orderkey=o_orderkey  
      AND o_orderstatus='X'  
      AND l_tax > 50;
```

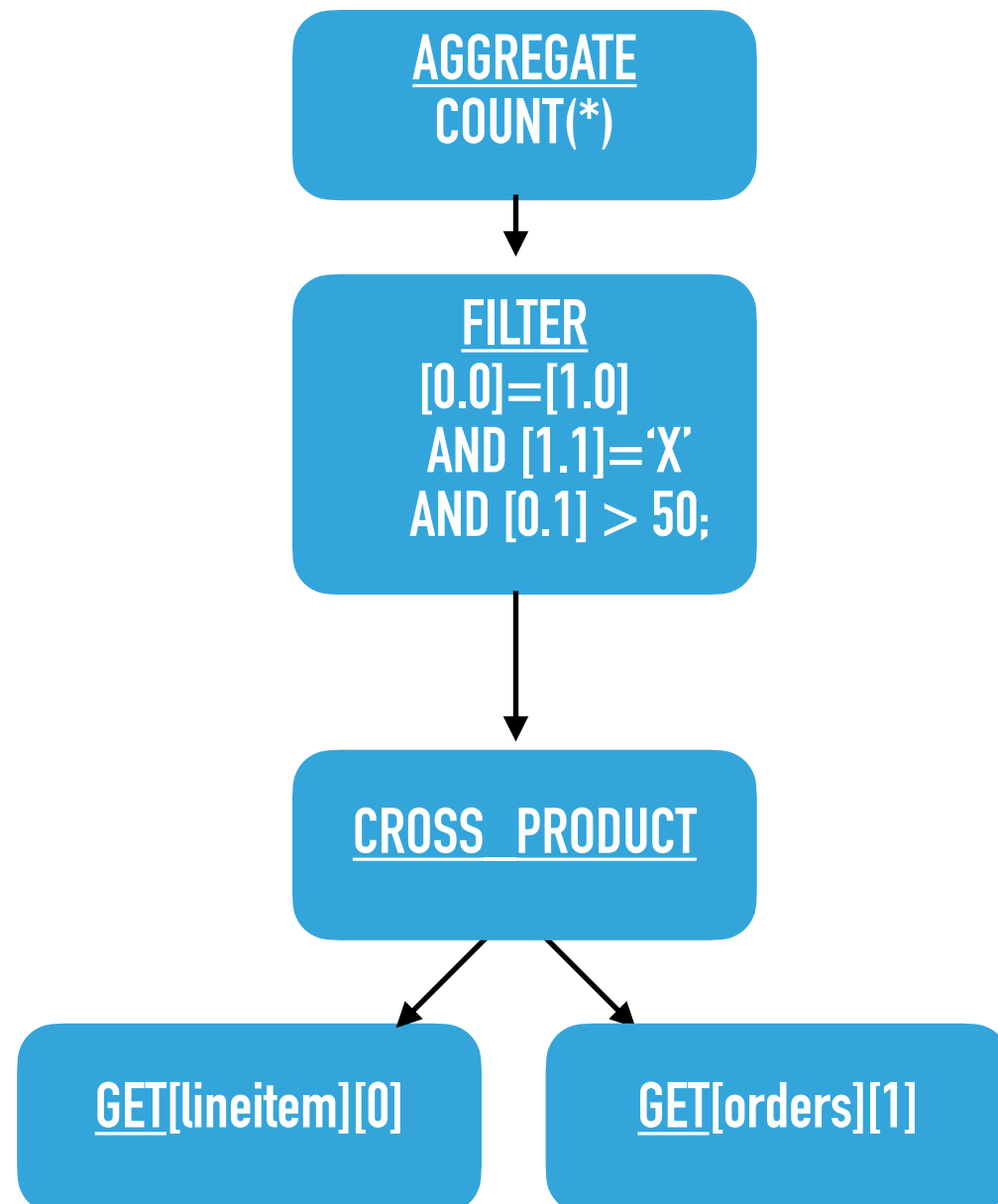
- ▶ After **Filter** is added
 - ▶ Filter has single big expression



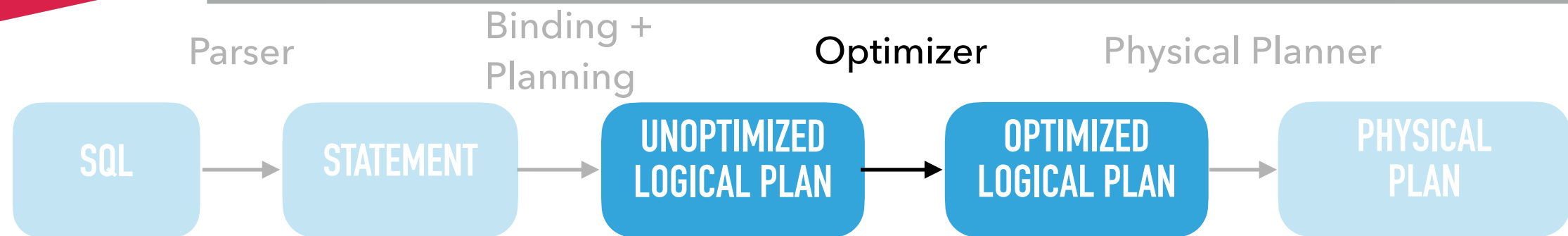
```
SELECT COUNT(*)  
FROM lineitem, orders  
WHERE l_orderkey=o_orderkey  
      AND o_orderstatus='X'  
      AND l_tax > 50;
```

- ▶ Finally add aggregate computation

Logical Query Tree



```
SELECT COUNT(*)
FROM lineitem, orders
WHERE l_orderkey=o_orderkey
AND o_orderstatus='X'
AND l_tax > 50;
```

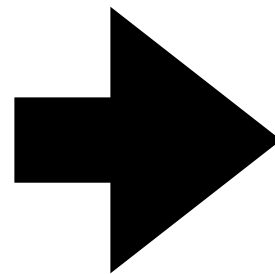
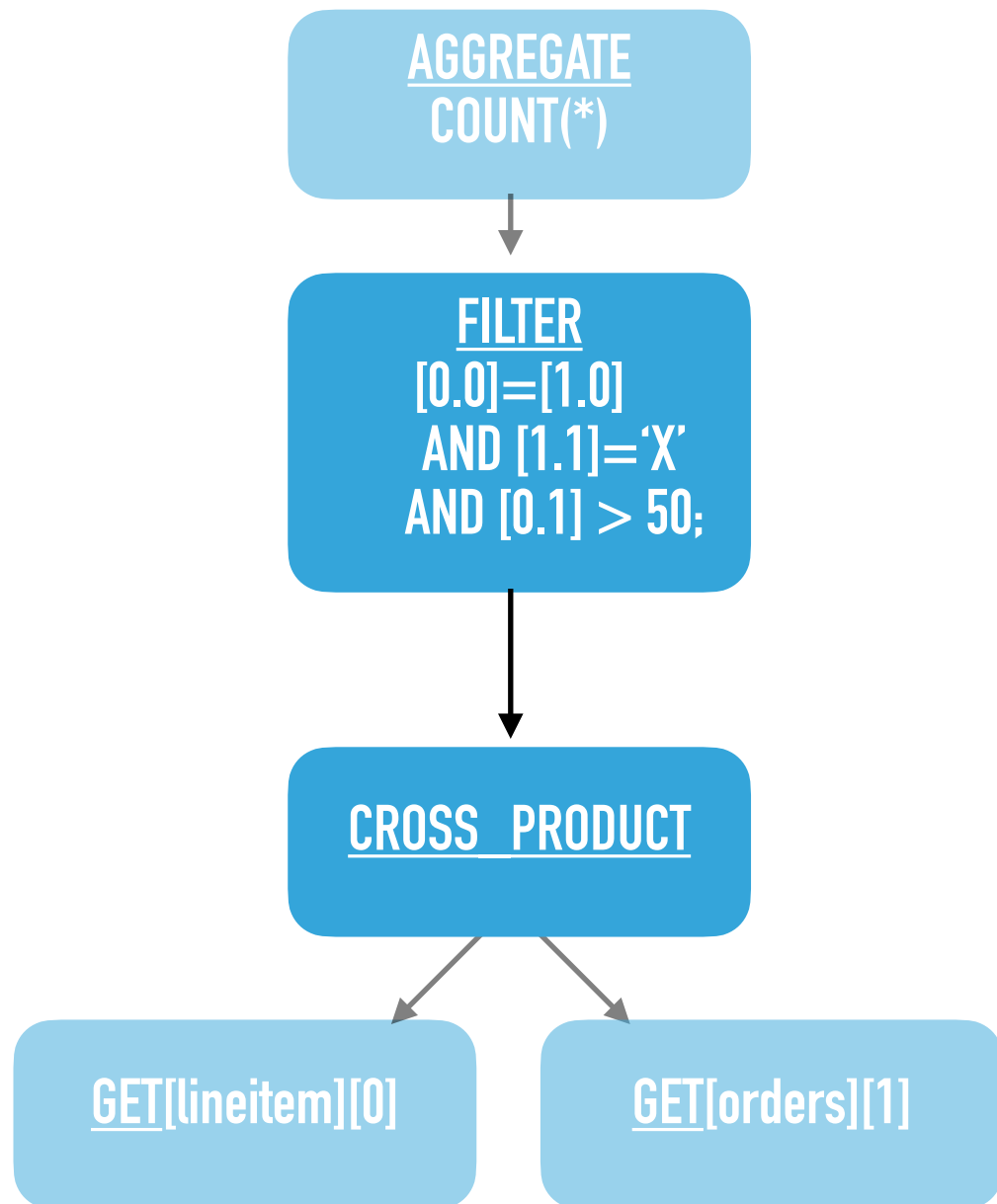


- ▶ **Optimizer:** transforms the logical query tree
- ▶ Created plan is logically equivalent
