

Mark Raasveldt & Pedro Holanda

DuckDB

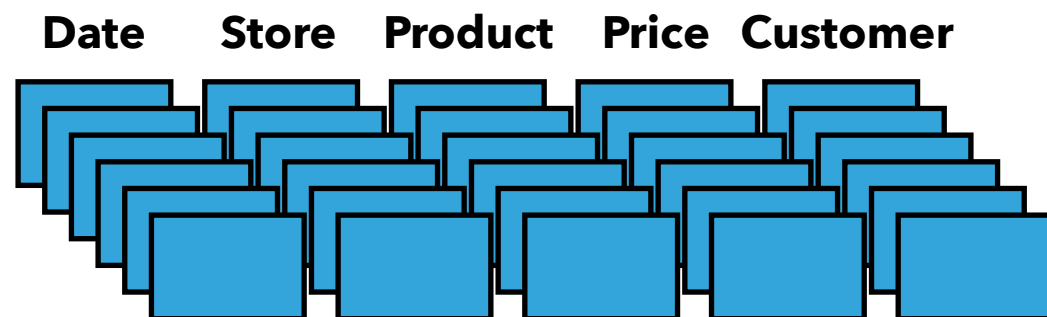
an Embeddable Analytical RDBMS

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

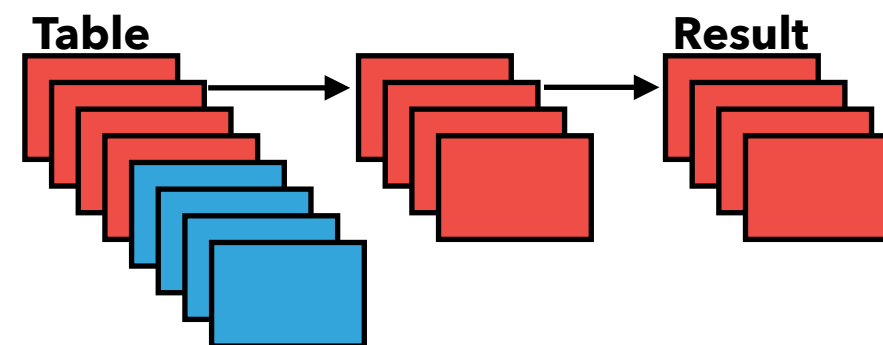
- ▶ **Before we start:** any Windows users here?
- ▶ If you do not have Visual Studio
Download the community edition:
- ▶ <https://visualstudio.microsoft.com/vs/community/>
- ▶ And CMake:
- ▶ <https://cmake.org/download/>
- ▶ This is needed for the practicum!

Internals at a Glance

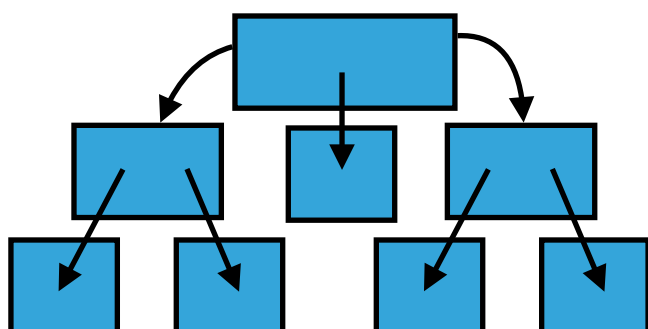
Column-Store



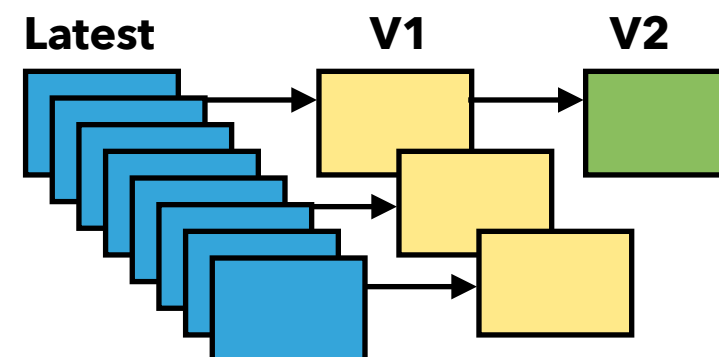
Vectorized Processing



ART Index



Multi-Version Concurrency Control



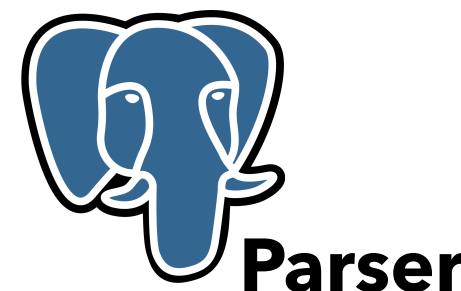
Single-File Storage



4KB

256KB

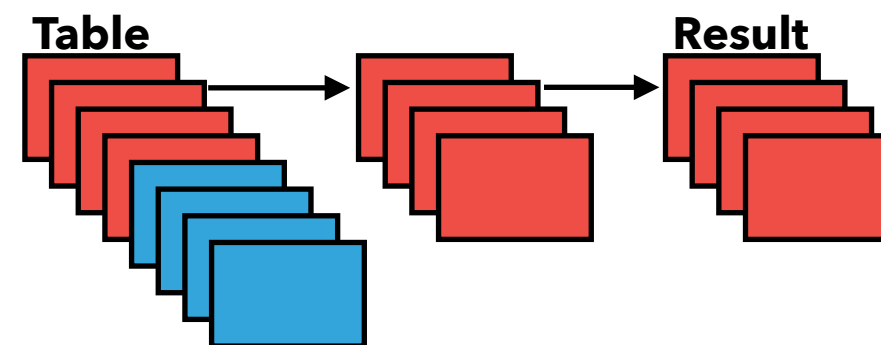
database.db



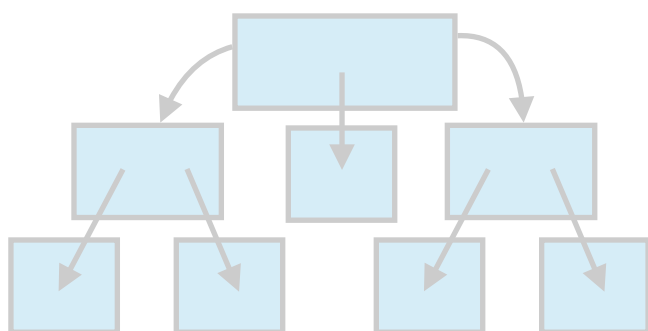
Column-Store



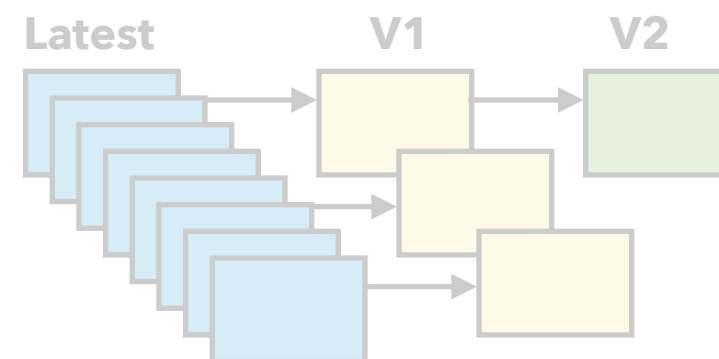
Vectorized Processing



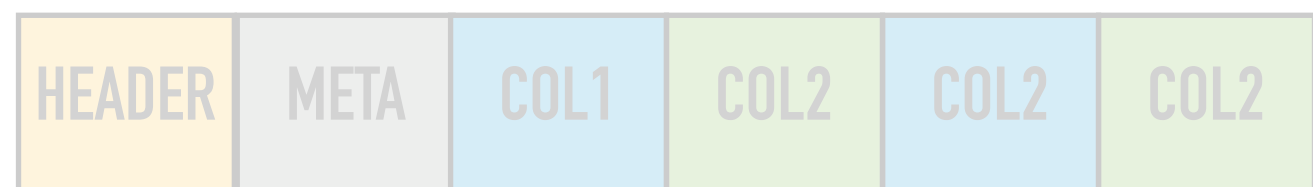
ART Index



Multi-Version Concurrency Control



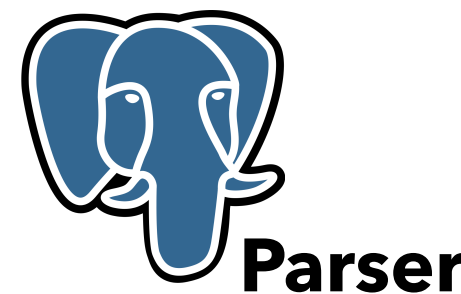
Single-File Storage



4KB

256KB

database.db



Query Processing Pipeline

► Life of a query

► How does the system go from query to result?

► We will focus on the following query:

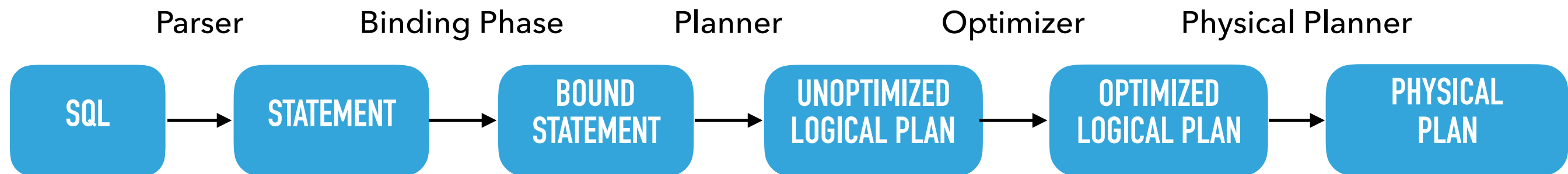
```
SELECT COUNT(*)  
FROM lineitem, orders  
WHERE l_orderkey=o_orderkey  
      AND o_orderstatus='X'  
      AND l_tax > 50;
```

► **Aggregate:** COUNT (*)

