

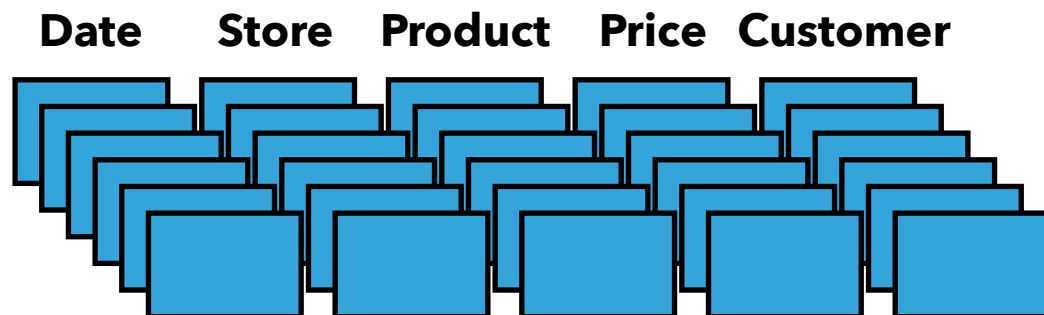
Mark Raasveldt & Pedro Holanda

DuckDB **an Embeddable Analytical RDBMS**

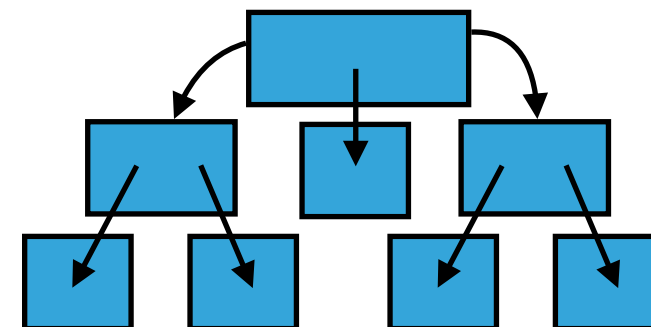
- ▶ Internals at a Glance
 - ▶ Column-Store
 - ▶ MVCC
 - ▶ ART Index
 - ▶ Storage
- ▶ Query processing pipeline
- ▶ Query execution
- ▶ **Hands-On**

Internals at a Glance

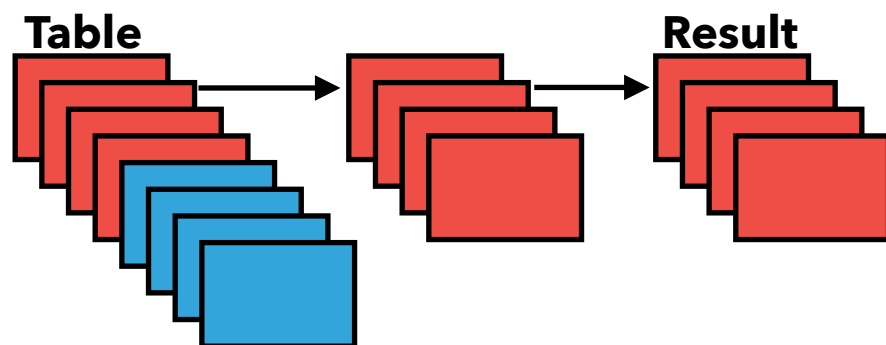
Column-Store



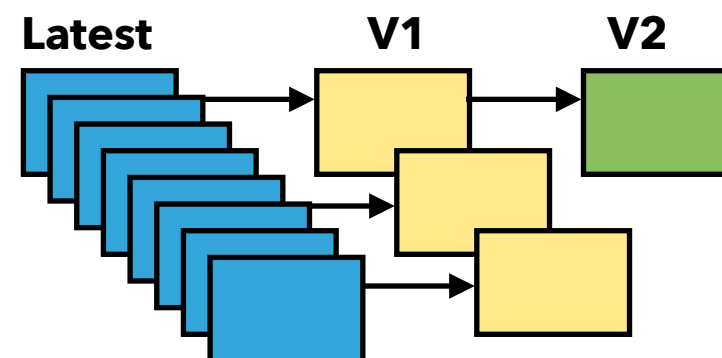
ART Index



Vectorized Processing



MVCC



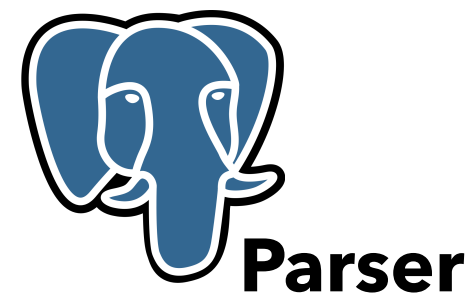
Single-File Storage



4KB

256KB

database.db

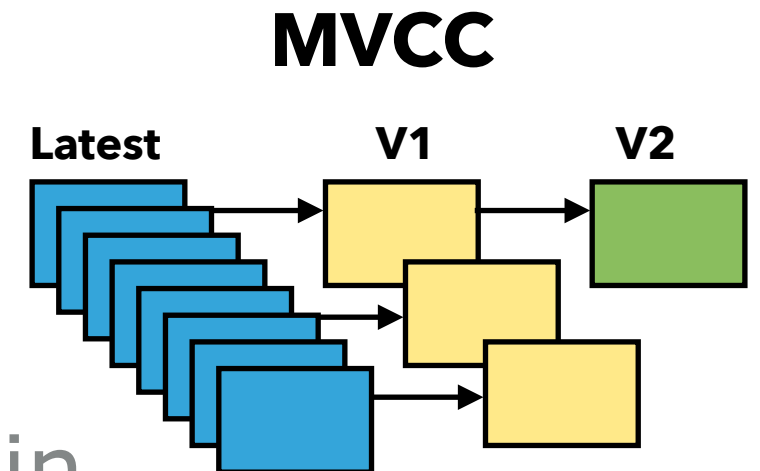


- ▶ DuckDB **vertically partitions** tables
 - ▶ This is also called a **column-store**
- ▶ **Advantages:**
 - ▶ Allows individual columns to be fetched/updated
 - ▶ Better compression ratio
 - ▶ Faster scans
- ▶ **Disadvantages:**
 - ▶ Less efficient fetch/update of individual rows



- [illegible]

- ▶ MVCC is used to provide **snapshot isolation**
- ▶ Data is updated **in place**
 - ▶ Avoids fragmentation
- ▶ Old versions are kept in version chain
 - ▶ Deleted once no longer required
- ▶ Version chain starts out empty
 - ▶ No overhead unless versions exist



Two transactions:

T1

T2

Date

0	1992-01-01
1	1992-01-02
2	1992-01-03
3	1992-01-04
4	1992-01-05
5	1992-01-06
6	1992-01-07
7	1992-01-08
8	1992-01-09
9	1992-01-10
10	1992-01-11

Two transactions:

T1

T2

Date

0	1992-01-01
1	1992-01-02
2	1992-01-03
3	1992-01-04
4	<u>1992-01-13</u>
5	1992-01-06
6	<u>1992-01-15</u>
7	1992-01-08
8	1992-01-09
9	1992-01-10
10	1992-01-11

V1

4 1992-01-05
6 1992-01-07

T1

Two transactions:

T1

T2

Date

0	1992-01-22
1	1992-01-27
2	1992-01-28
3	1992-01-04
4	1992-01-13
5	1992-01-06
6	1992-01-15
7	1992-01-08
8	1992-01-09
9	1992-01-10
10	1992-01-11

V2

1 1992-01-01
2 1992-01-02
3 1992-01-03

T2

V1

4 1992-01-05
6 1992-01-07

T1

Two transactions:

T1

T2

Date

0	<u>1992-01-22</u>
1	<u>1992-01-27</u>
2	<u>1992-01-28</u>
3	1992-01-04
4	<u>1992-01-13</u>
5	1992-01-06
6	<u>1992-01-15</u>
7	1992-01-08
8	1992-01-09
9	1992-01-10
10	1992-01-11

V2

1	1992-01-01
2	1992-01-02
3	1992-01-03

T2

V1

4	1992-01-05
6	1992-01-07

T1

T1

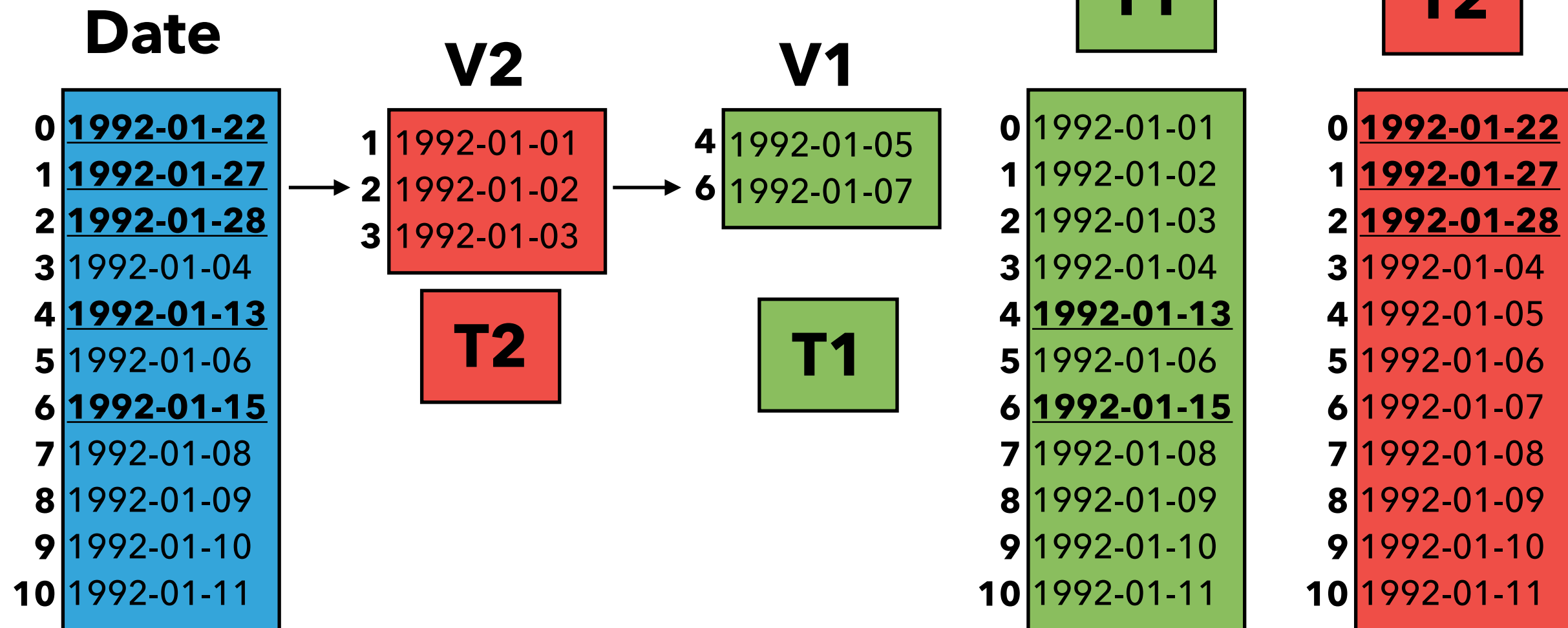
0	1992-01-01
1	1992-01-02
2	1992-01-03
3	1992-01-04
4	<u>1992-01-13</u>
5	1992-01-06
6	<u>1992-01-15</u>
7	1992-01-08
8	1992-01-09
9	1992-01-10
10	1992-01-11

T2

0	<u>1992-01-22</u>
1	<u>1992-01-27</u>
2	<u>1992-01-28</u>
3	1992-01-04
4	1992-01-05
5	1992-01-06
6	1992-01-07
7	1992-01-08
8	1992-01-09
9	1992-01-10
10	1992-01-11

After commit of

T1



After commit of

T2**Date**

0	1992-01-22
1	1992-01-27
2	1992-01-28
3	1992-01-04
4	1992-01-13
5	1992-01-06
6	1992-01-15
7	1992-01-08
8	1992-01-09
9	1992-01-10
10	1992-01-11

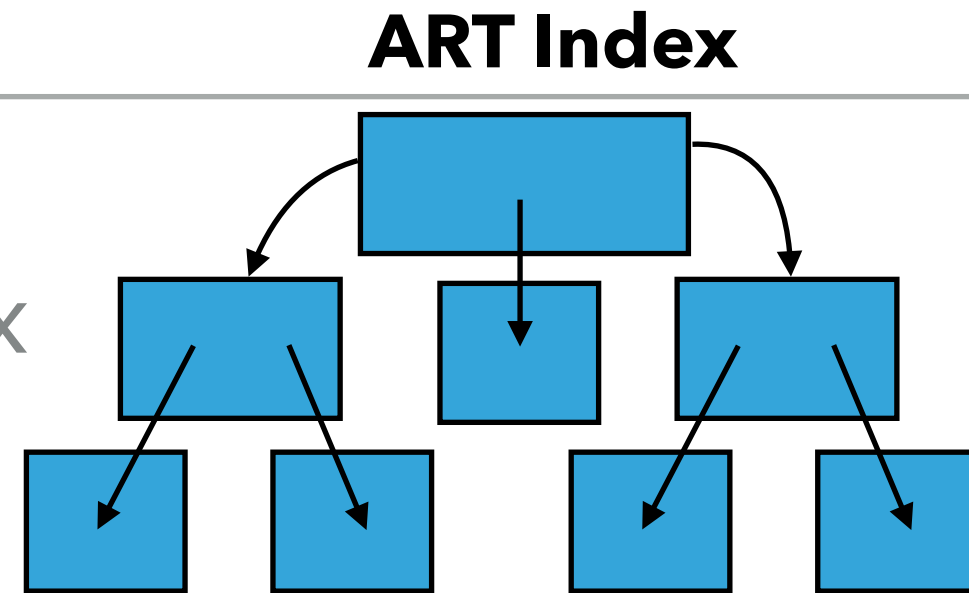
T1

0	1992-01-22
1	1992-01-27
2	1992-01-28
3	1992-01-04
4	1992-01-13
5	1992-01-06
6	1992-01-15
7	1992-01-08
8	1992-01-09
9	1992-01-10
10	1992-01-11