 - ▶ But (hopefully) faster

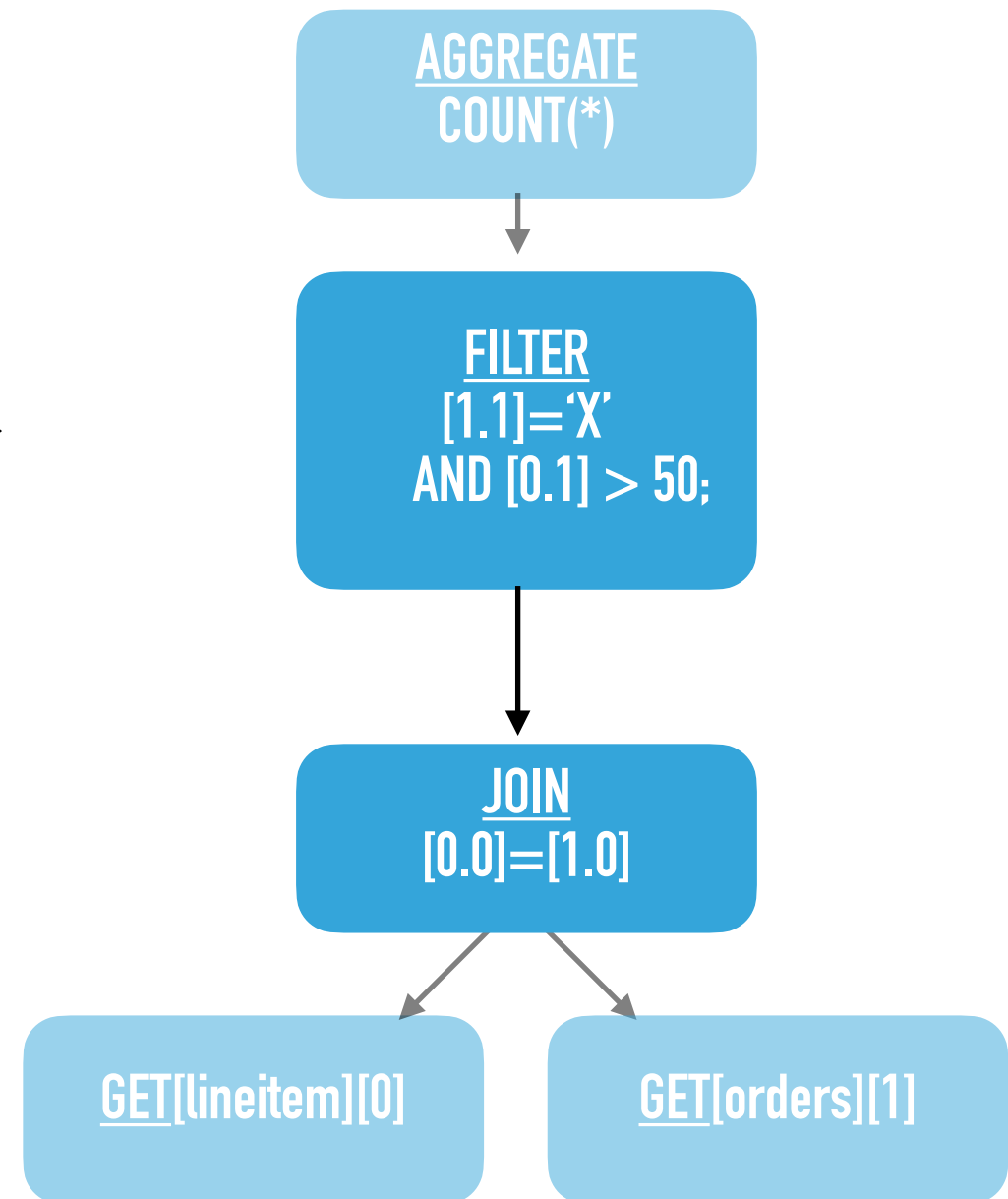
```
SELECT COUNT(*)  
FROM lineitem, orders  
WHERE l_orderkey=o_orderkey  
      AND o_orderstatus='X'  
      AND l_tax > 50;
```

- Pushdown filter into cross product: creates a join

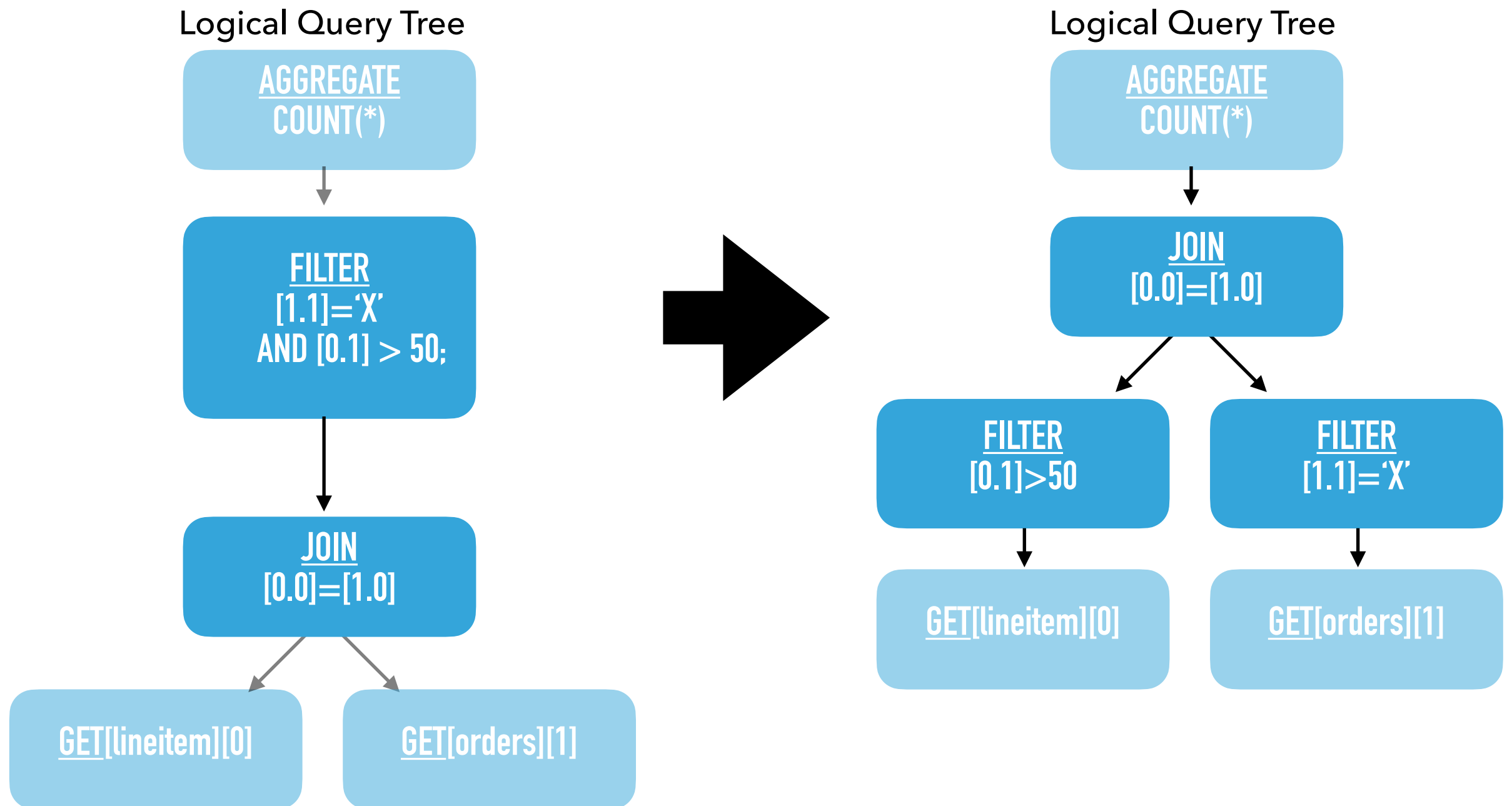
Logical Query Tree



Logical Query Tree



► Pushdown filters below the join



Parser &
Transformer

Binding Phase

Planner

Optimizer

Physical Planner

SQL

STATEMENT

BOUND
STATEMENTUNOPTIMIZED
LOGICAL PLANOPTIMIZED
LOGICAL PLANPHYSICAL
PLAN

- ▶ Many possible optimizations possible here
 - ▶ Join ordering, constant folding, CSE, subquery flattening, common subtree elimination, projection pushdown, etc...