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

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

- DuckDB uses a typical pipeline for query processing



Parser

Binding Phase

Planner

Optimizer

Physical Planner

SQL

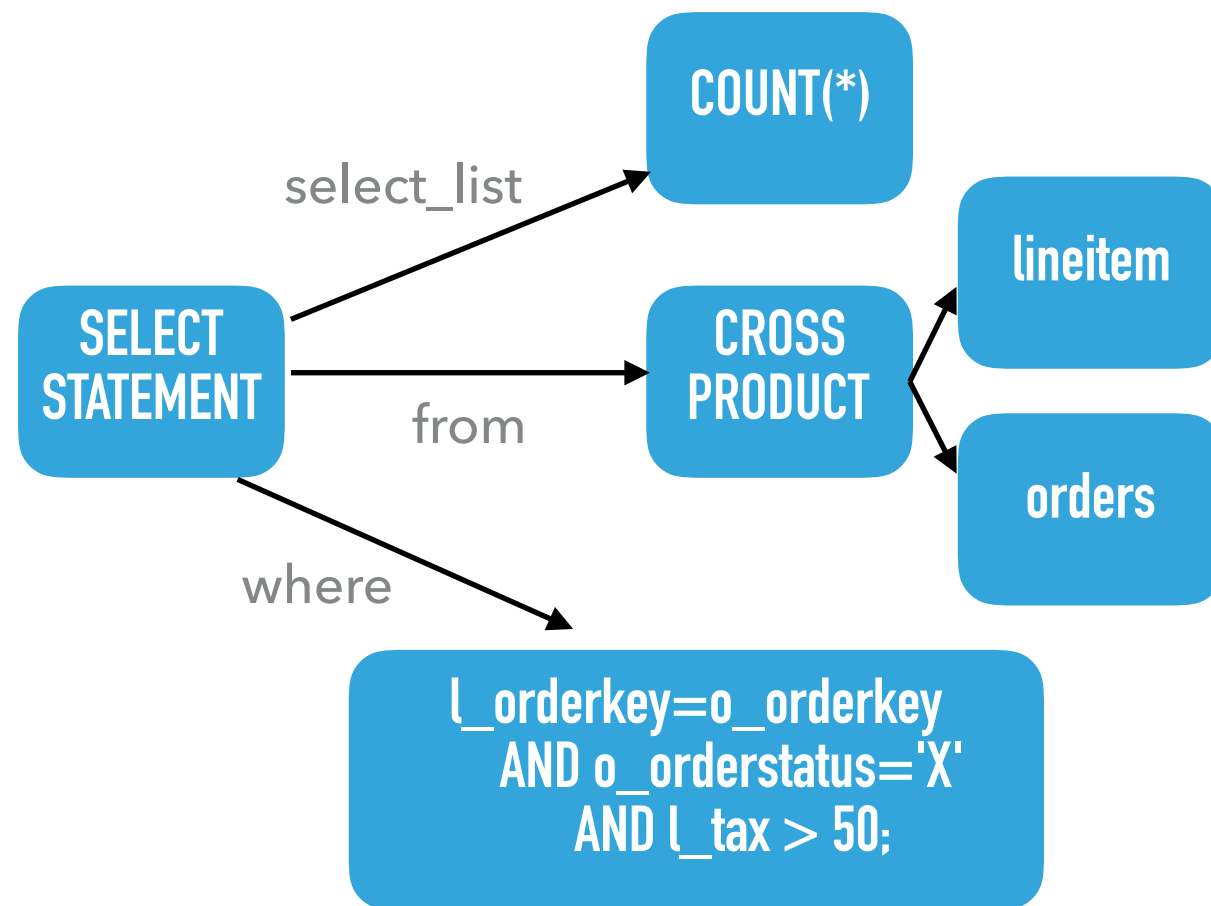
STATEMENT

BOUND
STATEMENTUNOPTIMIZED
LOGICAL PLANOPTIMIZED
LOGICAL PLANPHYSICAL
PLAN

- ▶ 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

Parser &
Transformer

Binding Phase

Planner

Optimizer

Physical Planner

SQL

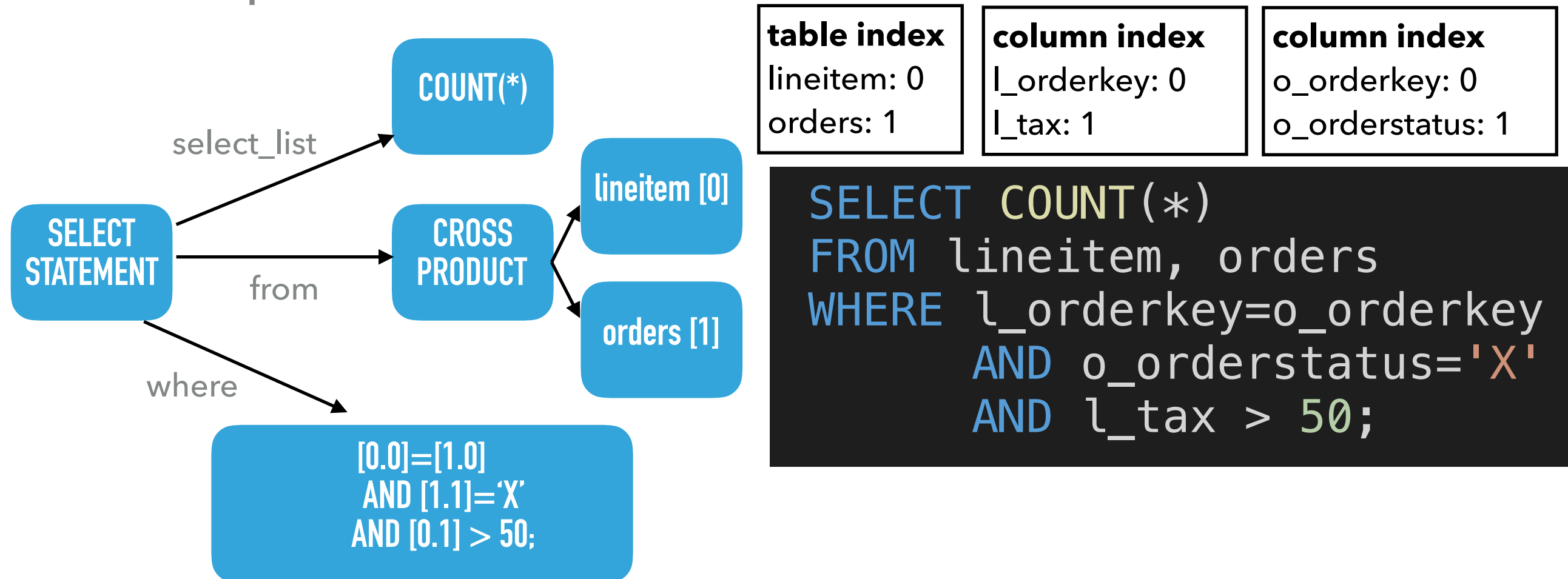
STATEMENT

BOUND
STATEMENTUNOPTIMIZED
LOGICAL PLANOPTIMIZED
LOGICAL PLANPHYSICAL
PLAN

- ▶ **Binding phase** has two purposes
 - ▶ Catalog lookup of table/column names
 - ▶ Type resolution of expressions
- ▶ Replace existing statement by **BoundStatement**

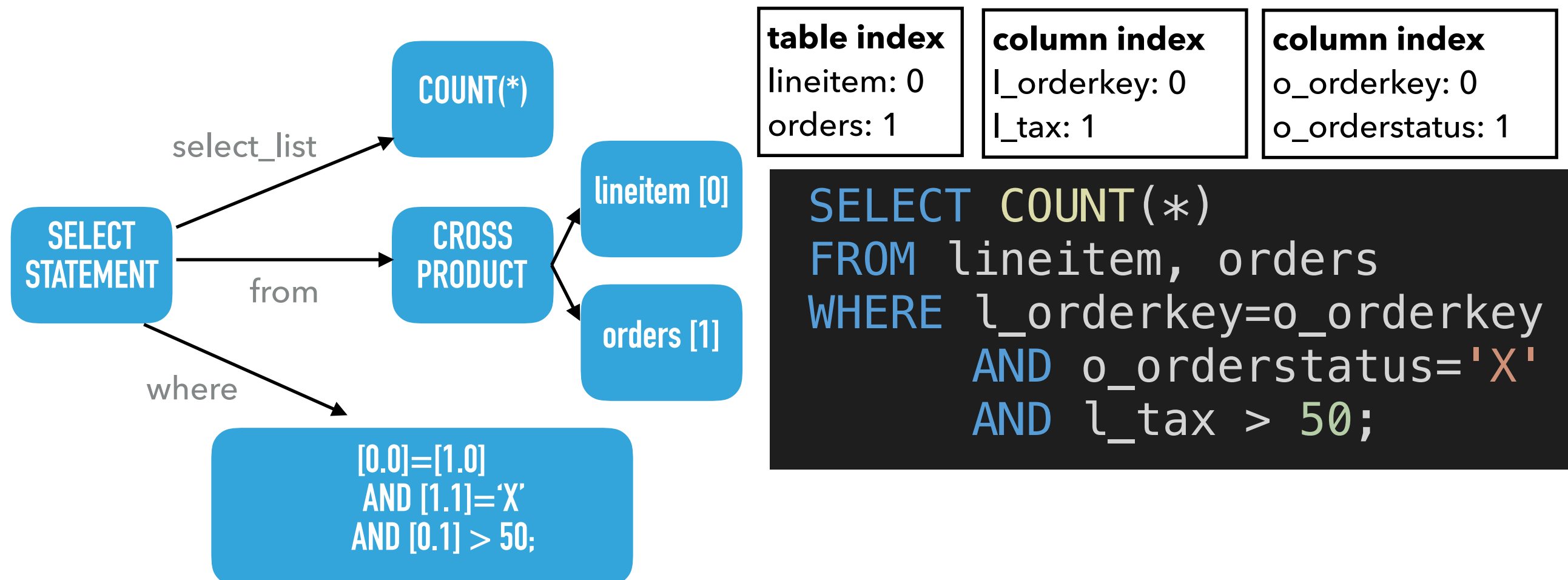
```
SELECT COUNT(*)  
FROM lineitem, orders  
WHERE l_orderkey=o_orderkey  
      AND o_orderstatus='X'  
      AND l_tax > 50;
```

- ▶ 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



Parser &
Transformer

Binding Phase

Planner

Optimizer

Physical Planner

SQL

STATEMENT

BOUND
STATEMENTUNOPTIMIZED
LOGICAL PLANOPTIMIZED
LOGICAL PLANPHYSICAL
PLAN

- ▶ **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;
```

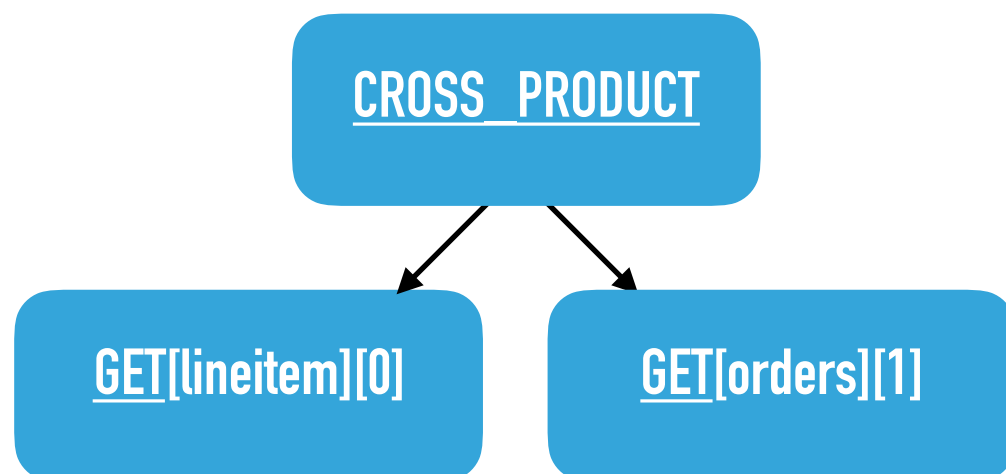

- ▶ Query tree starts with **tables**
- ▶ We have two tables: **lineitem** and **orders**
- ▶ These will result in two **LogicalGet** operations

GET[lineitem][0]

GET[orders][1]

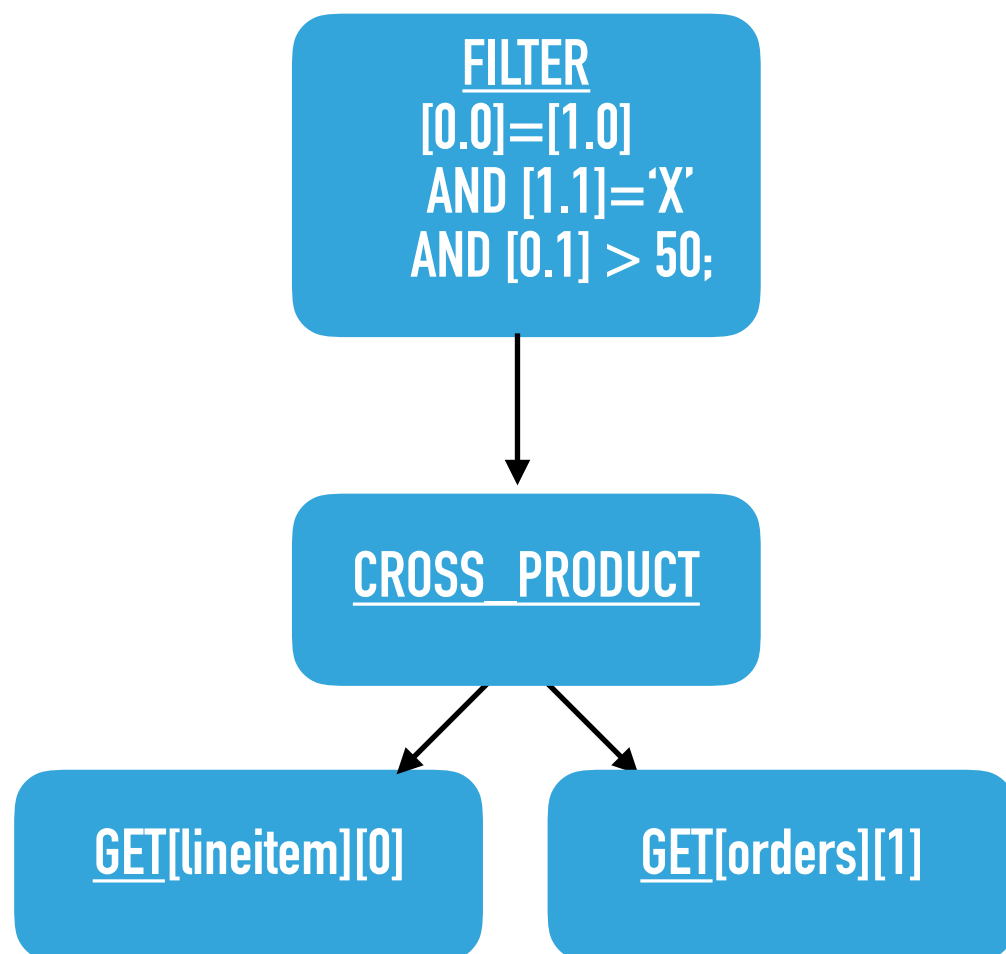
```
SELECT COUNT(*)  
FROM lineitem, orders  
WHERE l_orderkey=o_orderkey  
      AND o_orderstatus='X'  
      AND l_tax > 50;
```

- ▶ 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;
```