T2

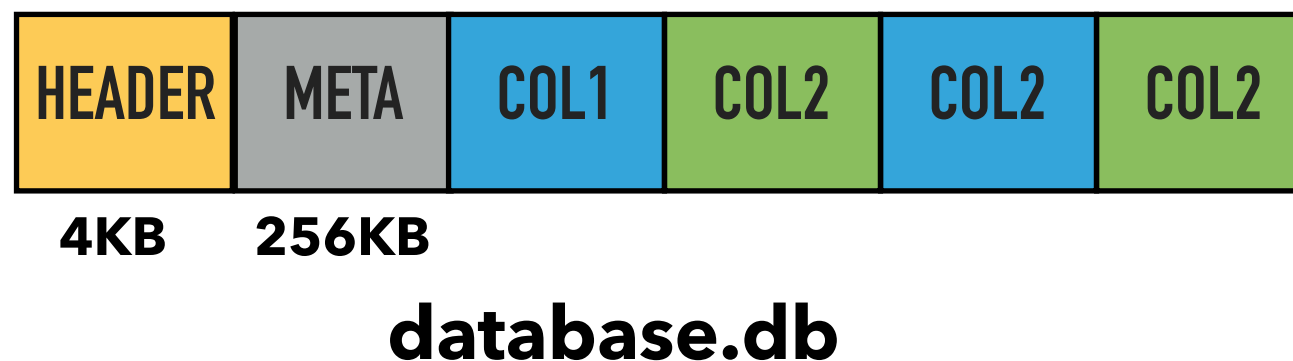
0	1992-01-22
1	1992-01-27
2	1992-01-28
3	1992-01-04
4	1992-01-13
5	1992-01-06
6	1992-01-15
7	1992-01-08
8	1992-01-09
9	1992-01-10
10	1992-01-11

- ▶ DuckDB has support for ART index
 - ▶ **Explicitly:** `CREATE INDEX`
 - ▶ **Automatically:** `PRIMARY KEY/UNIQUE`
- ▶ Only unclustered indexes are supported
 - ▶ In leaf nodes, ART index contains **row ids**
 - ▶ To reconstruct tuples, base table must be accessed
- ▶ ART index only speeds up **selective queries (<1%)**



- ▶ Data is stored in a **single file**
 - ▶ **Purpose:** user convenience and simplicity
- ▶ File starts with a small header (4KB)
- ▶ Rest is divided into **equal-sized blocks**
 - ▶ Currently *256KB*
- ▶ Blocks are referred to by **block_id**

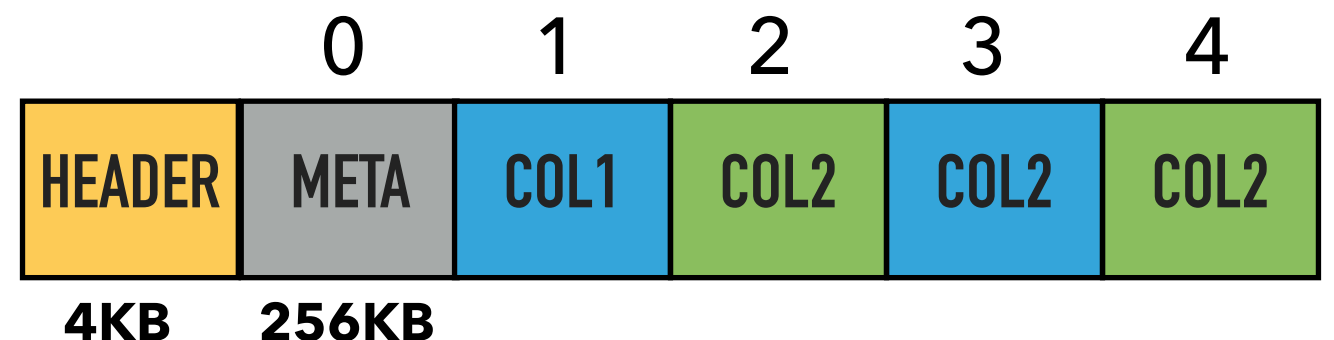
Single-File Storage



- ▶ Meta data contains **schema information**
 - ▶ Schema, table, column names
 - ▶ Column types
 - ▶ **Column Pointers:** Pointers to blocks where data is

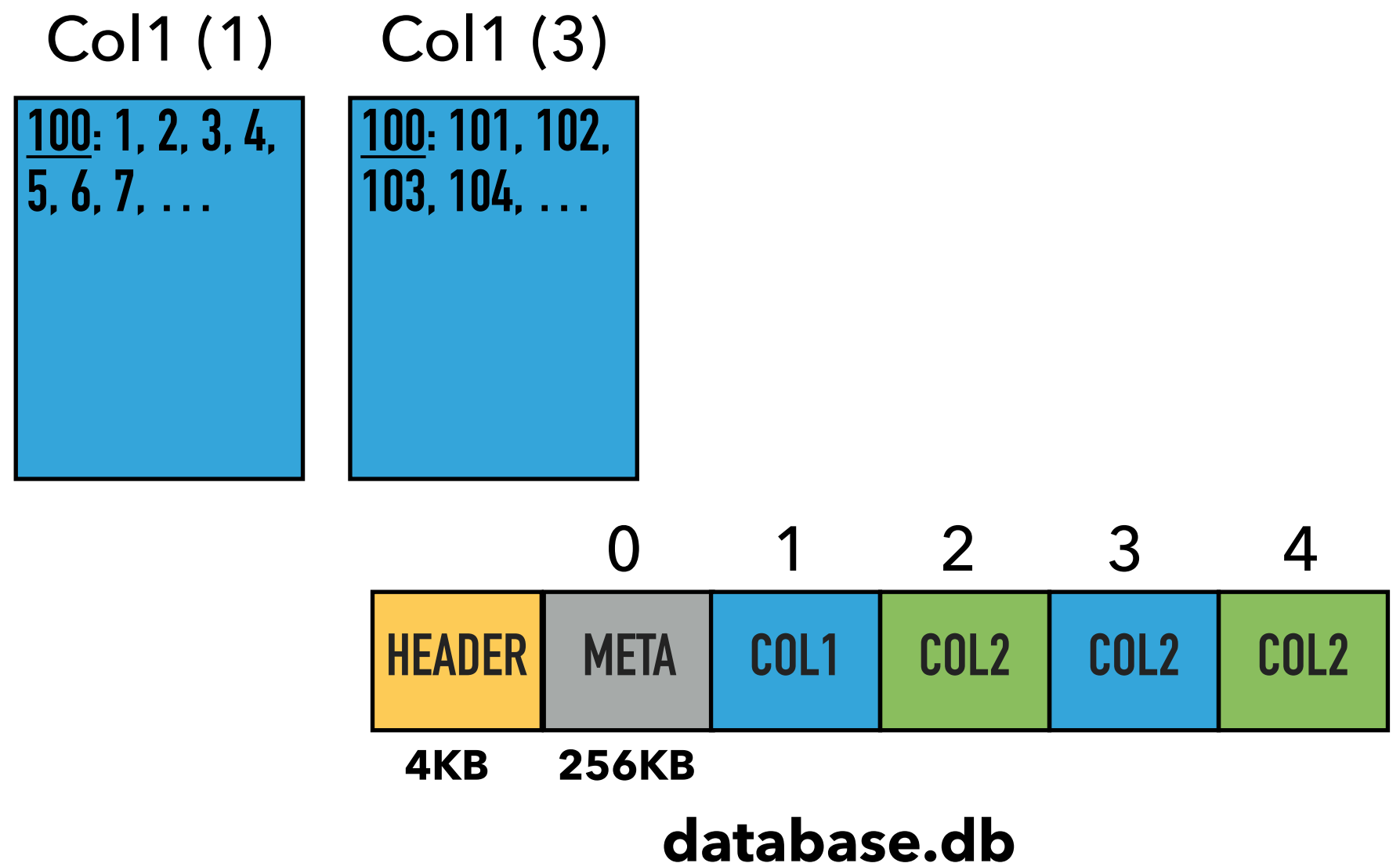
Meta (0)