- ▶ For this query only filter pushdown is necessary

```
SELECT COUNT(*)  
FROM lineitem, orders  
WHERE l_orderkey=o_orderkey  
      AND o_orderstatus='X'  
      AND l_tax > 50;
```


Parser &
Transformer

Binding Phase

Planner

Optimizer

Physical Planner

SQL

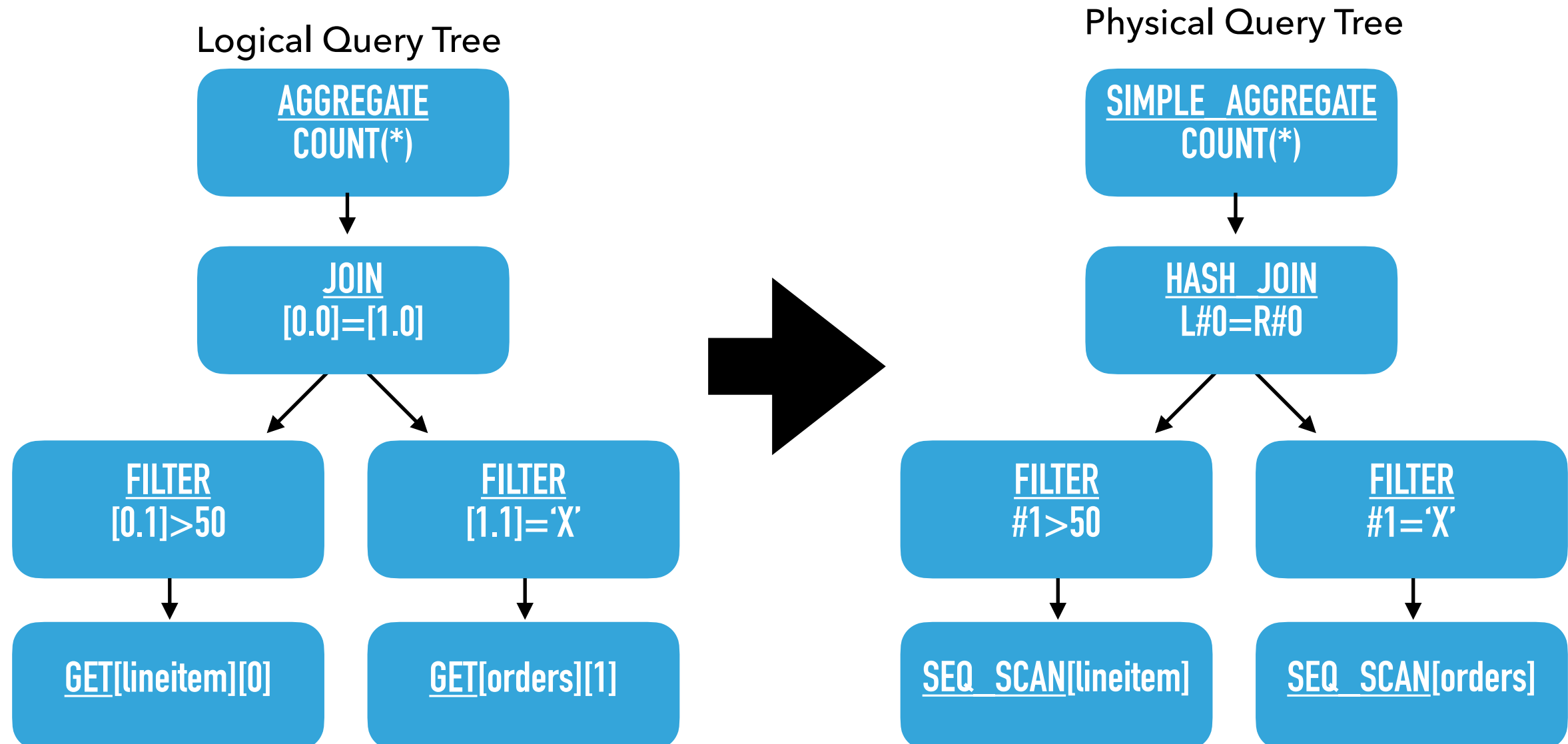
STATEMENT

BOUND
STATEMENTUNOPTIMIZED
LOGICAL PLANOPTIMIZED
LOGICAL PLANPHYSICAL
PLAN

- ▶ **Physical planner:** converts logical plan into physical (executable) plan
- ▶ Makes decision on implementations of operators
 - ▶ e.g. use a HashJoin, MergeJoin or NestedLoopJoin

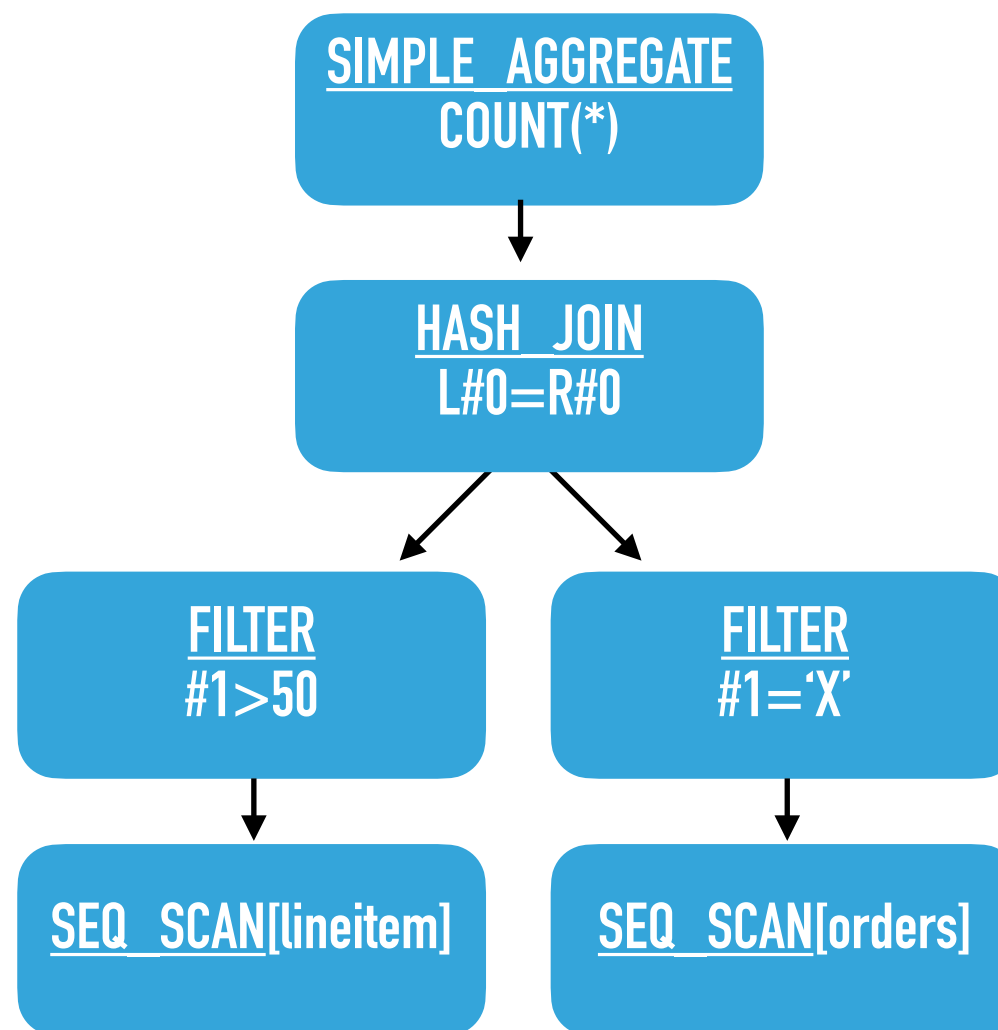
```
SELECT COUNT(*)  
FROM lineitem, orders  
WHERE l_orderkey=o_orderkey  
      AND o_orderstatus='X'  
      AND l_tax > 50;
```