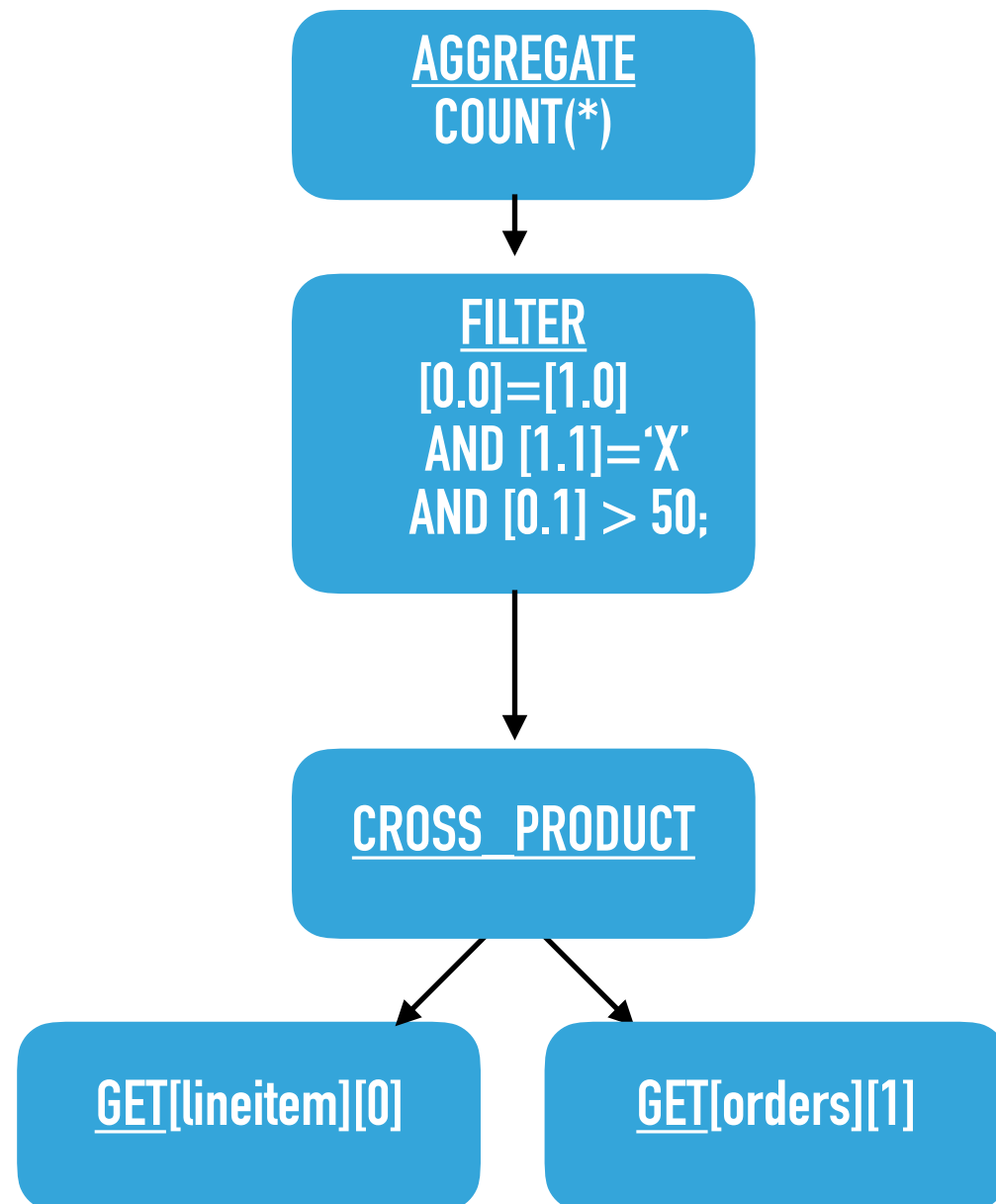
- ▶ After **Filter** is added
- ▶ Filter has single big expression



```
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;
```

Parser &
Transformer

Binding Phase

Planner

Optimizer

Physical Planner

SQL

STATEMENT

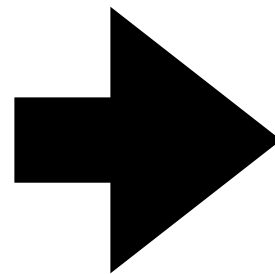
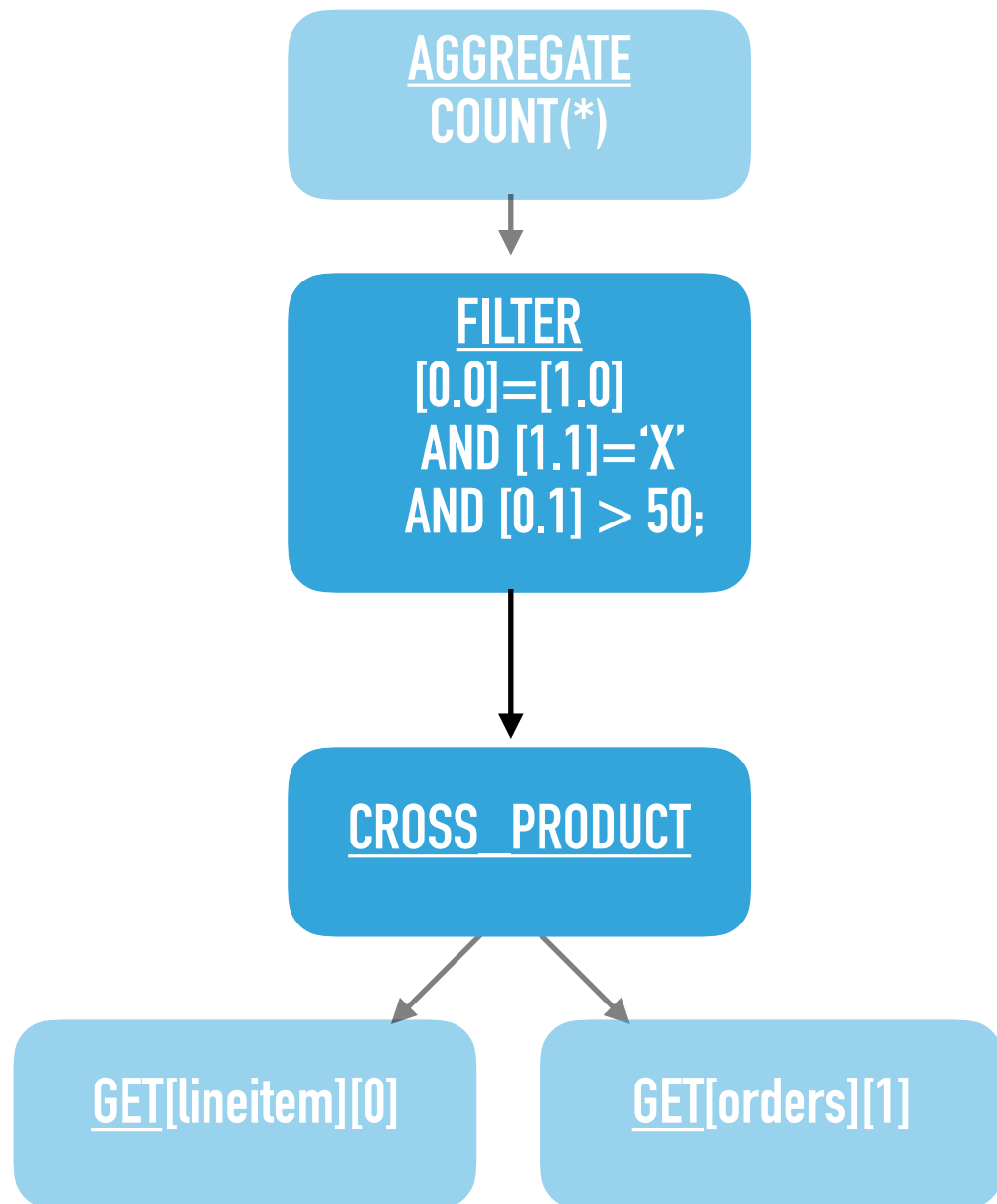
BOUND
STATEMENTUNOPTIMIZED
LOGICAL PLANOPTIMIZED
LOGICAL PLANPHYSICAL
PLAN

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

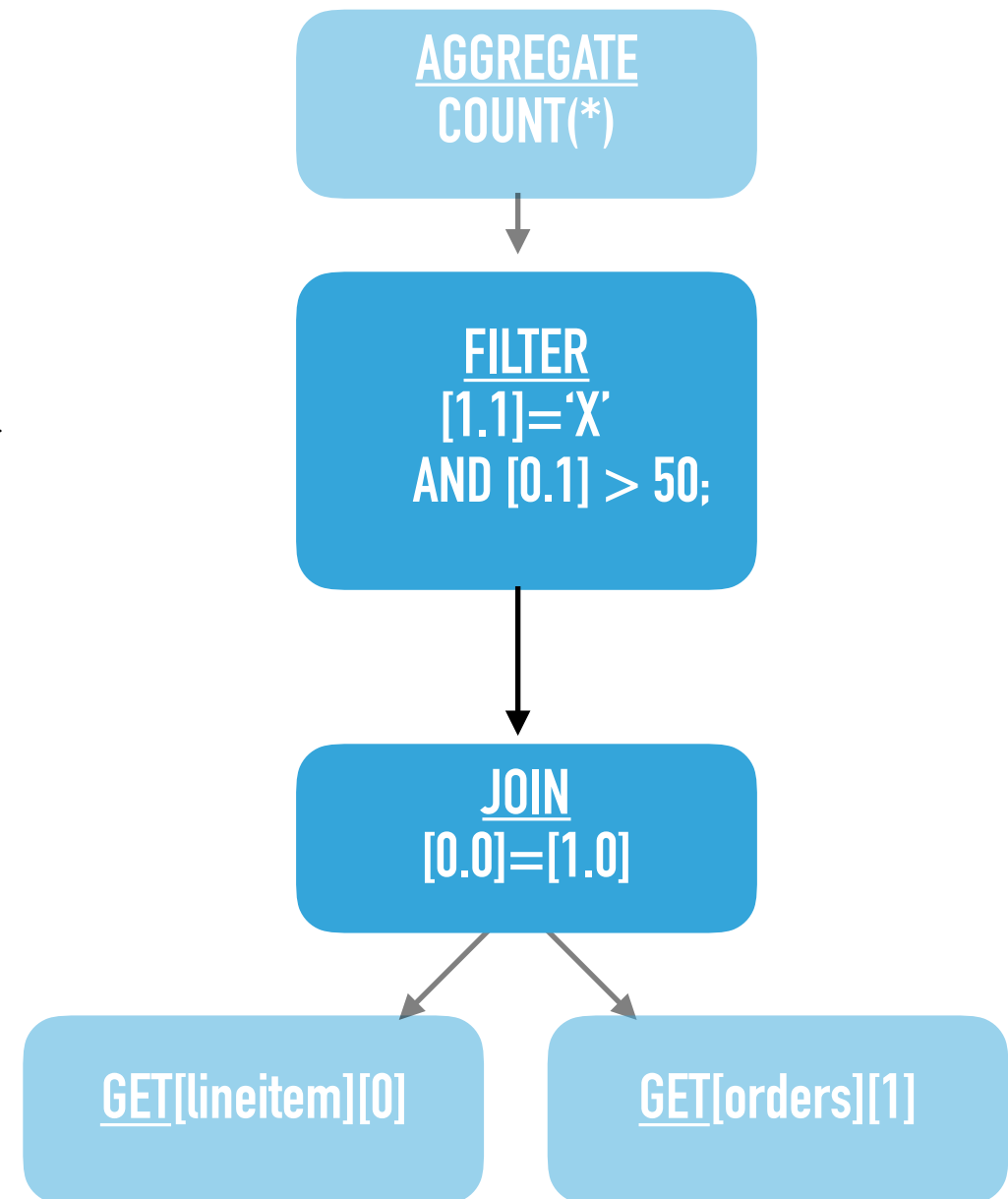
```
SELECT COUNT(*)  
FROM lineitem, orders  
WHERE l_orderkey=o_orderkey  
      AND o_orderstatus='X'  
      AND l_tax > 50;
```

- Pushdown filter into cross product: creates a join

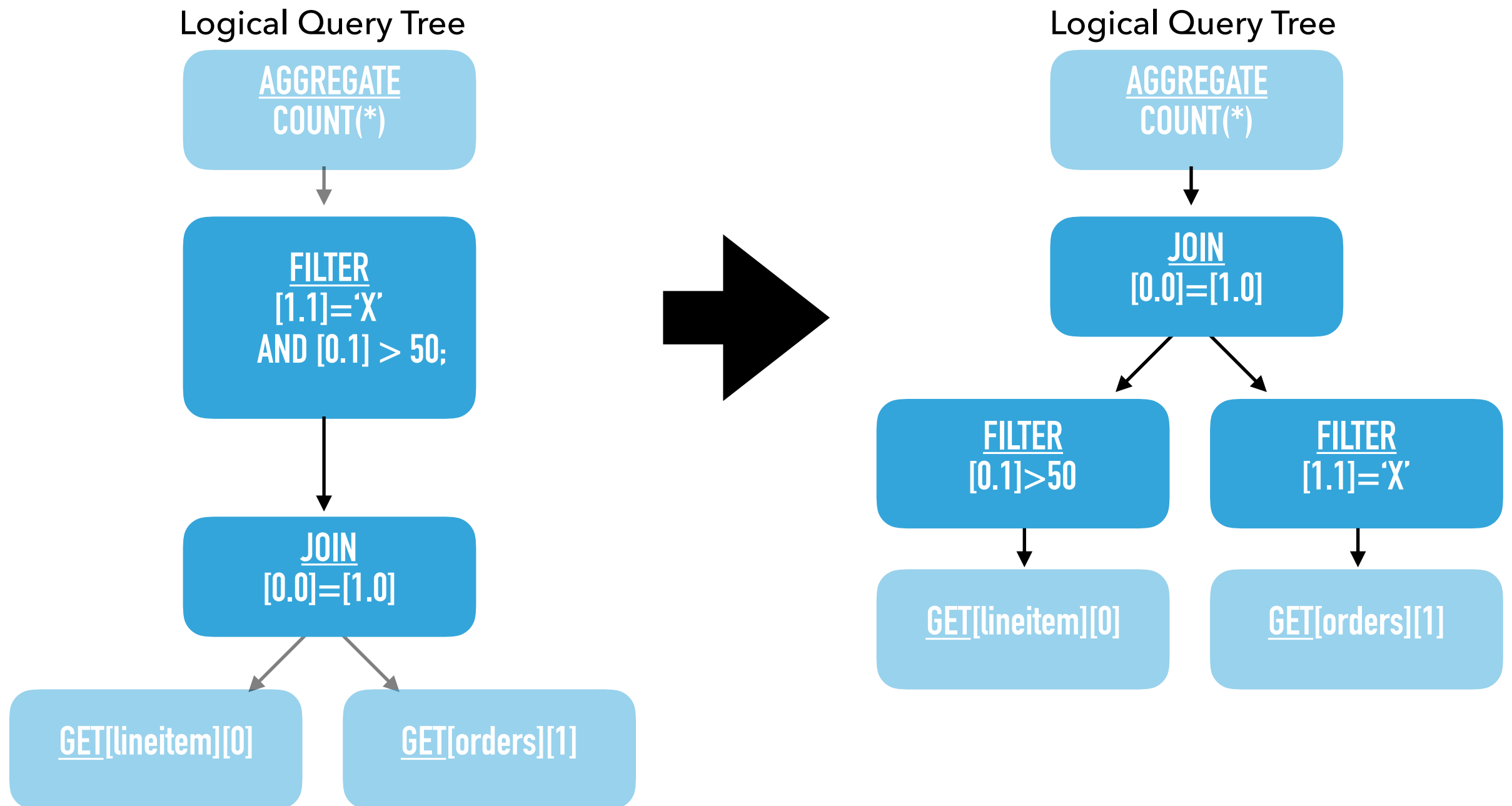
Logical Query Tree



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
 - ▶ Constant folding, CSE, subquery flattening, common subtree elimination, 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