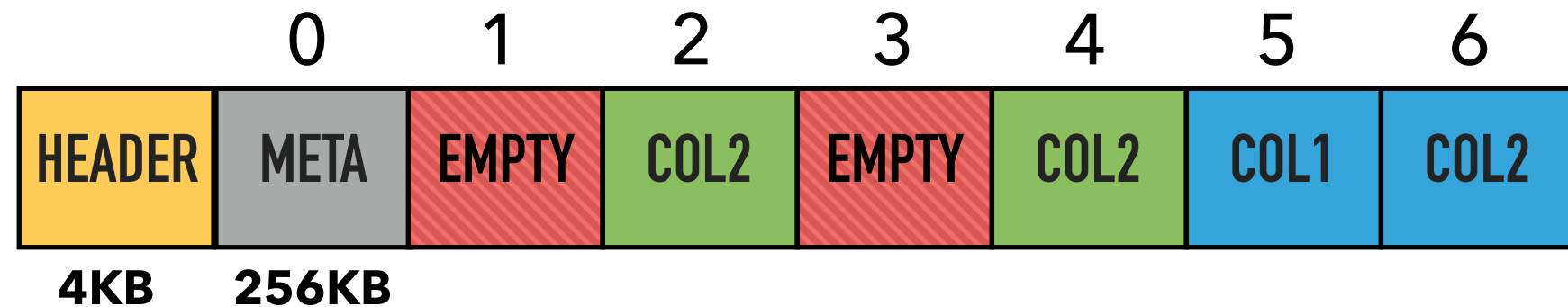
table
col1 INTEGER blocks: 1, 3
col2 VARCHAR blocks: 1, 4



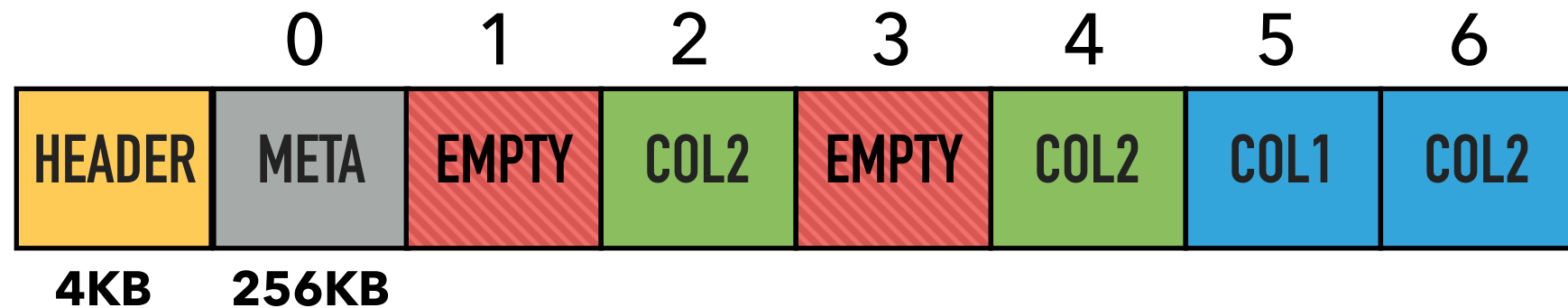
database.db

- ▶ Column data contains the data of that column
 - ▶ Tuple count + physical data
- ▶ Currently only uncompressed, compression WIP





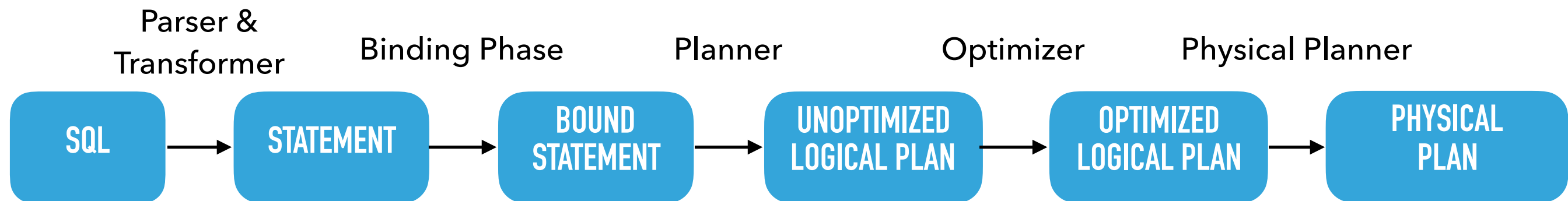
- ▶ Checkpointing rewrites **dirty blocks** first
- ▶ Header updated after writing successfully completes
- ▶ This ensures ACID properties
- ▶ **If writing fails during data write:**
 - ▶ Header is never updated
- ▶ Header is small enough to be updated atomically



- ▶ Checkpointing leaves **empty blocks**
 - ▶ Old data is not deleted until new blocks are written
- ▶ Next checkpoint will overwrite the empty blocks