- ▶ **SimpleAggregate:** no groups, no hash required
- ▶ **Hash Join:** Most effective for this equality join
- ▶ **Sequential Scan:** No index that helps us speed up



- ▶ Now we have the final query tree
- ▶ This is what we **execute** to run the query

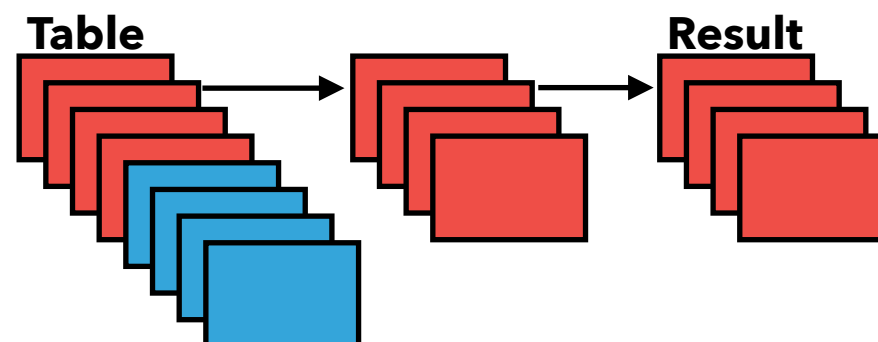
Physical Query Tree



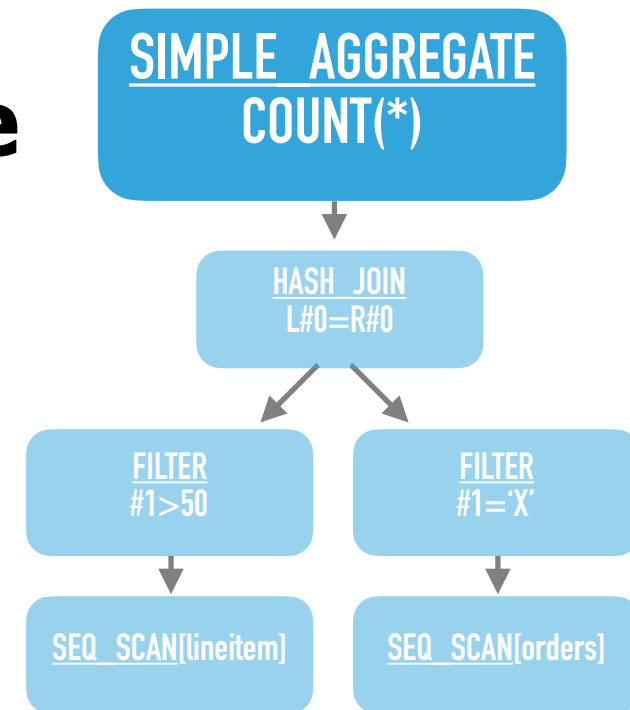
- ▶ Internals at a Glance
- ▶ Query processing pipeline
- ▶ **Query execution**
- ▶ Hands-On

- ▶ DuckDB uses a vectorized pull-based model
 - ▶ "vector volcano"
- ▶ Query starts by calling **GetChunk** on the root node
- ▶ Root node recursively calls **GetChunk** on children
- ▶ Scans fetch data from the base tables

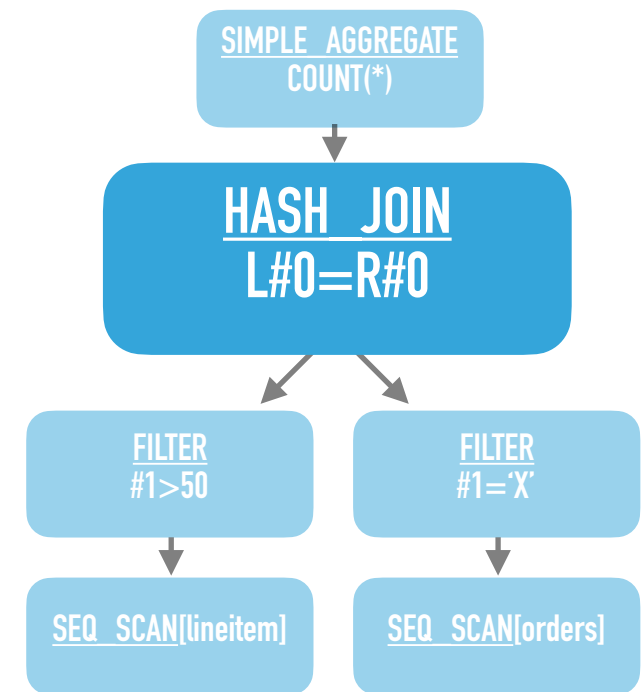
Vectorized Processing



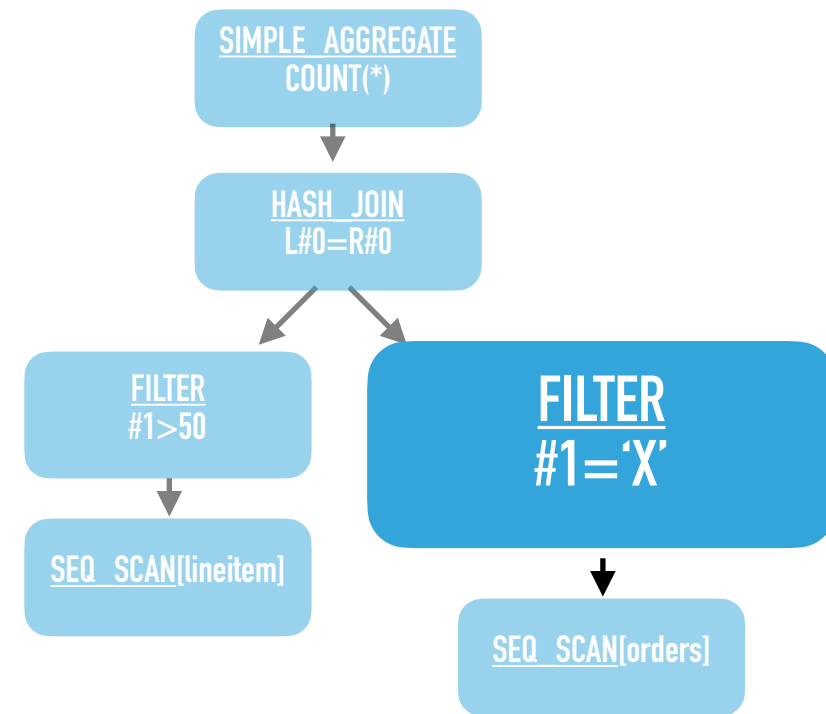