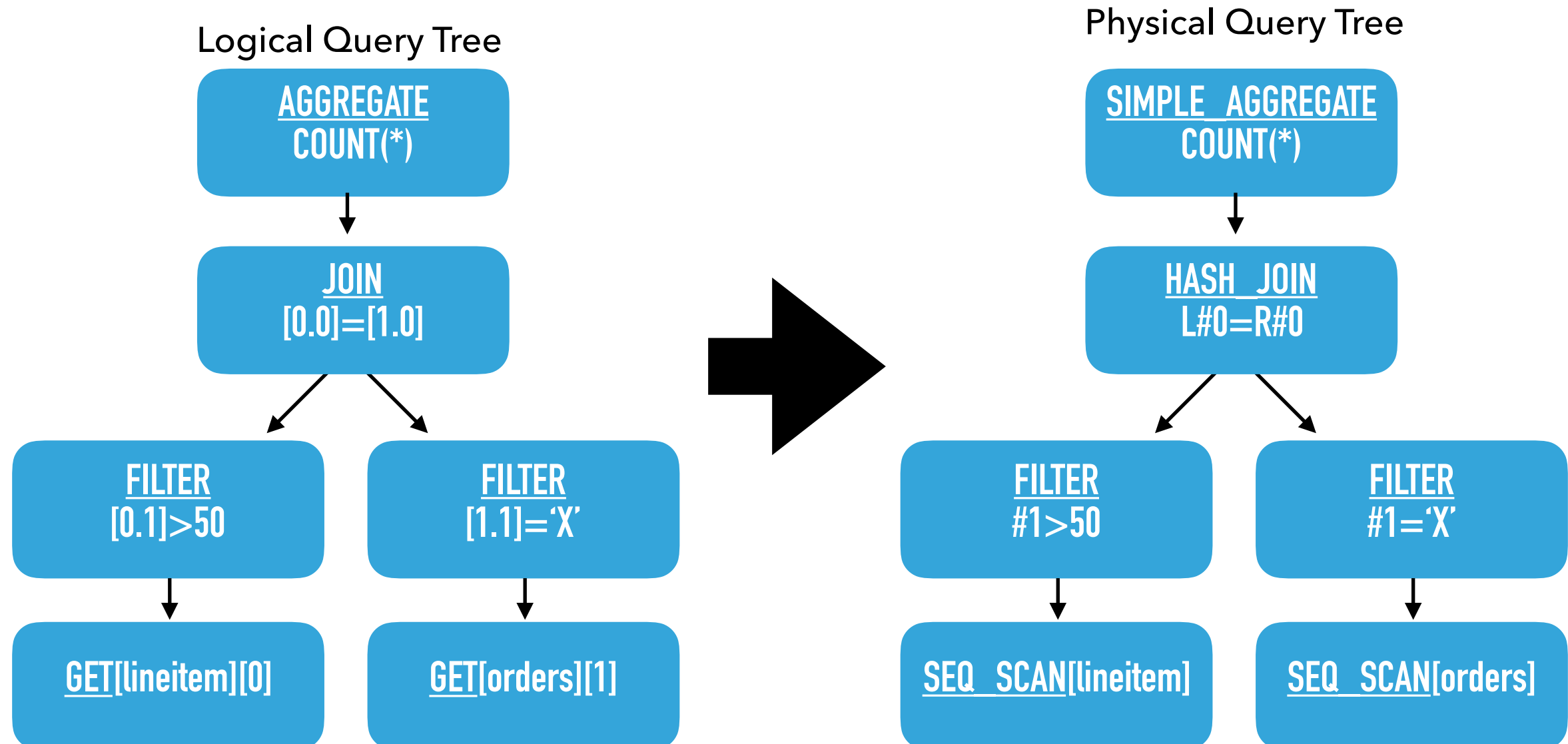
STATEMENT

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

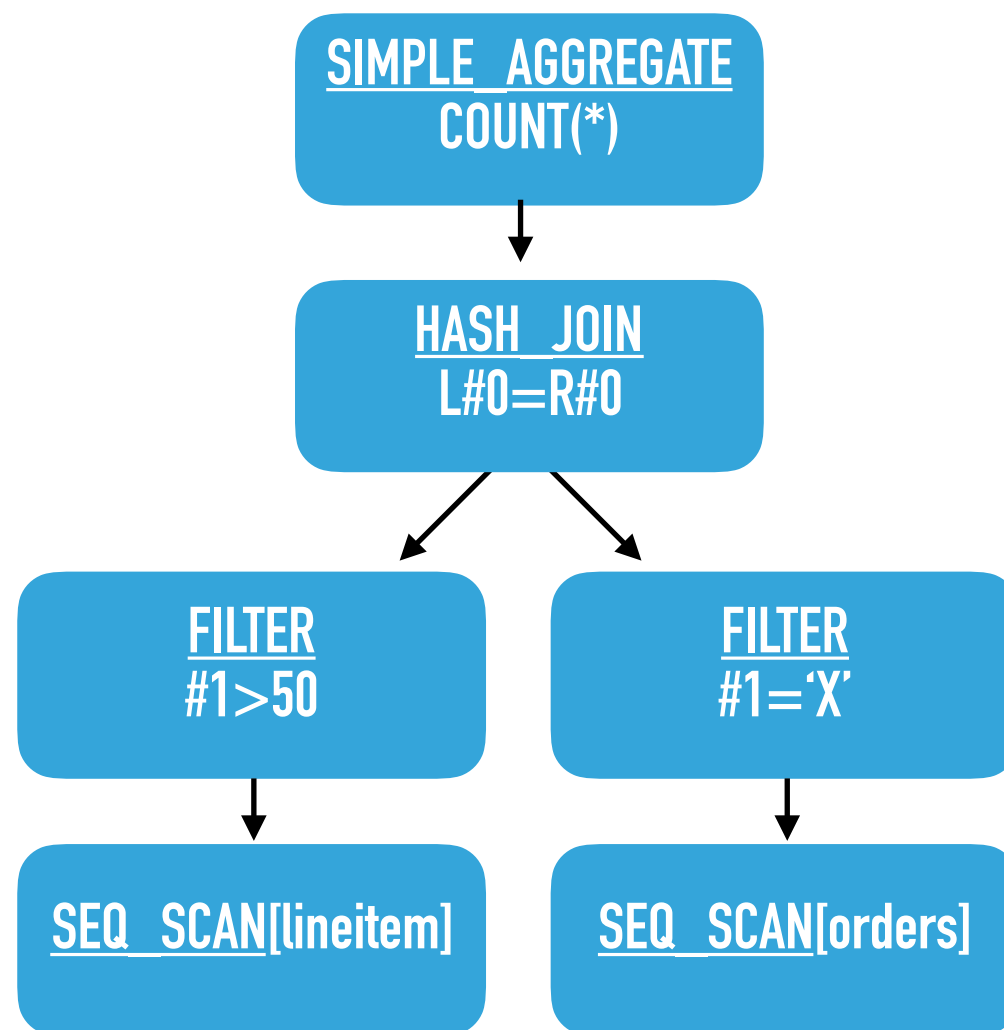
```
SELECT COUNT(*)  
FROM lineitem, orders  
WHERE l_orderkey=o_orderkey  
      AND o_orderstatus='X'  
      AND l_tax > 50;
```

- ▶ **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

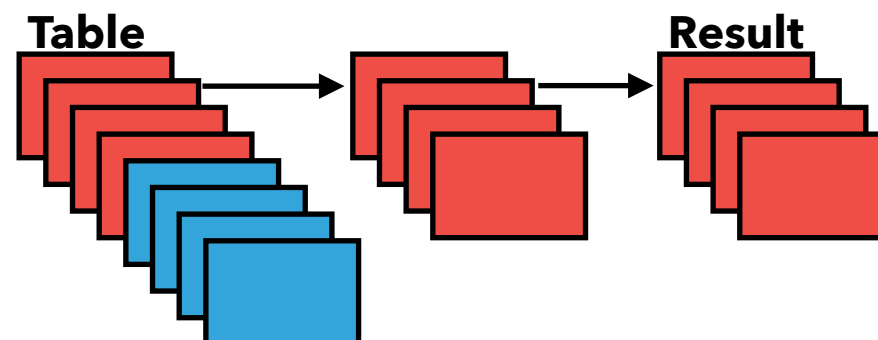
Physical Query Tree



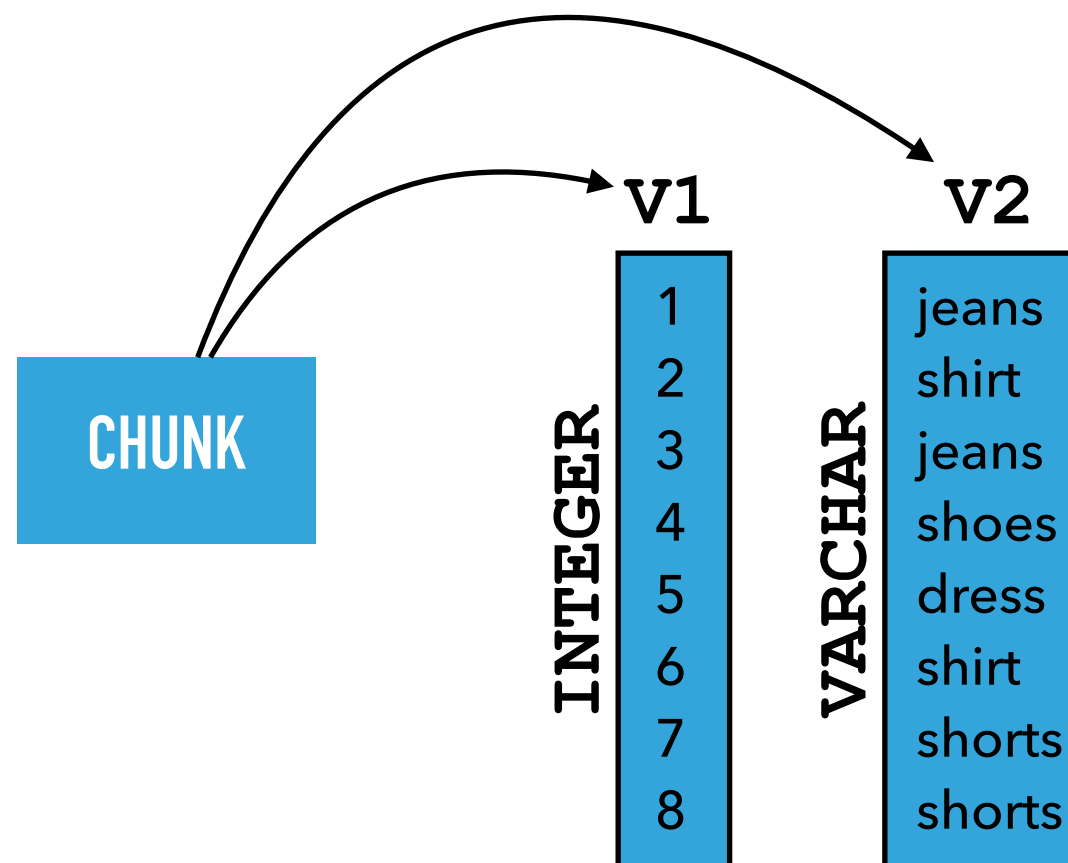
Query Execution

- ▶ 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



- ▶ Basic units: `Vector` and `DataChunk`
- ▶ `Vector` is a **column-slice**
 - ▶ Set of **up to 1024** values of a single type
- ▶ `DataChunk` is a **table-slice** (set of vectors)

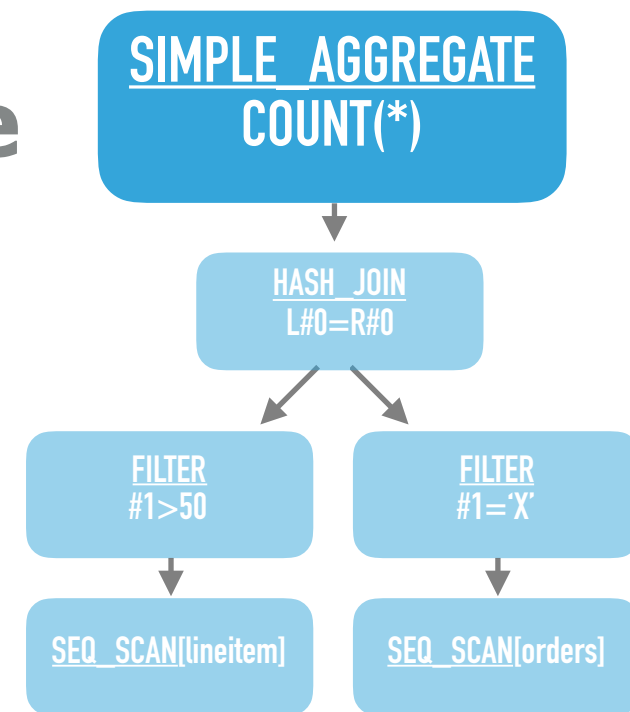


```
class Vector {  
public:  
→   typeId type;  
→   index_t count;  
→   data_ptr_t data;  
→   sel_t *sel_vector;  
→   nullmask_t nullmask;
```