Query Processing Pipeline

- DuckDB uses a typical pipeline for query processing



► Life of a Query

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

Parser &
Transformer

Binding Phase

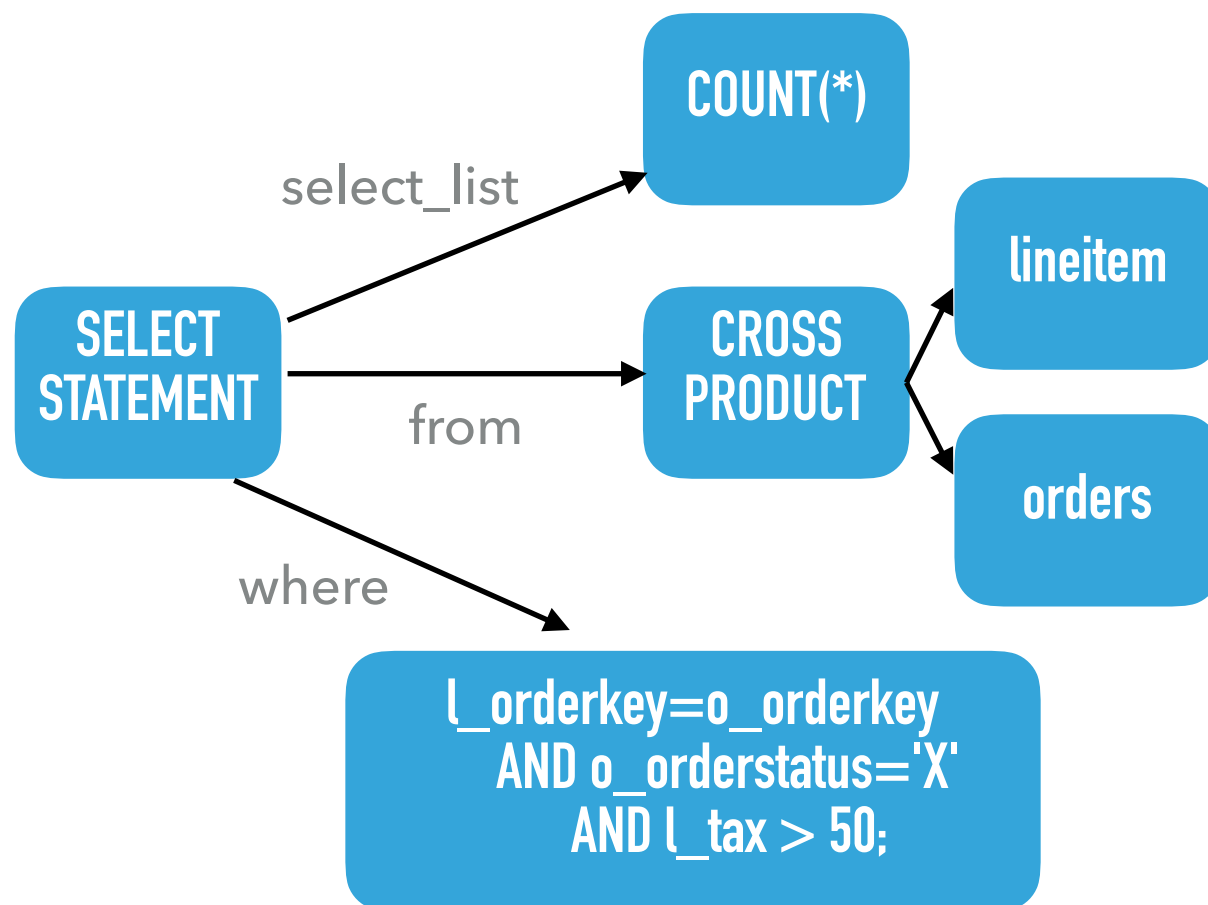
Planner

Optimizer

Physical Planner

SQL

STATEMENT

BOUND
STATEMENTUNOPTIMIZED
LOGICAL PLANOPTIMIZED
LOGICAL PLANPHYSICAL
PLAN

Tables and columns are strings
Nothing is resolved yet!

```
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

Catalog lookup is done here

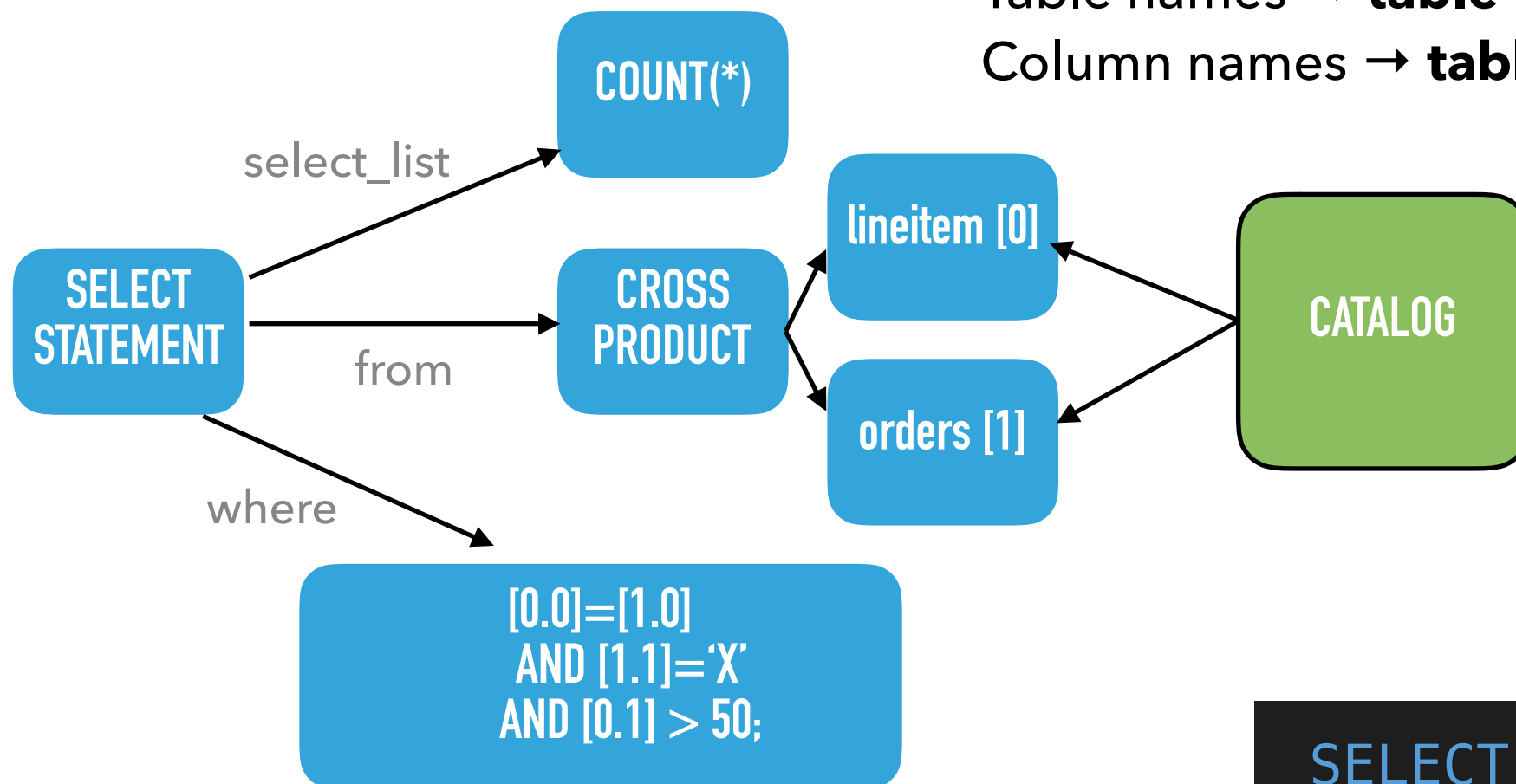
Table names → **table indexes**Column names → **table index + column index**

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='X'
      AND l_tax > 50;
  
```


Parser &
Transformer

Binding Phase

Planner

Optimizer

Physical Planner

SQL

STATEMENT

BOUND
STATEMENTUNOPTIMIZED
LOGICAL PLANOPTIMIZED
LOGICAL PLANPHYSICAL
PLAN

Logical Query Tree

AGGREGATE
COUNT(*)FILTER
[0.0]=[1.0]
AND [1.1]='X'
AND [0.1] > 50;CROSS_PRODUCTGET[lineitem][0]GET[orders][1]