- ▶ Start with root node: **SimpleAggregate**
 - ▶ Aggregate without groups
- ▶ Immediately calls **GetChunk** on child



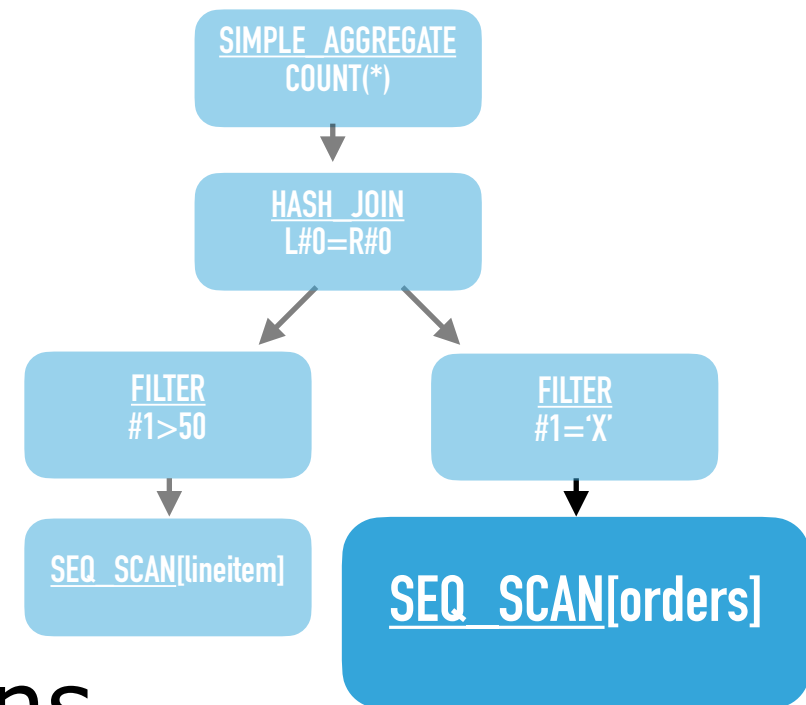
- ▶ **Hash Join**
- ▶ Start by **building HT**
- ▶ Call **GetChunk** on right node



- ▶ **Filter**
- ▶ Again, pull a chunk from child



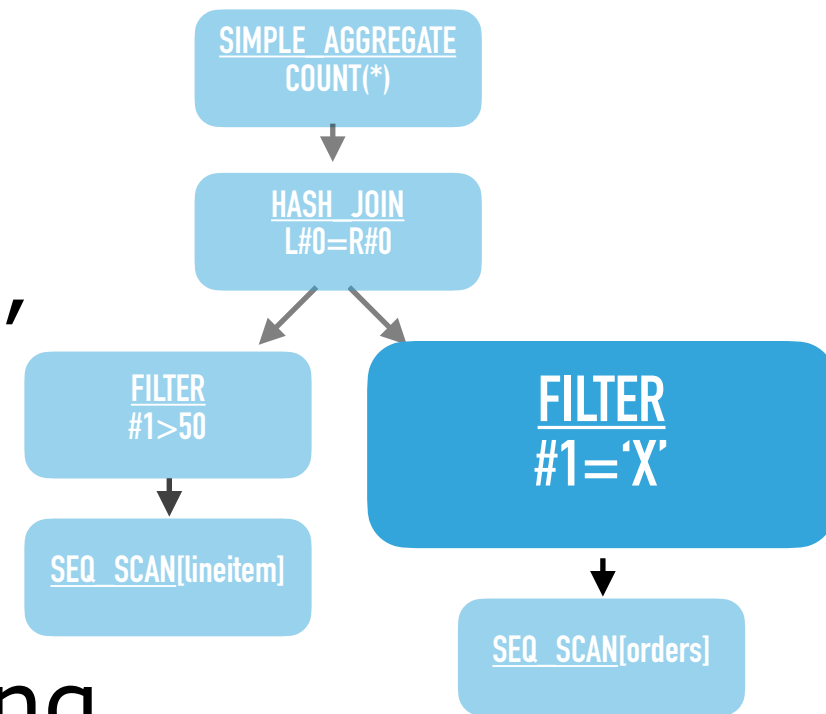
- ▶ **Sequential Scan**
- ▶ Finally we can start executing
- ▶ Scan the base table
- ▶ Return a DataChunk with two columns
 - ▶ `o_orderkey` and `o_orderstatus`



DataChunk			
	V1		V2
INTEGER	1	VARCHAR	N
	2		X
	3		N

► Filter

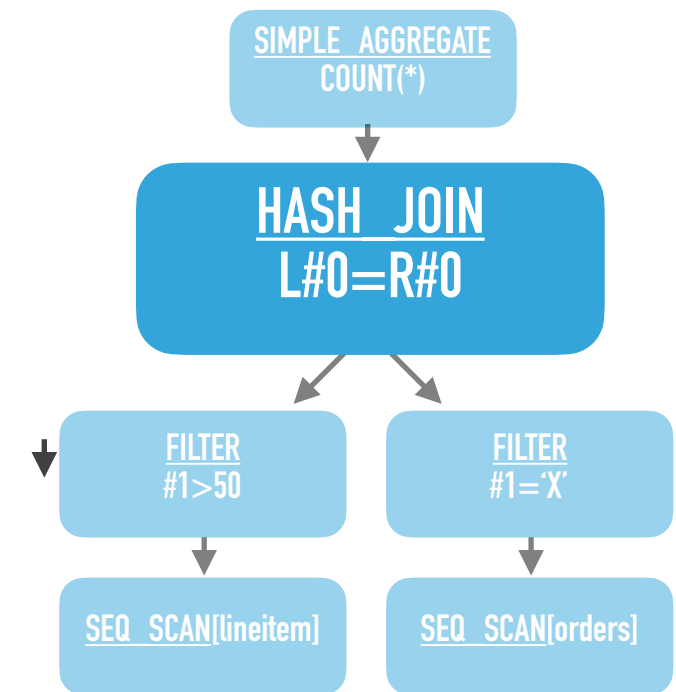
- Now we can perform the filter $\#1 = 'X'$
 - Only the second tuple passes
- **Selection vector** pointing to surviving tuple is created
- Note that no data is copied or changed