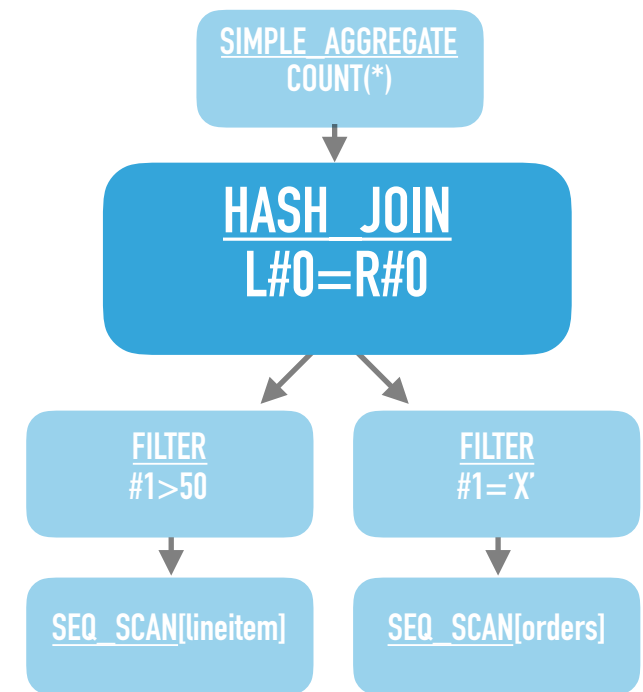
```
class DataChunk {  
public:  
→   index_t column_count;  
→   unique_ptr<Vector[]> data;
```

- ▶ **nullmask**: bitmap indicating which values are NULL
- ▶ **sel_vector**: optional selection vector indicating **which** values to use in the vector

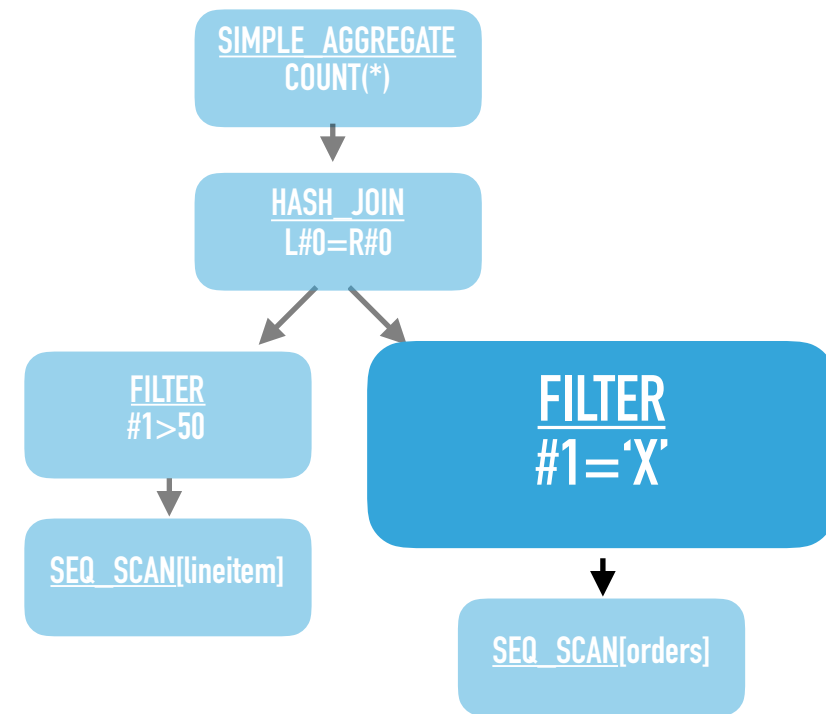
- ▶ Start with root node: **SimpleAggregate**
 - ▶ Aggregate without groups
- ▶ Immediately calls **GetChunk** on child



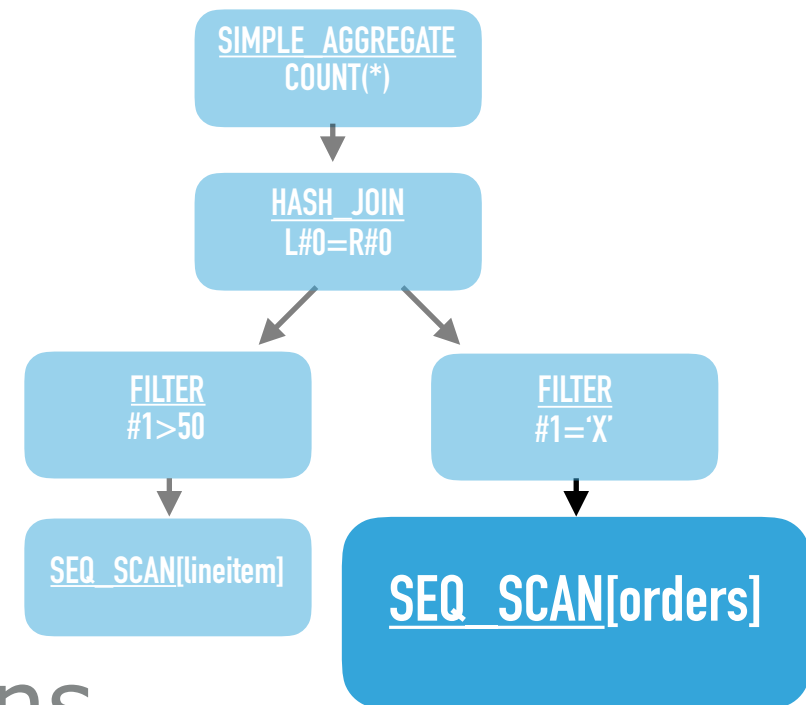
- ▶ Hash Join
- ▶ Start by **building HT**
- ▶ Call **GetChunk** on right node



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



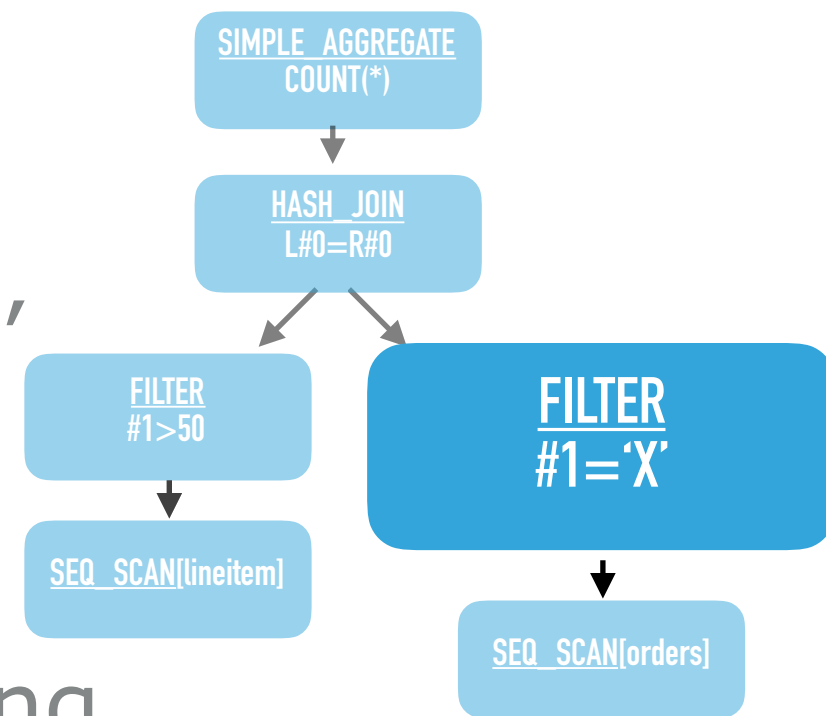
- ▶ **Sequential Scan**
- ▶ Finally we can start executing
- ▶ Scan the base table
- ▶ Return a DataChunk with two columns
 - ▶ `o_orderkey` and `o_orderstatus`



DataChunk			
	V1		V2
INTEGER	1	VARCHAR	N
	2		X
	3		N

► Filter

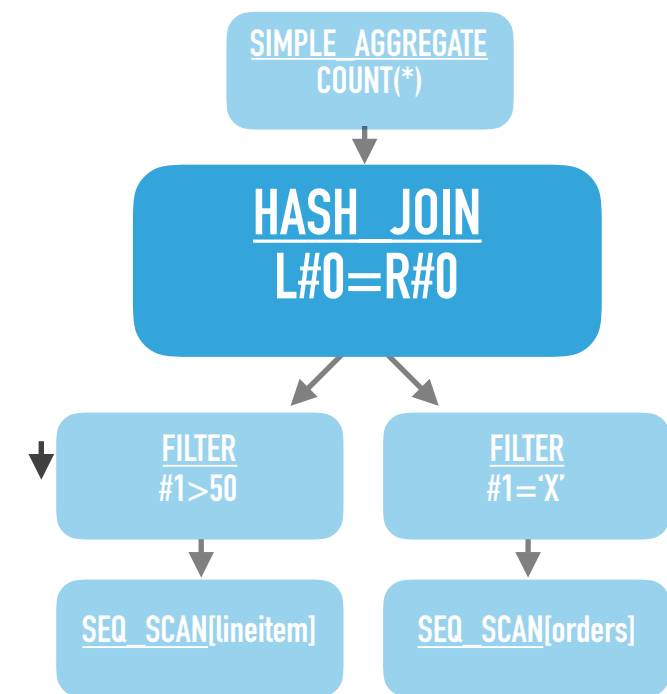
- Now we can perform the filter $\#1 = 'X'$
 - Only the second tuple passes
- **Selection vector** pointing to surviving tuple is created
- Note that no data is copied or changed



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				

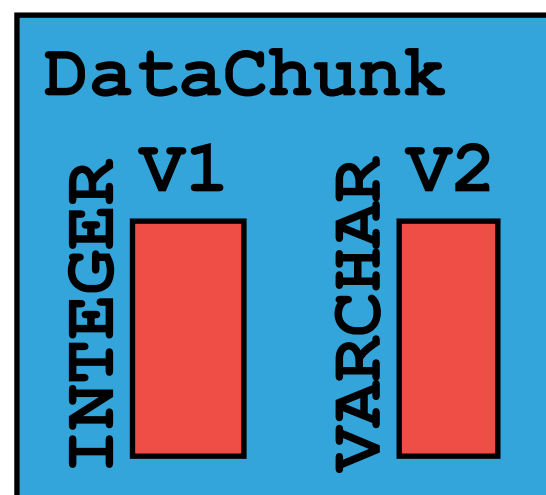
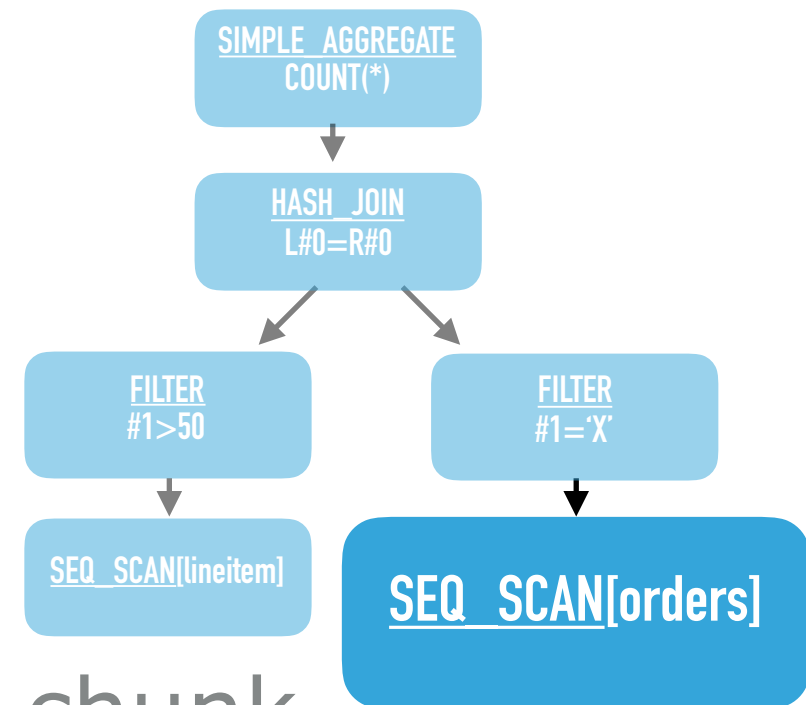
Keys	Payload	Next
2	X	0