Internally, nodes are called **Logical**...
e.g. LogicalAggregate, LogicalFilter

```
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

Logical Query Tree

AGGREGATE
COUNT(*)JOIN
[0.0]=[1.0]FILTER
[0.1]>50FILTER
[1.1]='X'GET[lineitem][0]GET[orders][1]

Optimizer **transforms** logical query tree
Created plan is equivalent but (hopefully) faster

```
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 Query Tree

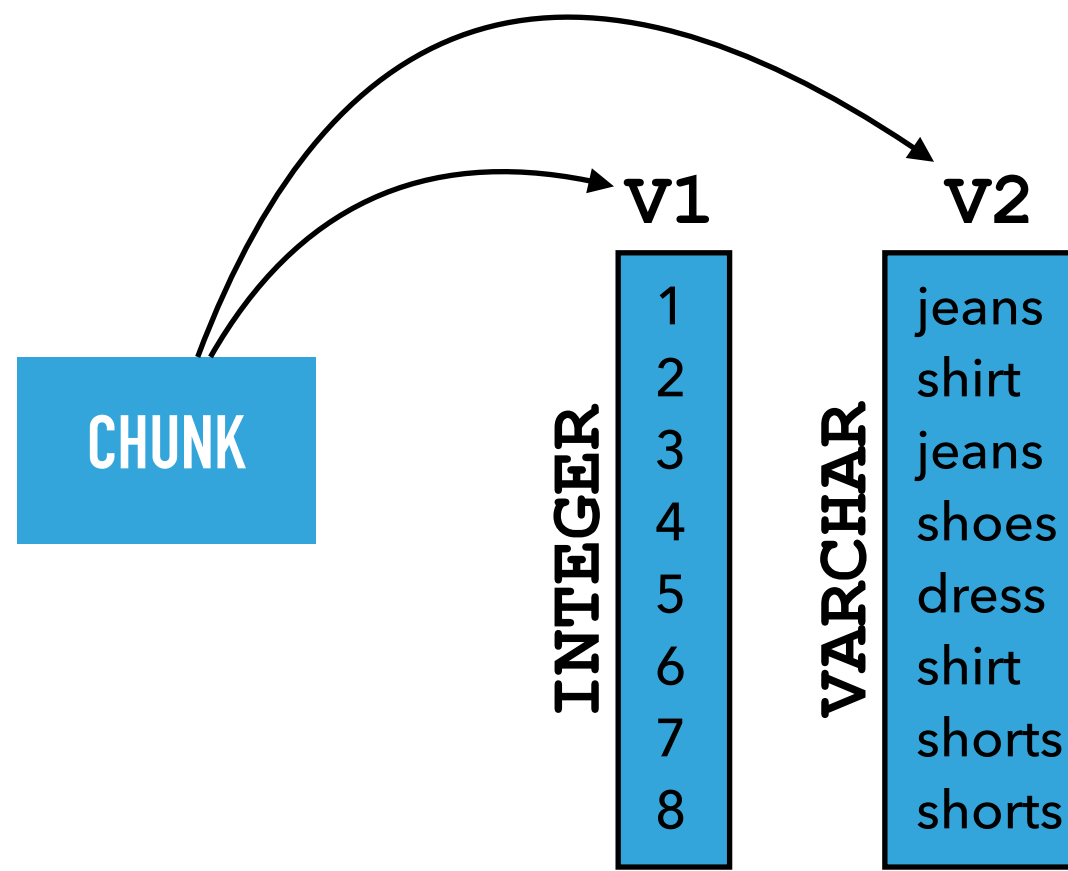
SIMPLE AGGREGATE
COUNT(*)HASH JOIN
L#0=R#0FILTER
#1>50SEQ_SCAN[lineitem]FILTER
#1='X'SEQ_SCAN[orders]

Physical planner creates **physical query tree**
Here implementations of operators are chosen

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

Query Execution

- ▶ 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)



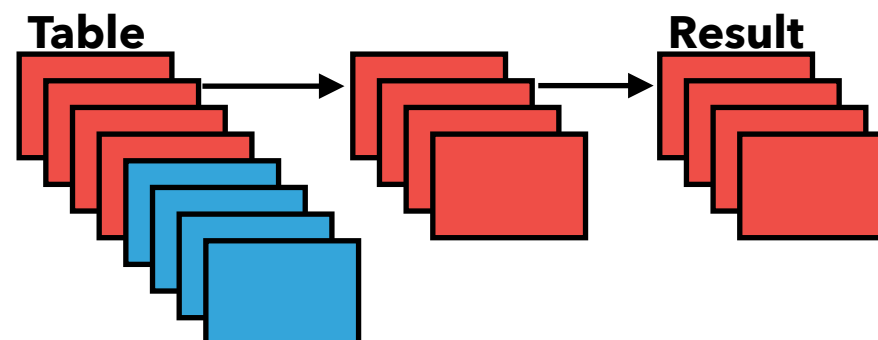
```
class Vector {  
public:  
→   typeId type;  
→   index_t count;  
→   data_ptr_t data;  
→   sel_t *sel_vector;  
→   nullmask_t nullmask;
```

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

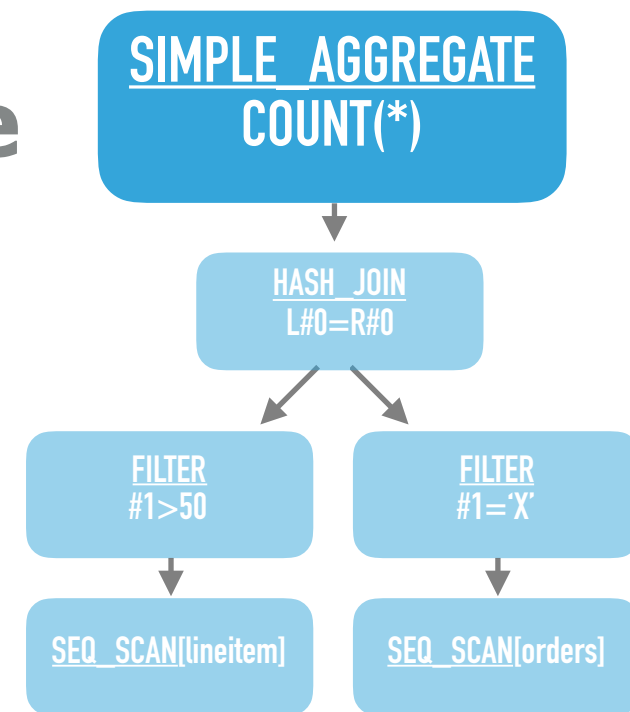
- ▶ **nullmask**: bitmap indicating which values are NULL
- ▶ **sel_vector**: optional selection vector indicating **which** values to use in the vector