DataChunk				
	V1		V2	SEL
INTEGER	1		N	1
	2		X	
	3		N	

► Hash Join

- Now we have our first input chunk
- We input it into the HT
- Now we fetch another chunk from RHS



DataChunk				
	V1		V2	SEL
INTEGER	1		N	1
	2		X	
	3		N	
VARCHAR				

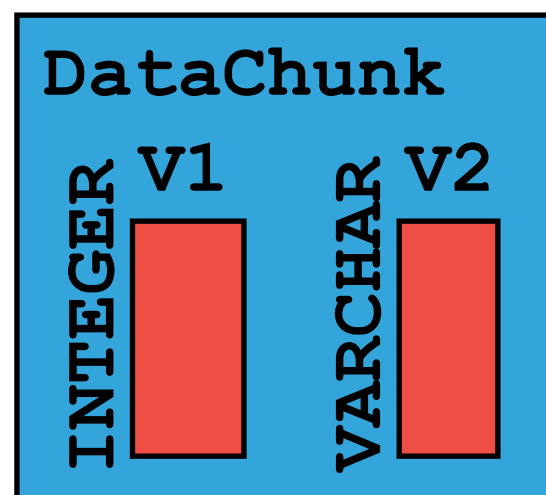
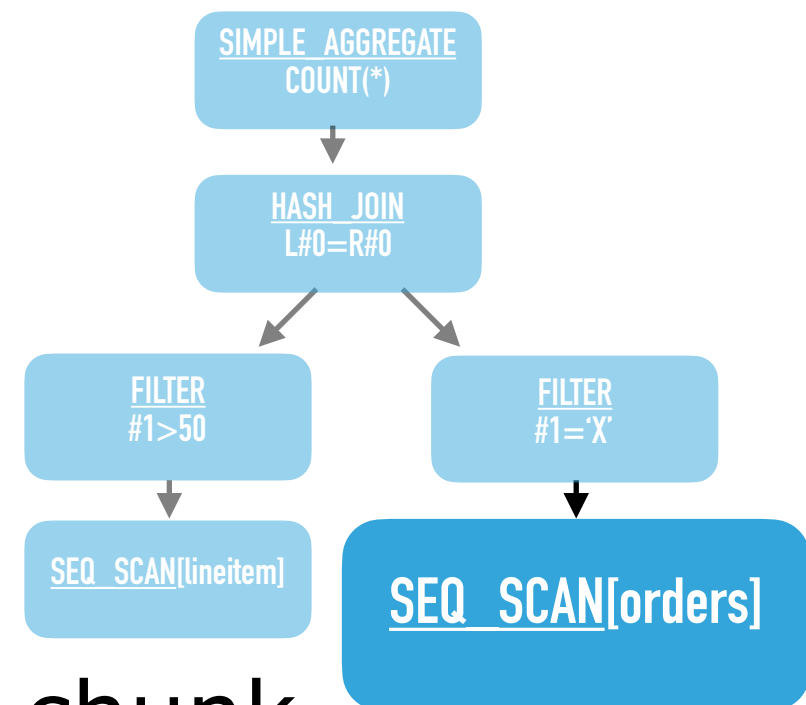
Keys	Payload	Next
2	X	0

► Sequential Scan

► The filter again calls `GetChunk`

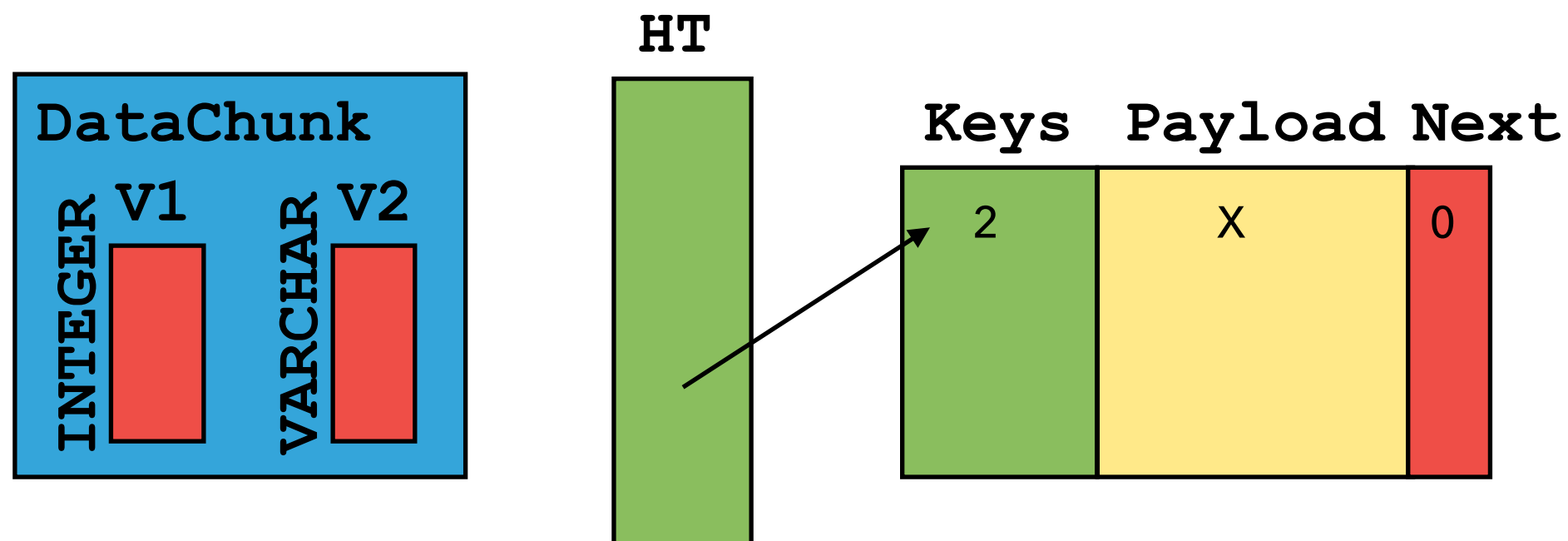
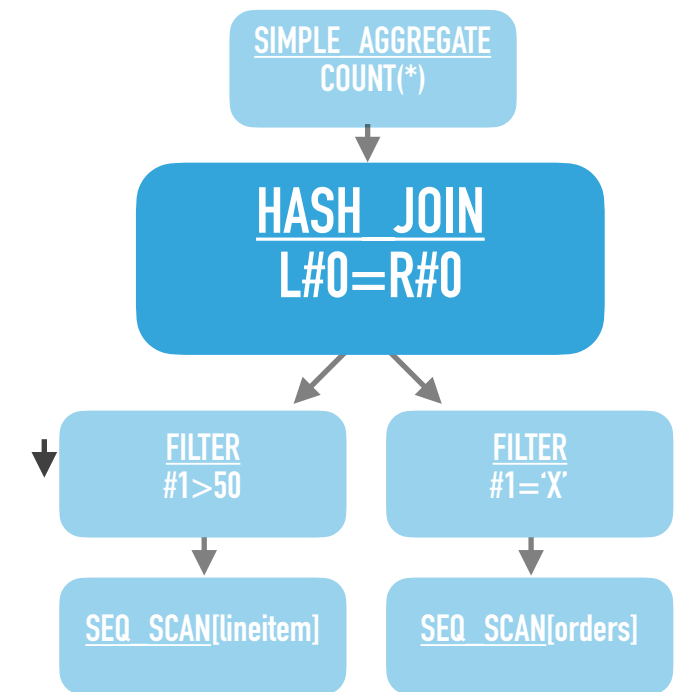
► Scan base table again:

► The scan is finished, return empty chunk

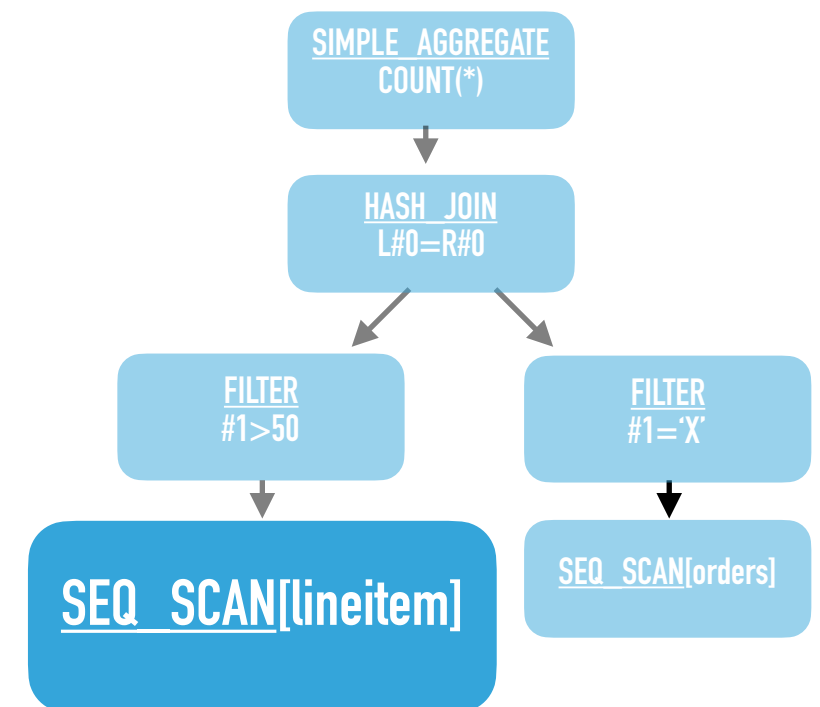


► Hash Join

- HT receives second input chunk
 - But it is empty!
- The RHS is exhausted
- Finish building HT and call **GetChunk** on LHS

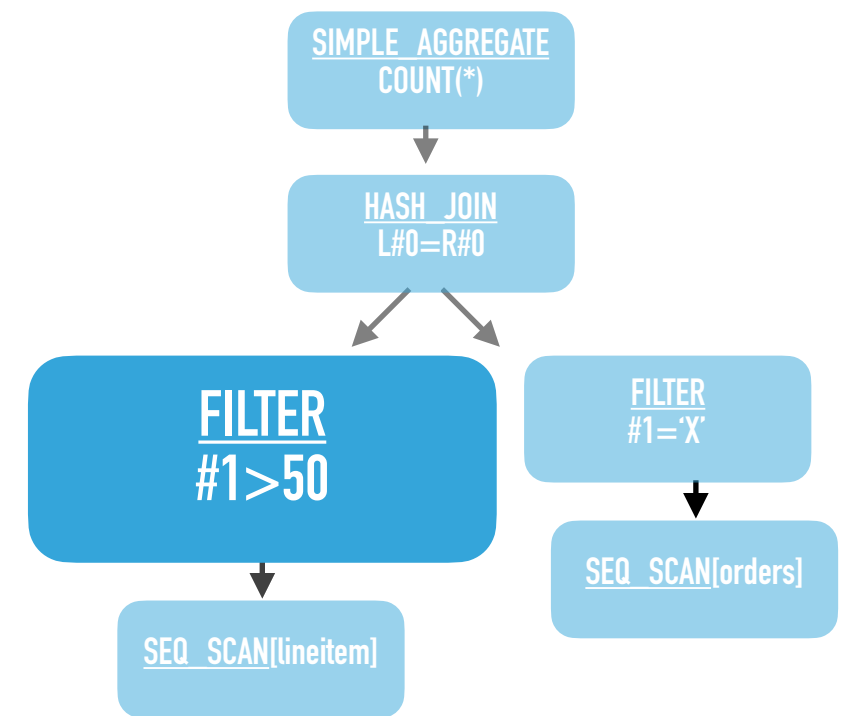