► Sequential Scan

► The filter again calls `GetChunk`

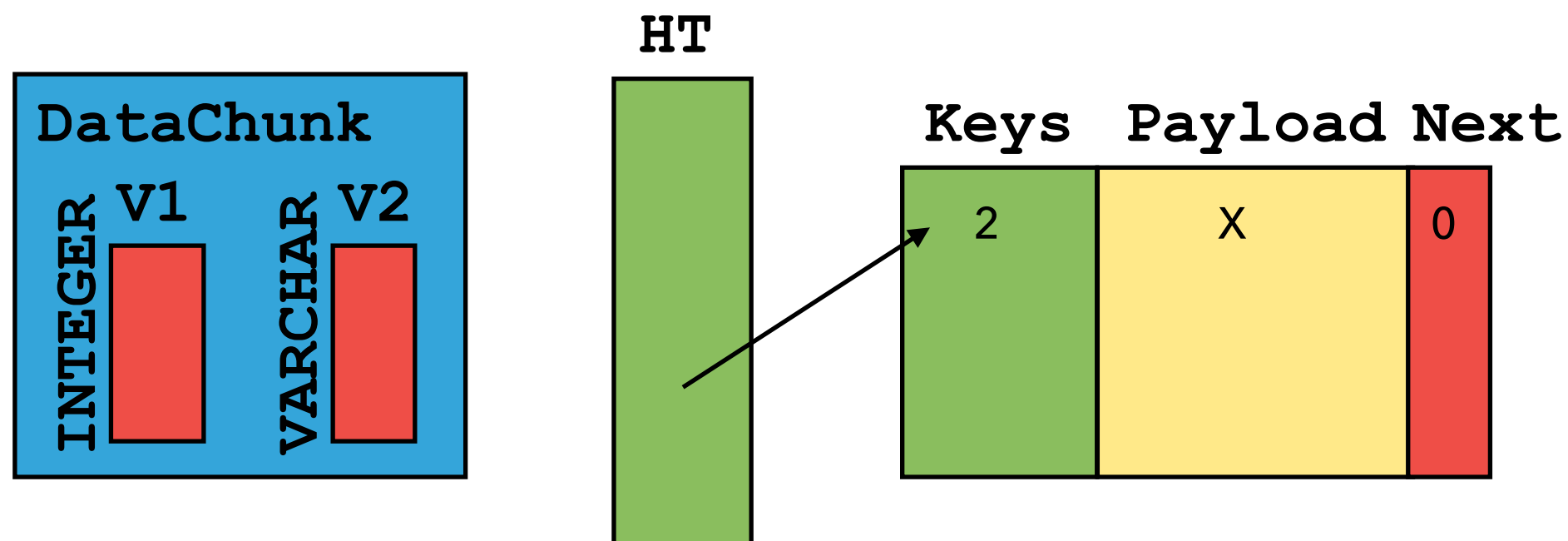
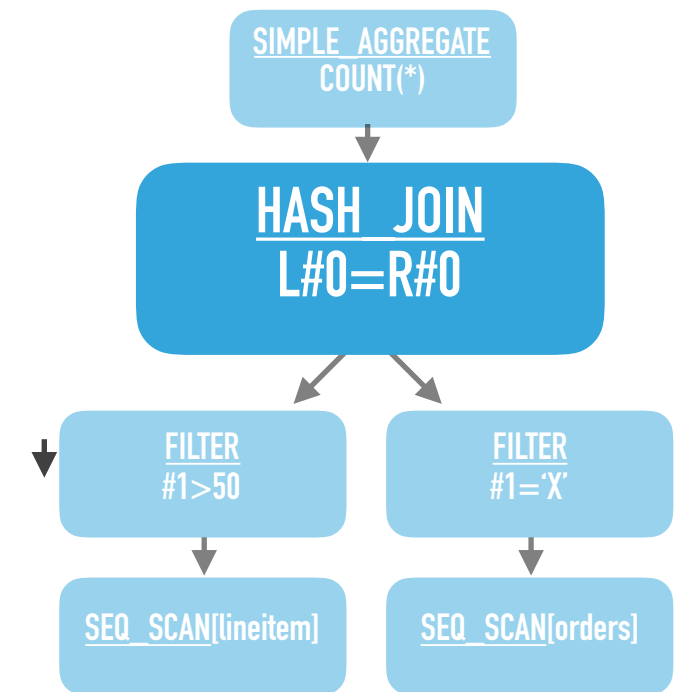
► Scan base table again:

► The scan is finished, return empty chunk

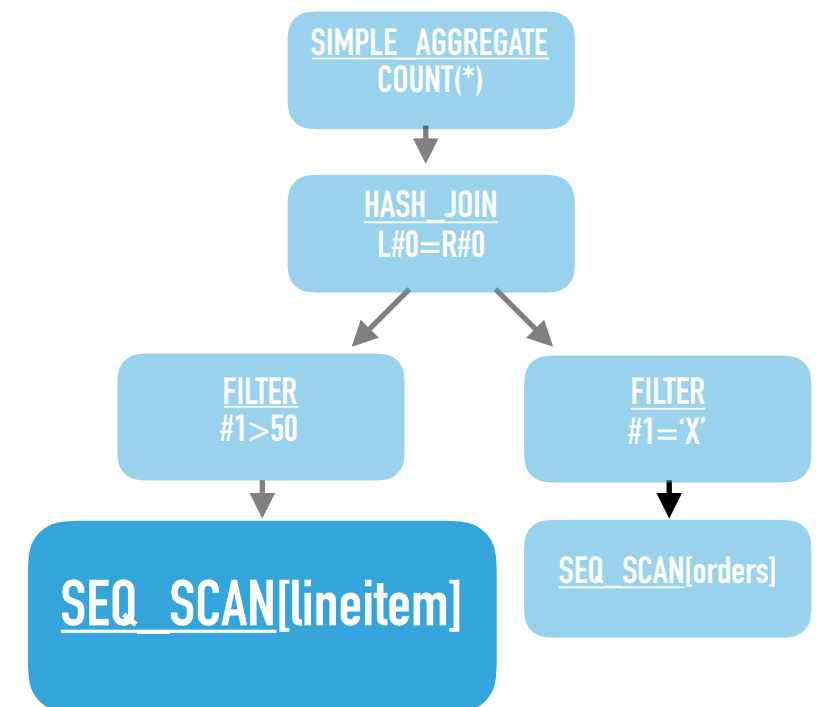


► Hash Join

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

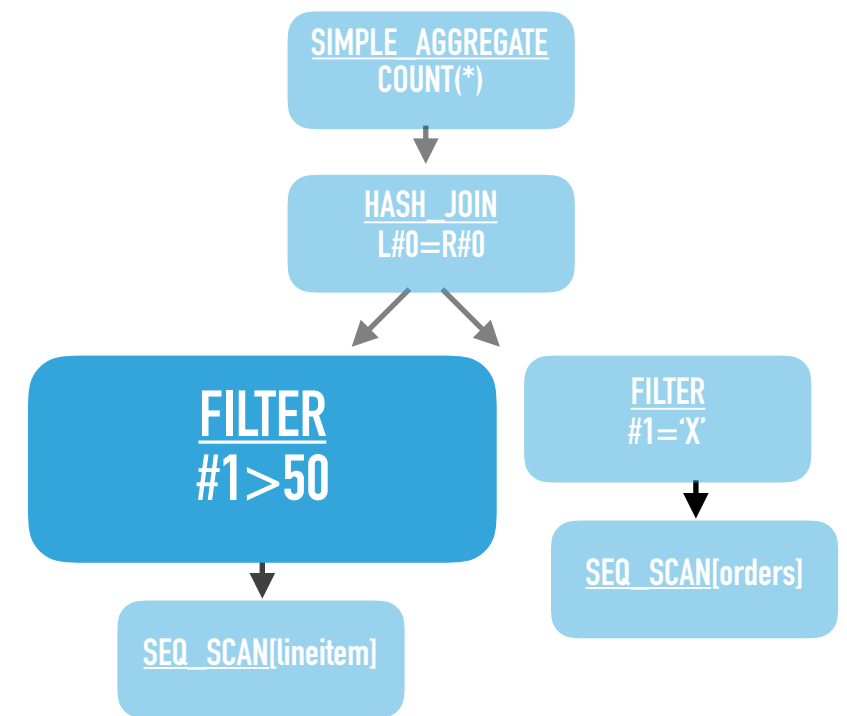


- ▶ **Sequential Scan**
- ▶ We arrive at scan on lineitem
- ▶ DataChunk with two columns
 - ▶ `l_orderkey` and `l_tax`



DataChunk			
	V1		V2
INTEGER	1	VARCHAR	80
	2		60
	2		20

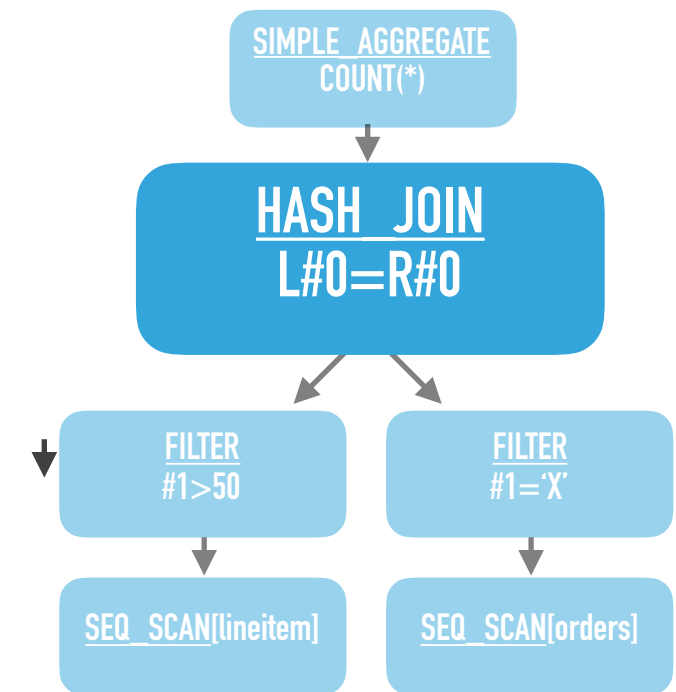
- ▶ **Filter**
- ▶ Performs the filter **#1>50**
- ▶ Again, add a selection vector



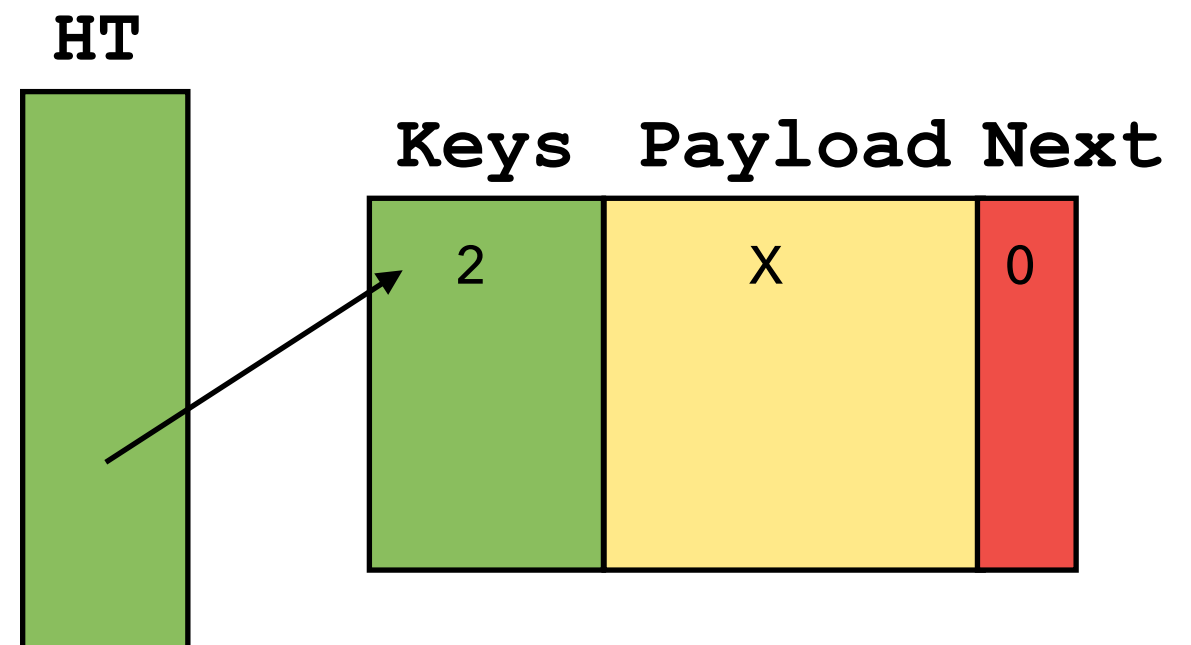
DataChunk				
	V1		V2	SEL
INTEGER	1	VARCHAR	80	0 1
	2		60	
	2		20	

► Hash Join

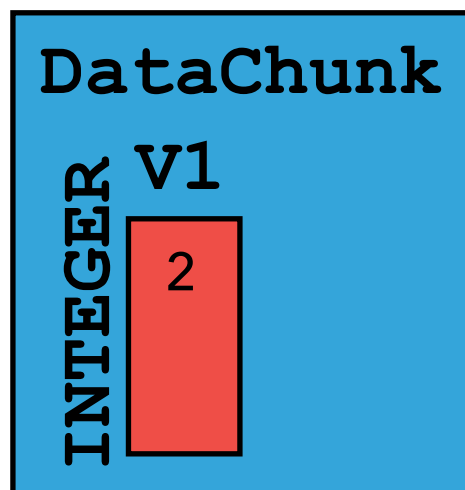
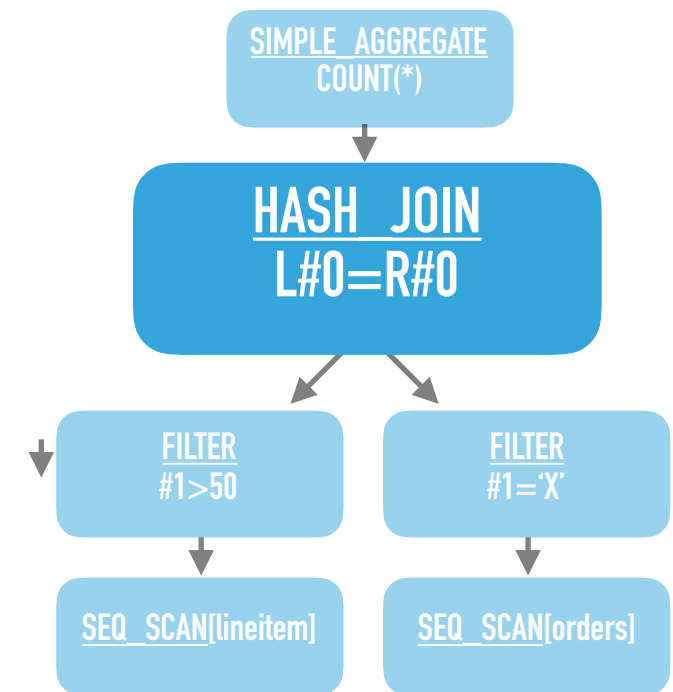
- Now it is time to probe the HT
- We compute the hash for each tuple
- Then lookup in the HT