- ▶ 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

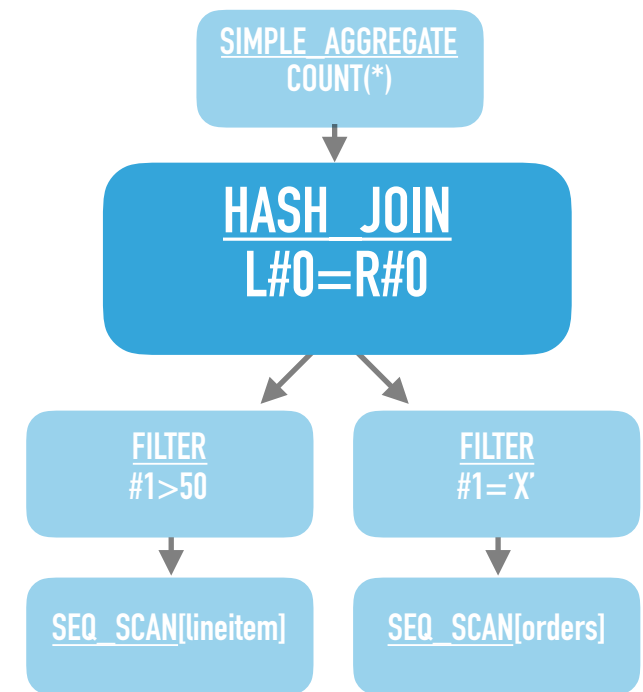
Vectorized Processing



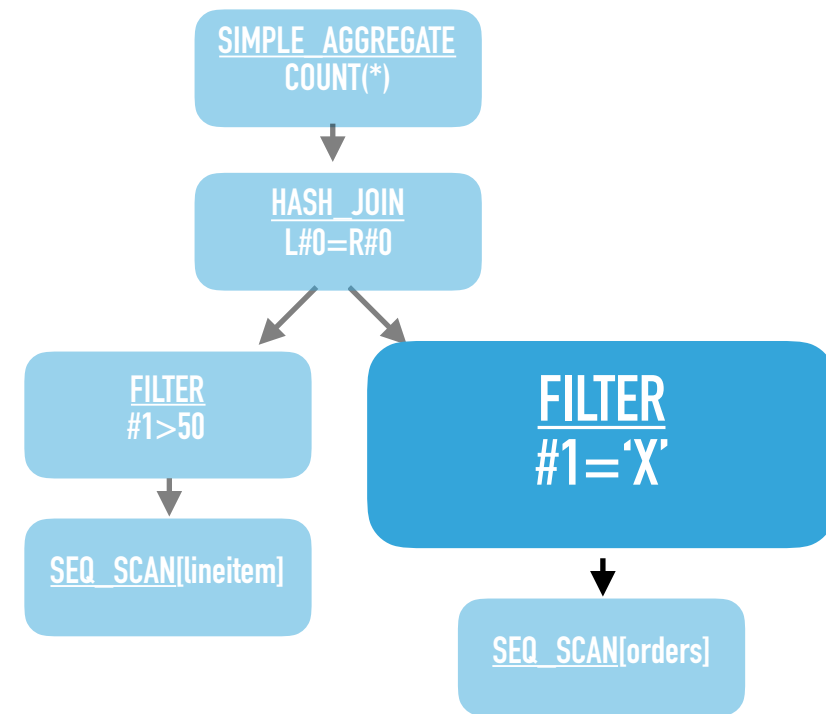
- ▶ Start with root node: **SimpleAggregate**
 - ▶ Aggregate without groups
- ▶ Immediately calls **GetChunk** on child



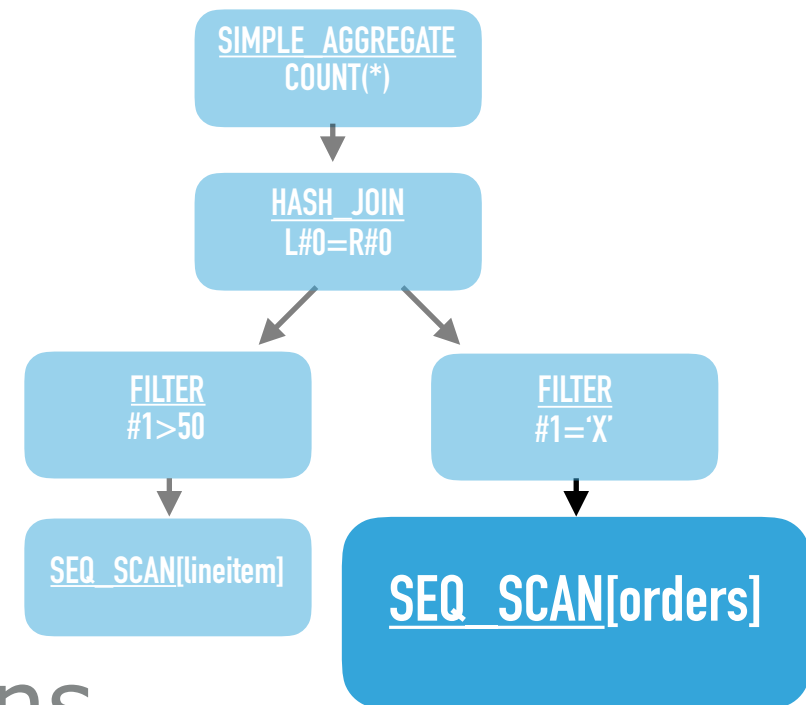
- ▶ Hash Join
- ▶ Start by **building HT**
- ▶ Call **GetChunk** on right node