- ▶ **Sequential Scan**
- ▶ We arrive at scan on lineitem
- ▶ DataChunk with two columns
 - ▶ `l_orderkey` and `l_tax`



DataChunk			
	V1		V2
INTEGER	1	VARCHAR	80
	2		60
	2		20

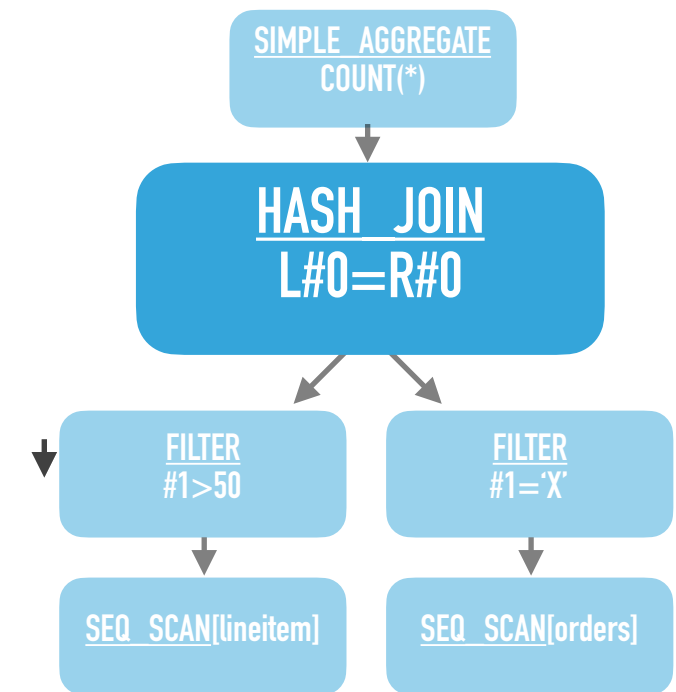
- ▶ **Filter**
- ▶ Performs the filter **#1>50**
- ▶ Again, add a selection vector



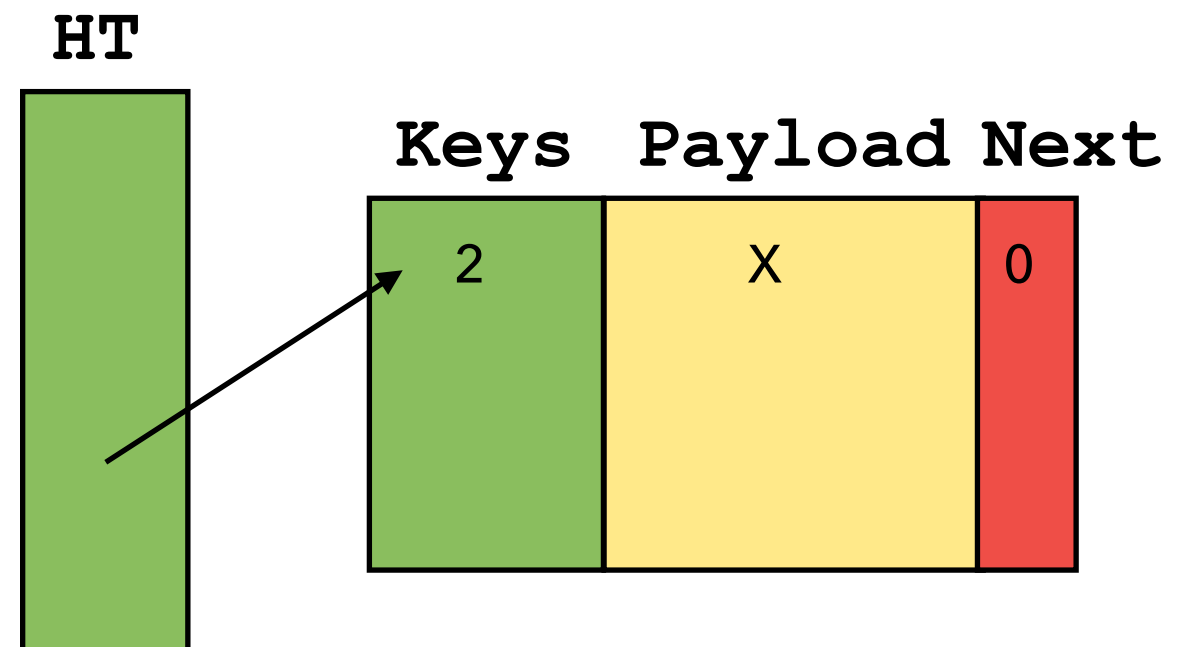
DataChunk				
	V1		V2	SEL
INTEGER	1	VARCHAR	80	0
	2		60	1
	2		20	

► Hash Join

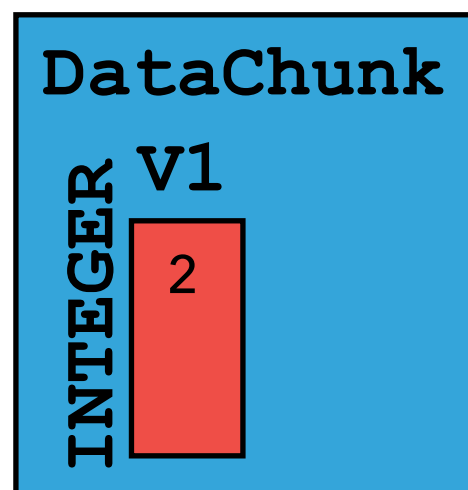
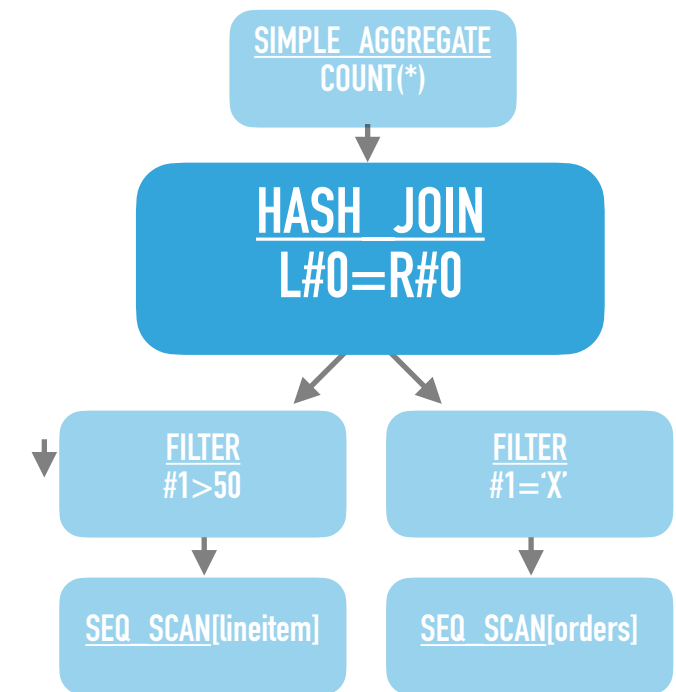
- Now it is time to probe the HT
- We compute the hash for each tuple
- Then lookup in the HT



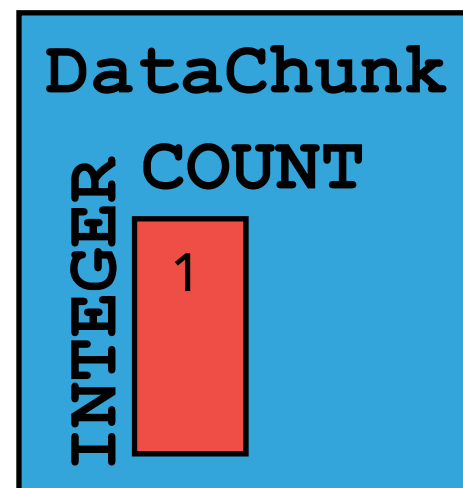
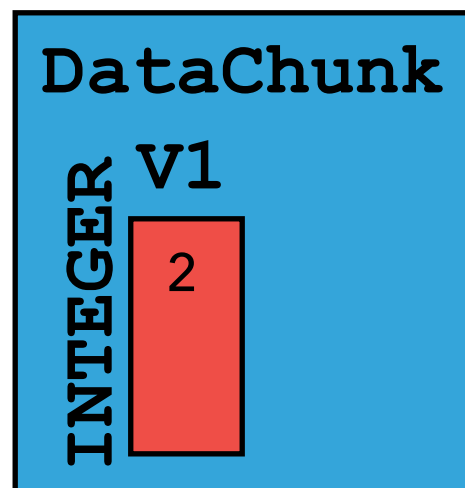
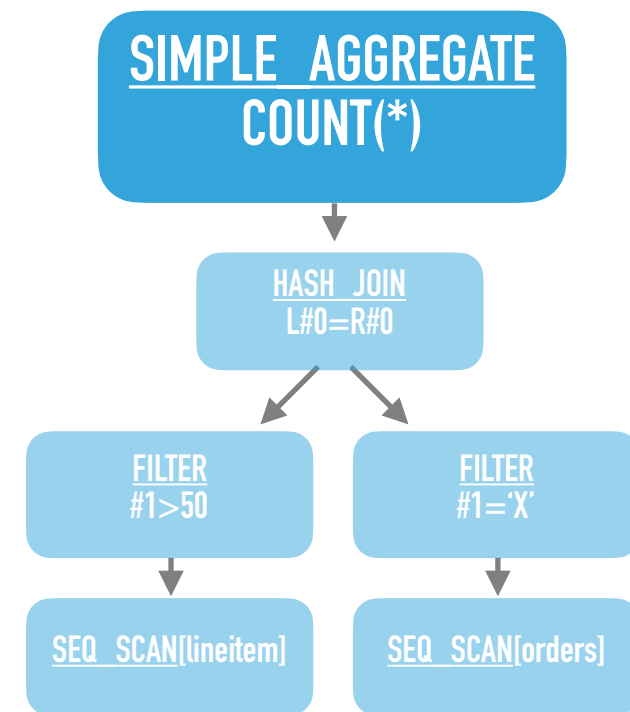
DataChunk				
INTEGER	V1	VARCHAR	V2	SEL
	1		80	0
	2		60	1