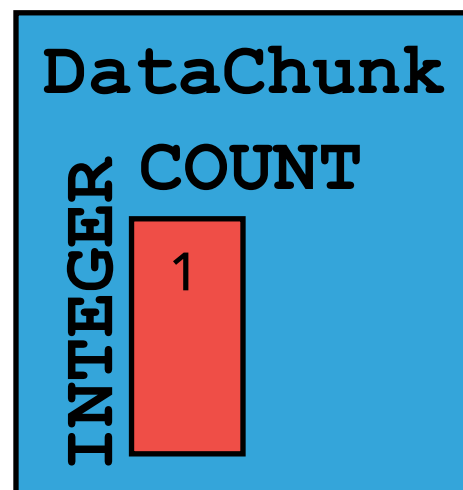
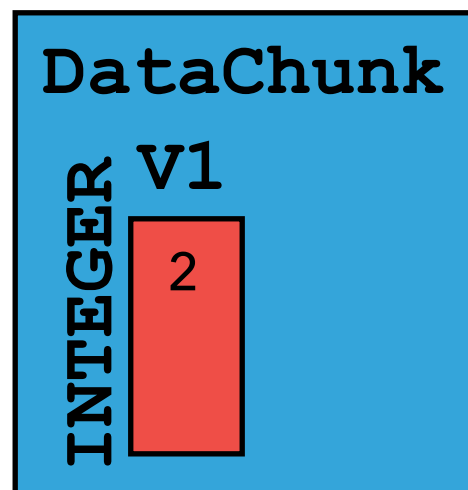
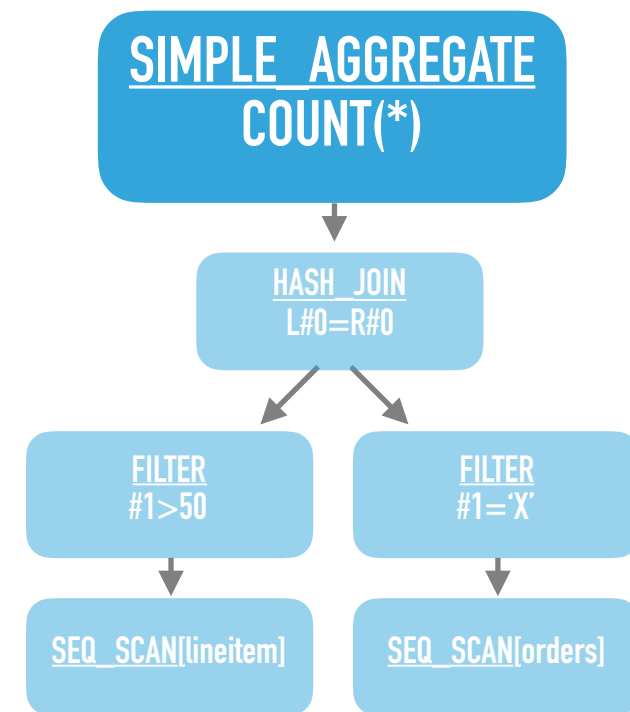
DataChunk				
INTEGER	V1	VARCHAR	V2	SEL
	1		80	0
	2		60	1
	2		20	



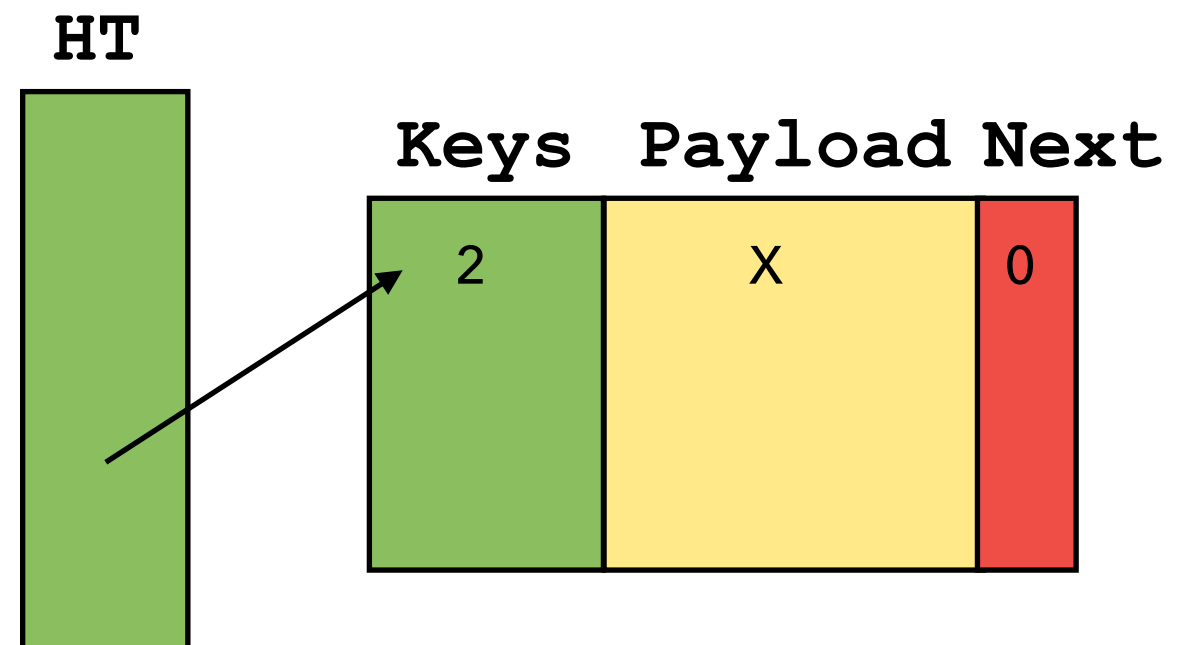
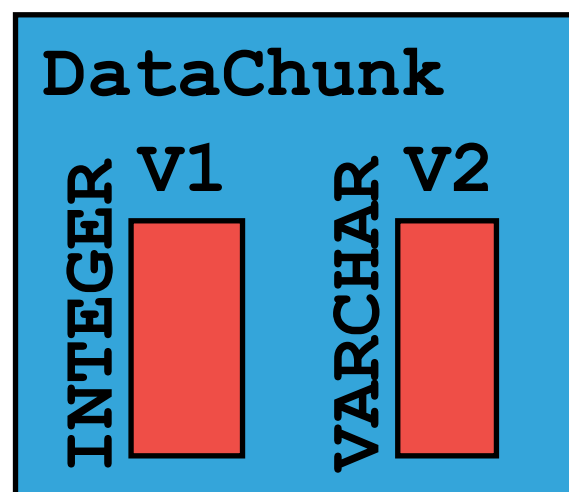
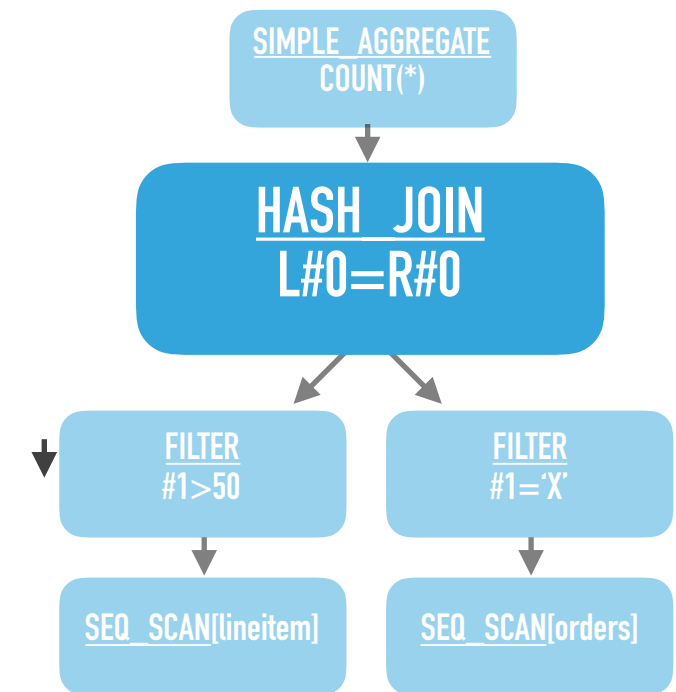
- ▶ **Hash Join**
- ▶ We get one hit on our join!
- ▶ The hash join now produces the result
- ▶ We return this to the aggregate



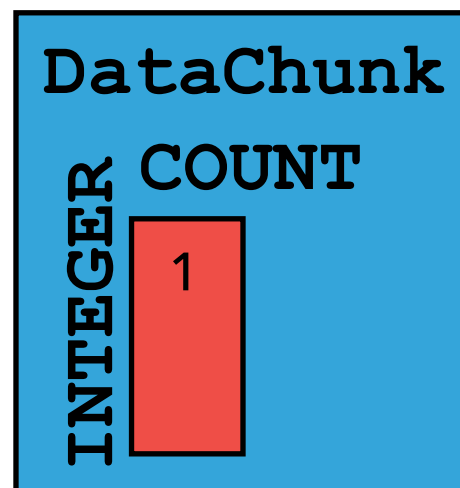
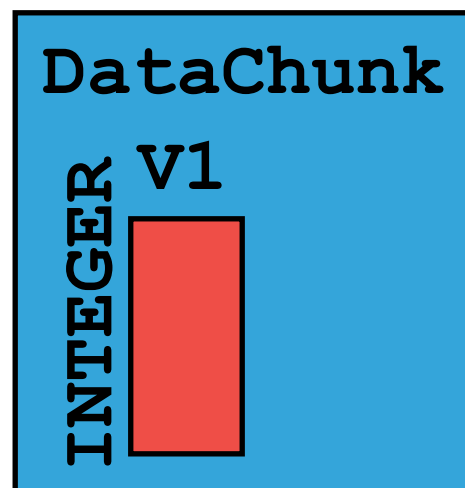
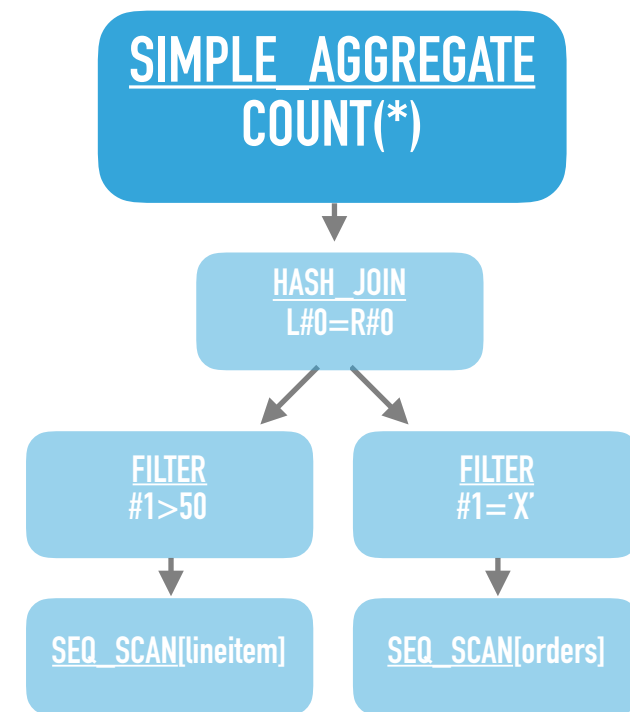
- ▶ The aggregate takes our input chunk
- ▶ Updates the aggregate
- ▶ Then fetches from the child again



- ▶ 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!



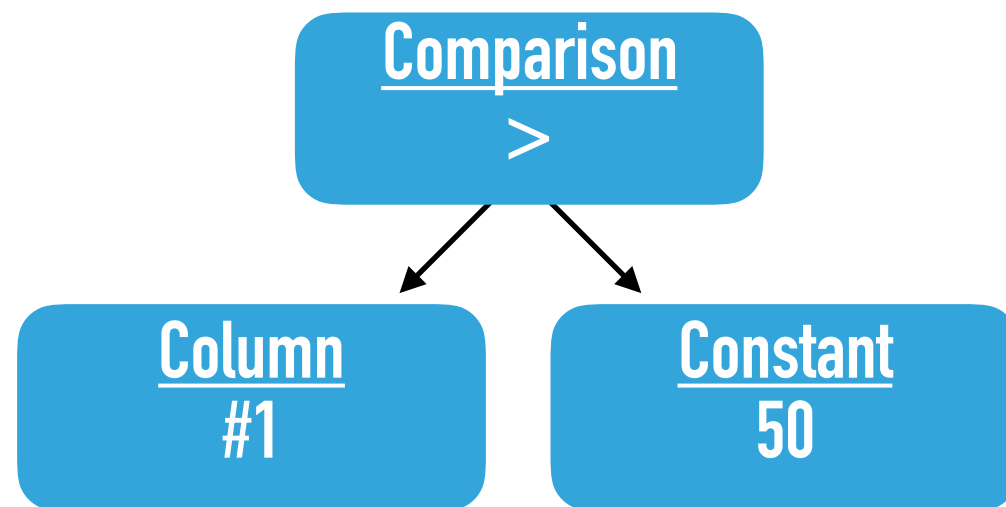
- ▶ Aggregate gets an empty chunk
- ▶ Returns the final result of our query