	2		20	



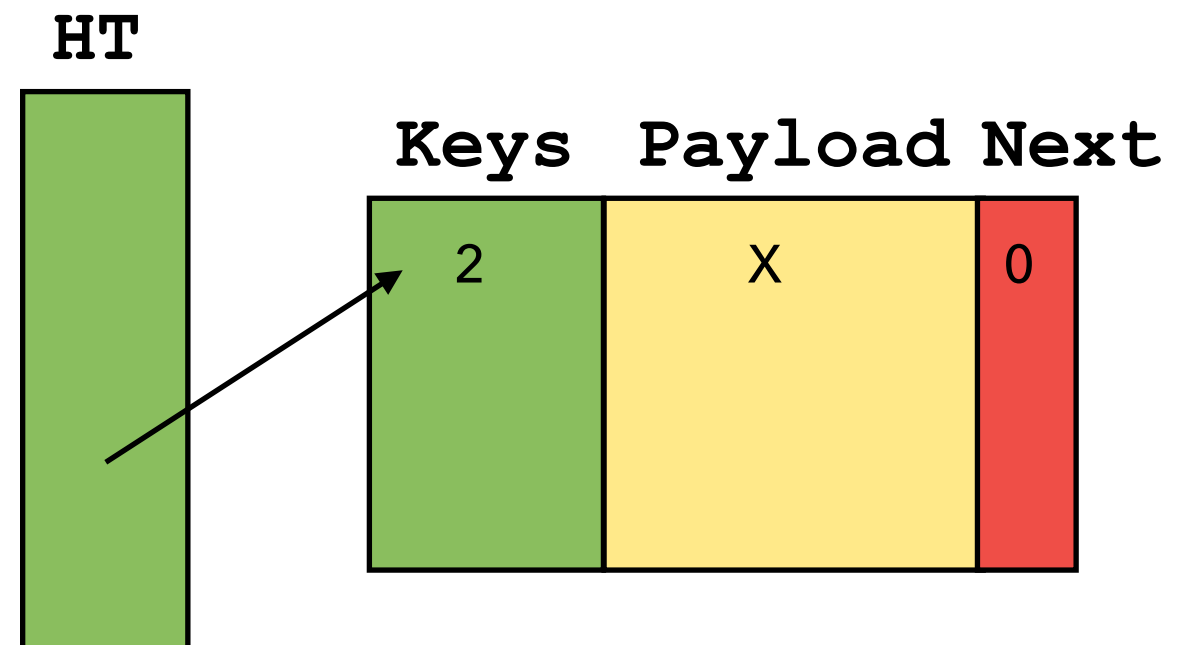
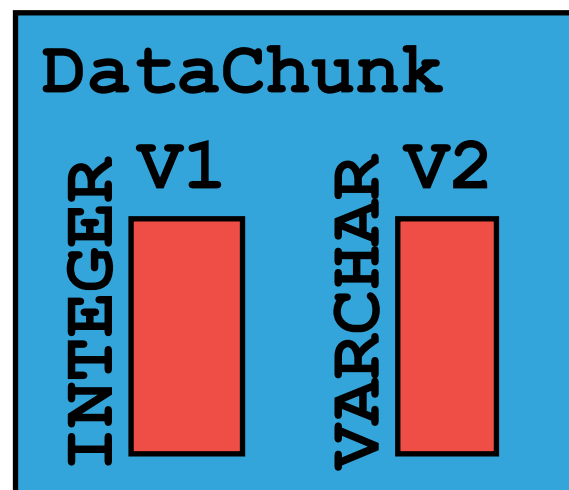
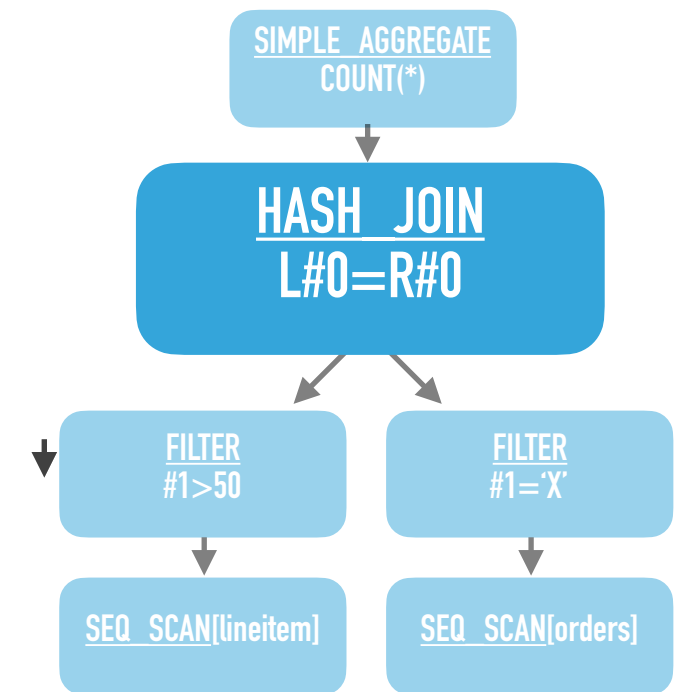
- ▶ **Hash Join**
- ▶ We get one hit on our join!
- ▶ The hash join now produces the result
- ▶ We return this to the aggregate



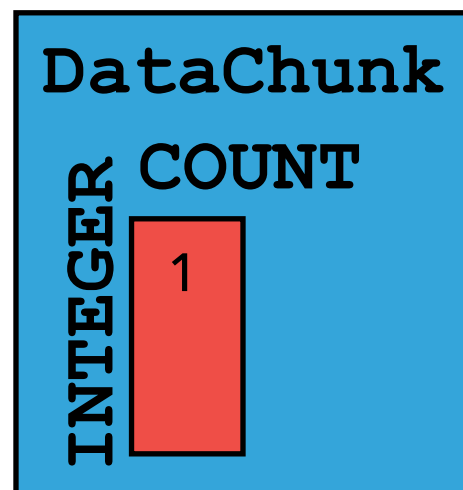
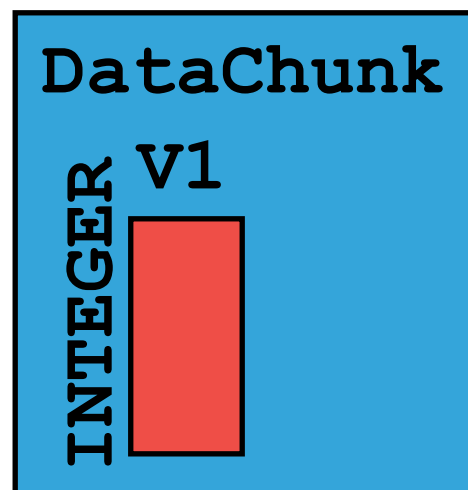
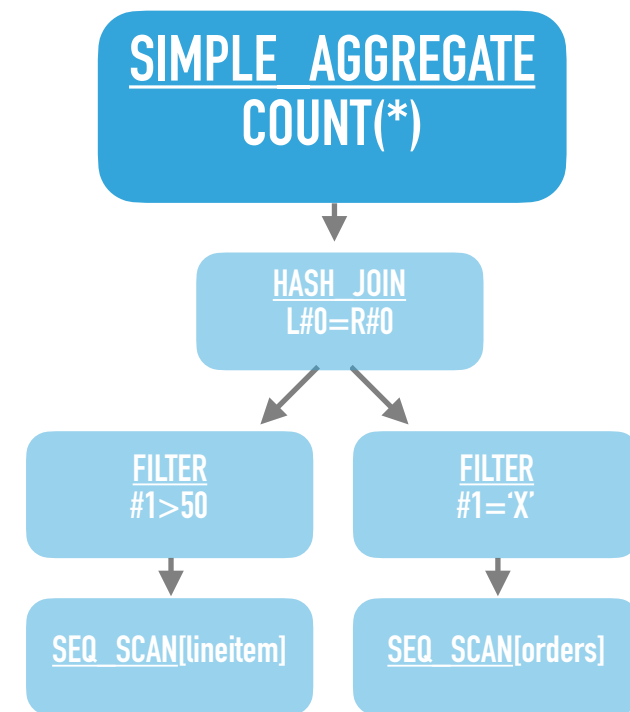
- ▶ The aggregate takes our input chunk
- ▶ Updates the aggregate
- ▶ Then fetches from the child again



- ▶ We go back to the hash join
- ▶ Fetch from probe side again
- ▶ This time, input chunk is empty
- ▶ Now the hash join is entirely finished!



- ▶ Aggregate gets an empty chunk
- ▶ Returns the final result of our query



- ▶ 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!