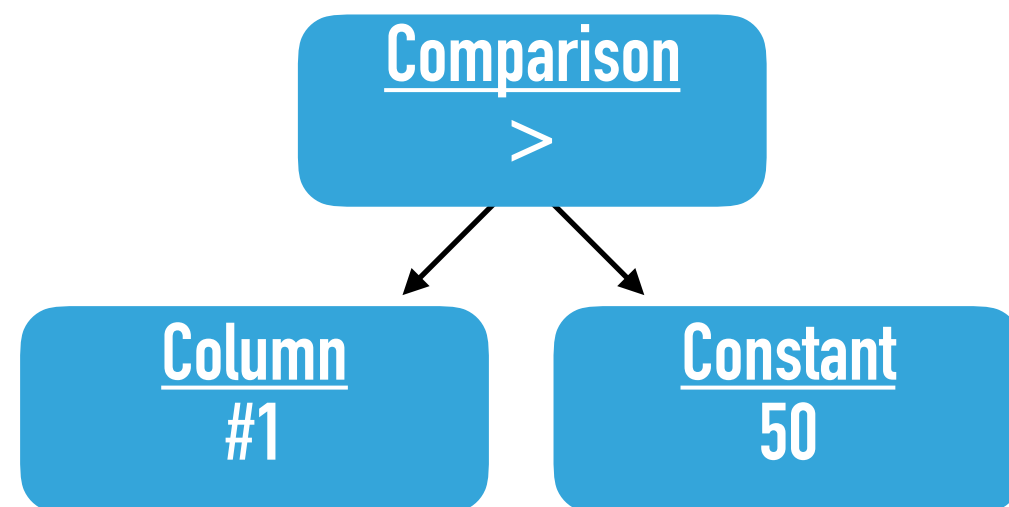
Expression Execution

- ▶ **Expressions** exist within the query tree nodes
 - ▶ Filter has a set of **filter predicates**
 - ▶ Projection has **projection list**
- ▶ Represented as **expression tree**

FILTER
#1>50



- ▶ **ExpressionExecutor** runs the expressions
- ▶ This occurs as part of the execution of the node
- ▶ Expressions are executed in vectorized fashion



FILTER
#1 > 50

Comparison

>

Column
#1

Constant
50

DataChunk

INTEGER

V1

1
2
2

VARCHAR

V2

80
60
20

C1

80
60
20

Column reference
#1 fetches second
column from input

FILTER
#1>50

Comparison

>

Column
#1

Constant
50

DataChunk

INTEGER V1
1
2
2

VARCHAR V2
80
60
20

C1

80
60
20

C2

50

Constant is a
single value

FILTER
#1 > 50

Comparison

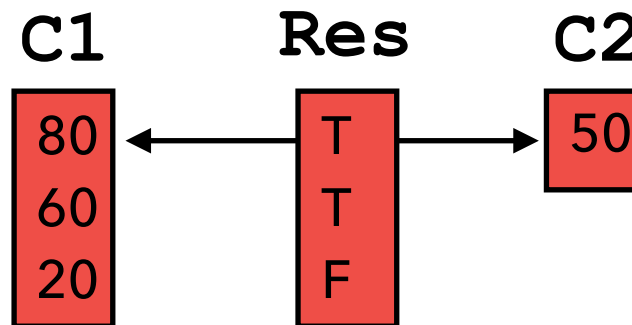
>

Column
#1

Constant
50

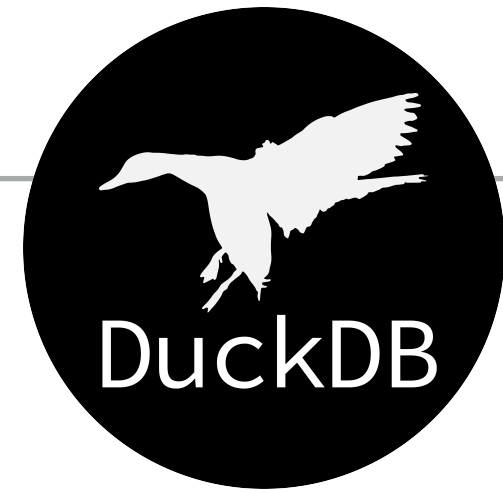
DataChunk

	V1	V2
INTEGER	1	80
	2	60
	2	20



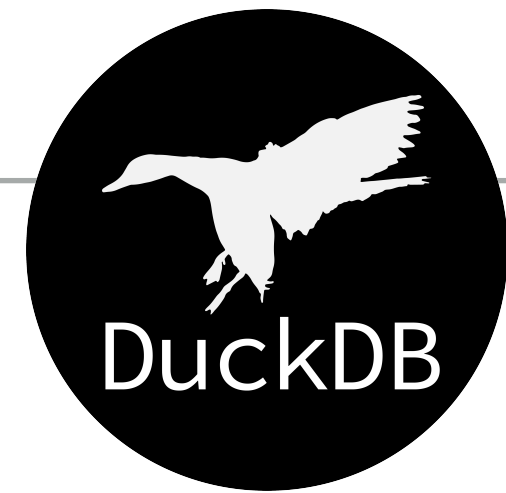
Comparison
runs and returns
matching tuples

Hands On

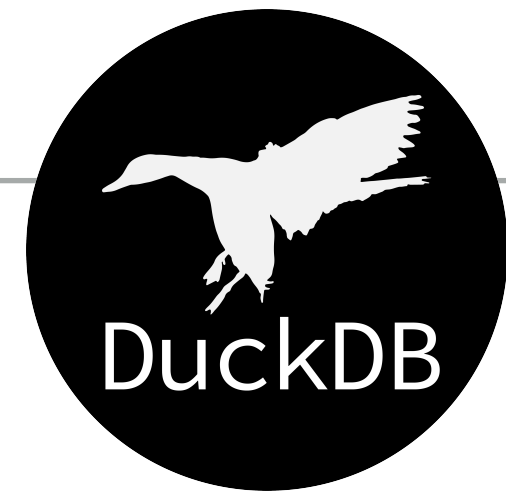


- ▶ **Assignment:** Implement a function in DuckDB
- ▶ Open issues for functions from other systems:
- ▶ <https://github.com/cwida/duckdb/issues/193>
- ▶ Implement one of those
 - ▶ For those that are successful, submit a PR!

Set Up & Testing



- ▶ **Set up:**
- ▶ 1. Download the source code
 - ▶ `git clone https://github.com/cwida/duckdb`
- ▶ 2. Compile the source code
 - ▶ First download CMake if you don't have it
 - ▶ **Linux/OSX:** `make debug`
 - ▶ **Windows:** Use CMake to generate a Visual Studio project, then build it from Visual Studio

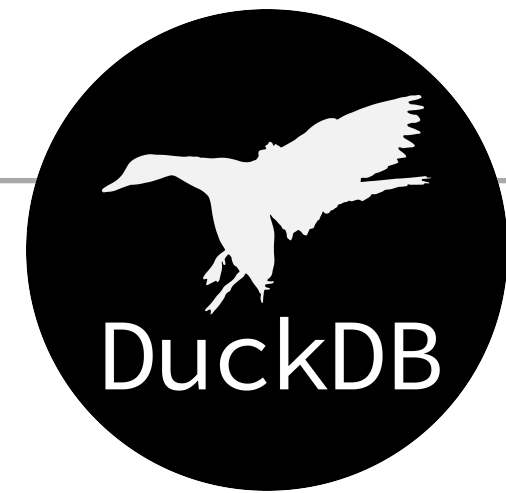


- ▶ Tests are in the `test` directory
 - ▶ We use the Catch framework for tests
- ▶ Tests look like this:

```
TEST_CASE("Test scalar queries", "[scalarquery]") {  
→   unique_ptr<QueryResult> result;  
→   DuckDB db(nullptr);  
→   Connection con(db);  
→   con.EnableQueryVerification();  
  
→   result = con.Query("SELECT 42");  
→   REQUIRE(CHECK_COLUMN(result, 0, {42}));  
  
→   result = con.Query("SELECT 42 + 1");  
→   REQUIRE(CHECK_COLUMN(result, 0, {43}));  
  
→   result = con.Query("SELECT 2 * (42 + 1), 35 - 2");  
→   REQUIRE(CHECK_COLUMN(result, 0, {86}));  
→   REQUIRE(CHECK_COLUMN(result, 1, {33}));  
}
```

Create in-memory database

Run queries
& verify result



- ▶ Tests can be run as follows:
- ▶ **Linux/OSX:**
- ▶ `build/debug/test/unittest "Test scalar queries"`
- ▶ **Windows**
- ▶ Run `unittest project`
- ▶ Command line parameter: `"Test scalar queries"`

```
TEST_CASE("Test scalar queries", "[scalarquery]") {  
→   unique_ptr<QueryResult> result;  
→   DuckDB db(nullptr);  
→   Connection con(db);  
→   con.EnableQueryVerification();  
→  
→   result = con.Query("SELECT 42");  
}
```

Function Definition

- ▶ Each function has different **overloads**

Addition

+

- ▶ e.g. addition operator:
 - ▶ `+(SMALLINT,SMALLINT)`
 - ▶ `+(INTEGER,INTEGER)`
 - ▶ `+(BIGINT,BIGINT)`
 - ▶ ...
- ▶ **Binder** chooses which version to use

- ▶ Set of permitted **implicit casts**

Addition

+

- ▶ TINYINT → SMALLINT, INTEGER, BIGINT, FLOAT, DOUBLE
- ▶ SMALLINT → INTEGER, BIGINT, FLOAT, DOUBLE
- ▶ INTEGER → BIGINT, FLOAT, DOUBLE
- ▶ BIGINT → FLOAT, DOUBLE
- ▶ FLOAT → DOUBLE

Addition

+

- ▶ Binder prefers to cast as little as possible
- ▶ e.g. TINYINT + INTEGER has multiple eligible options
- ▶ INTEGER + INTEGER will be chosen
 - ▶ Requires only one implicit cast
- ▶ Other options require two casts:
 - ▶ BIGINT + BIGINT, FLOAT+FLOAT, DOUBLE + DOUBLE

- ▶ The same binding rules apply to **functions**
 - ▶ `substring(string, start, length)`
 - ▶ Three parameters: `VARCHAR`, `INTEGER`, `INTEGER`
- ▶ Binder will automatically insert `CAST` if required
 - ▶ e.g. `TINYINT` → `INTEGER`
- ▶ In the code for **substring** we only need to implement the case with parameters `VARCHAR`, `INTEGER`, `INTEGER`

- ▶ Code: how to add a function definition

```
set.AddFunction(ScalarFunction(  
→  "substring", . . . . . // name of function  
→  {  SQLType::VARCHAR, . . // argument list  
→    |  SQLType::INTEGER,  
→    |  SQLType::INTEGER },  
→  SQLType::VARCHAR, . . . . // return type  
→  substring_function)); // pointer to function implementation
```

- ▶ Function code is implemented in `substring_function`

Creating a Simple Function

- ▶ Create a simple function:
 - ▶ `add_one (INTEGER) -> INTEGER`
- ▶ This function adds one to its integer input
- ▶ Returns the result

- ▶ **Step one: Create tests**
- ▶ Navigate to `test/sql/function`
- ▶ Create a new file: `test_add_one.cpp`
- ▶ Add it to `CMakeLists.txt` in that folder

► Step one: Create tests

Table + selection vector tests

► Step one: Create tests

```
> build/debug/test/unittest "Test add one function"
Query failed with message: Catalog: Function with name add_one does not exist!

~~~~~

unittest is a Catch v2.4.0 host application.
Run with -? for options

-----

Test add one function
-----

/Users/myth/Programs/duckdb/test/sql/function/test_add_one.cpp:7
.....

/Users/myth/Programs/duckdb/test/sql/function/test_add_one.cpp:18: FAILED:
    REQUIRE( CHECK_COLUMN(result, 0, {2}) )
with expansion:
    false

=====

test cases: 1 | 1 failed
assertions: 3 | 2 passed | 1 failed
```

- ▶ **Step two: Create the function**
- ▶ Navigate to `src/function/scalar`
- ▶ All function implementations are here
- ▶ In `math` directory, create new file: `add_one.cpp`
 - ▶ And add it to the `CMakeLists.txt`

- ▶ **Step two: Create the function**
- ▶ Add code to register function:

```
void AddOne::RegisterFunction(BuiltinFunctions &set) {  
→   set.AddFunction(ScalarFunction("add_one",  
→       { SQLType::INTEGER },  
→       SQLType::INTEGER,  
→       add_one_function));  
}
```

- ▶ **Step two: Create the function**
- ▶ Now add actual function code:

```
static void add_one_function(  
    ... ExpressionExecutor &exec,  
    Vector inputs[] index_t input_count,  
    BoundFunctionExpression &expr, Vector &result {  
    result.Initialize(TypeId::INTEGER);  
    VectorOperations::UnaryExec<int32_t, int32_t>(  
        inputs[0], result, [&](int32_t input) {  
            return input + 1;  
        });  
}
```

Initialize result

Loop over input & compute result

- ▶ **Step two: Create the function**
- ▶ Finally add some more bookkeeping code:
- ▶ `include/function/scalar/math_functions.hpp`

```
struct AddOne {  
    static void RegisterFunction(BuiltinFunctions &set);  
};
```

- ▶ `function/scalar/math_functions.cpp`

```
void BuiltinFunctions::RegisterMathFunctions() {  
    Register<AddOne>();  
}
```

- ▶ **Step two: Create the function**
- ▶ Now run the tests

```
> build/debug/test/unittest "Test add one function"
```

```
=====
All tests passed (6 assertions in 1 test case)
```

- ▶ Everything passes!

- ▶ Time to implement your own function
- ▶ **Advice:** Start with the `add_one` function
- ▶ Once that works, move on to different functions

▶ Suggestions:

▶ `RTRIM (VARCHAR) -> VARCHAR` [MySQL]

▶ Remove spaces on right side of string

▶ `REVERSE (VARCHAR) -> VARCHAR` [MySQL]

▶ Reverse characters of a string

▶ `REPEAT (VARCHAR, INTEGER) -> VARCHAR` [MySQL]

▶ Repeat the specified string a number of times

▶ `INSTR (VARCHAR, VARCHAR) -> BOOL` [SQLite]

▶ Returns true if second string is part of first string

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