- ▶ **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

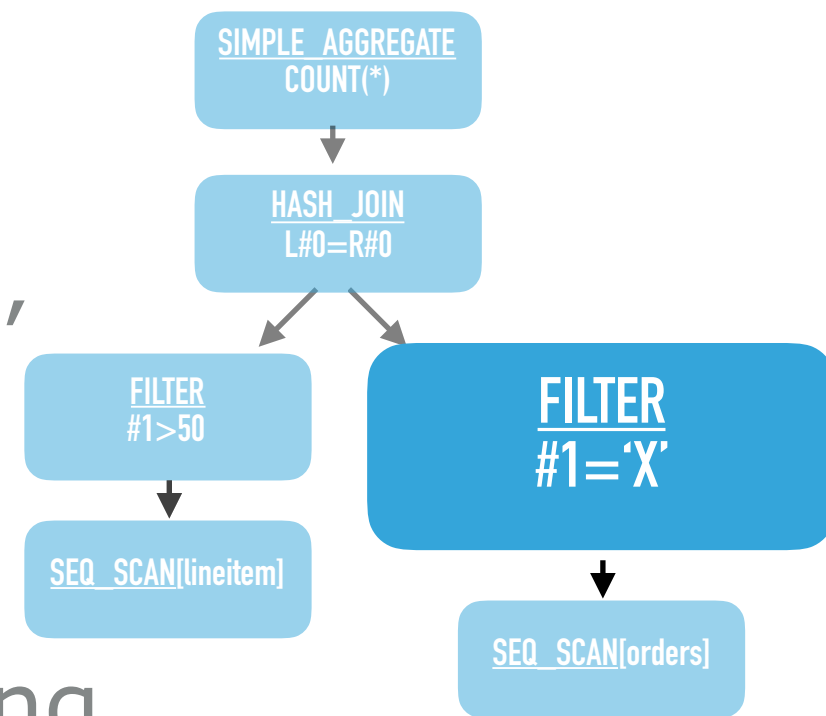
► **Filter**

► Now we can perform the filter $\#1 = 'X'$

► Only the second tuple passes

► **Selection vector** pointing to surviving tuple is created

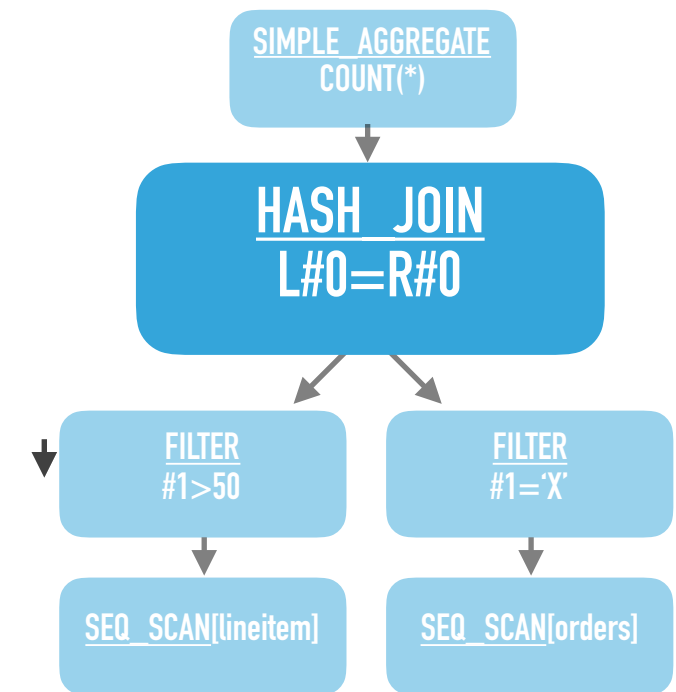
► Note that no data is copied or changed



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

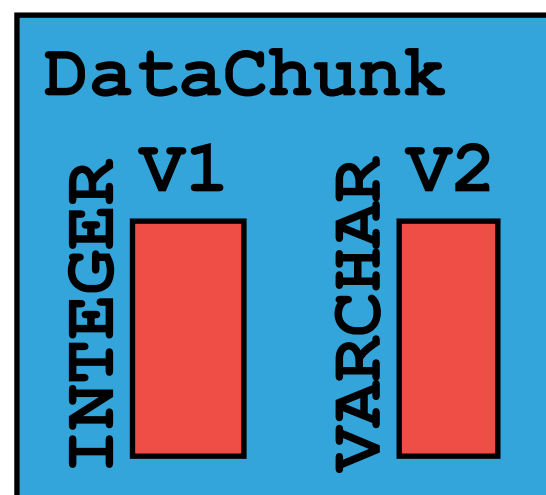
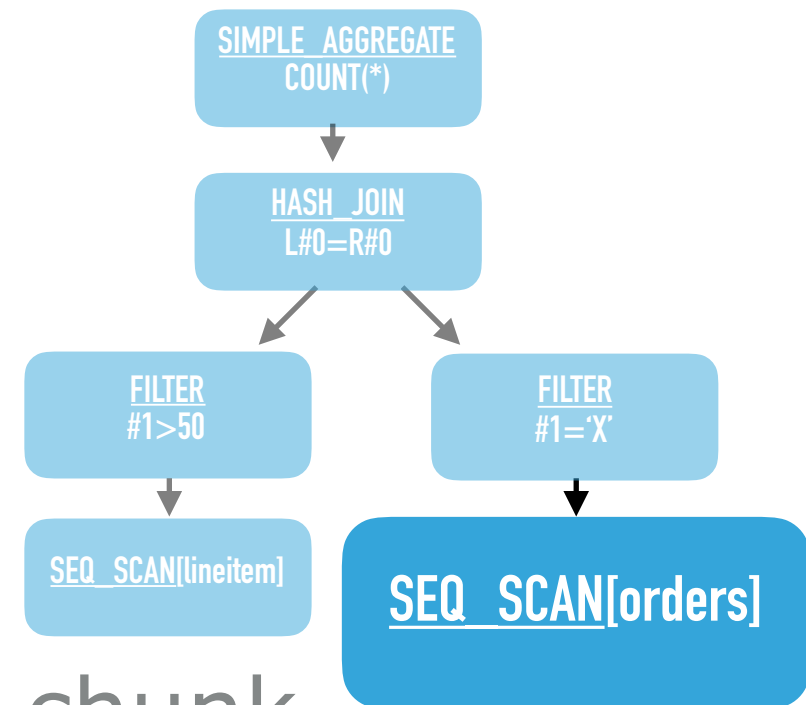
Keys	Payload	Next
2	X	0

► Sequential Scan

► The filter again calls `GetChunk`

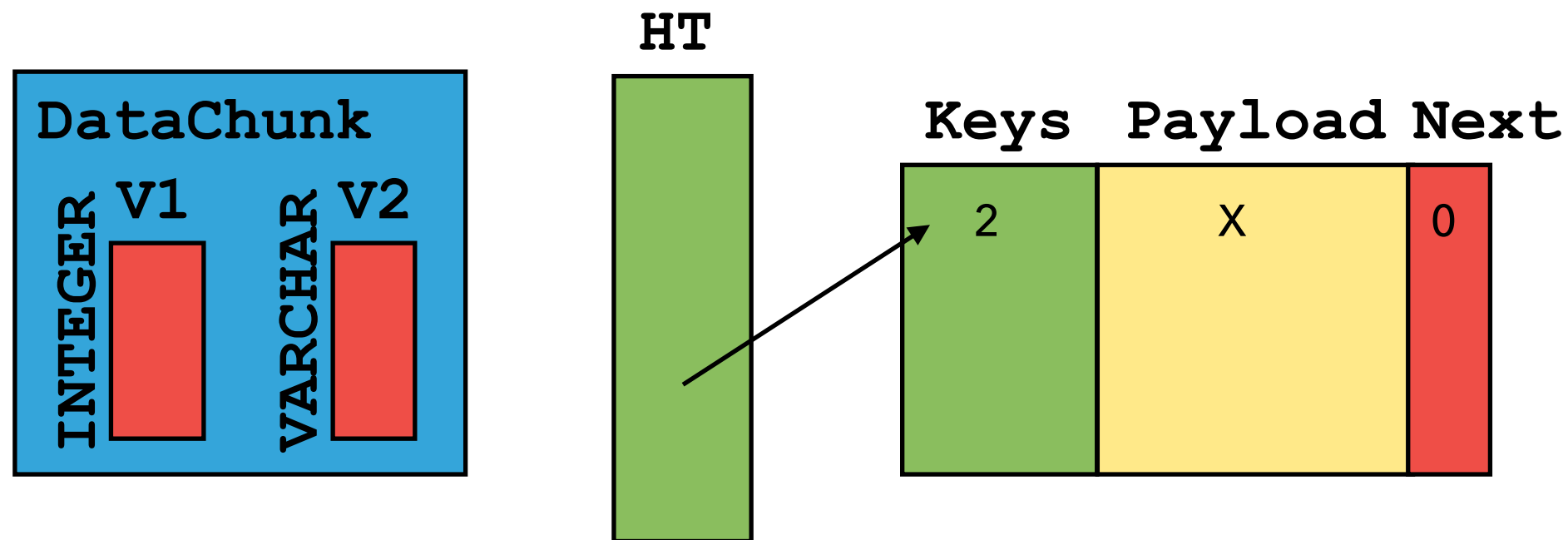
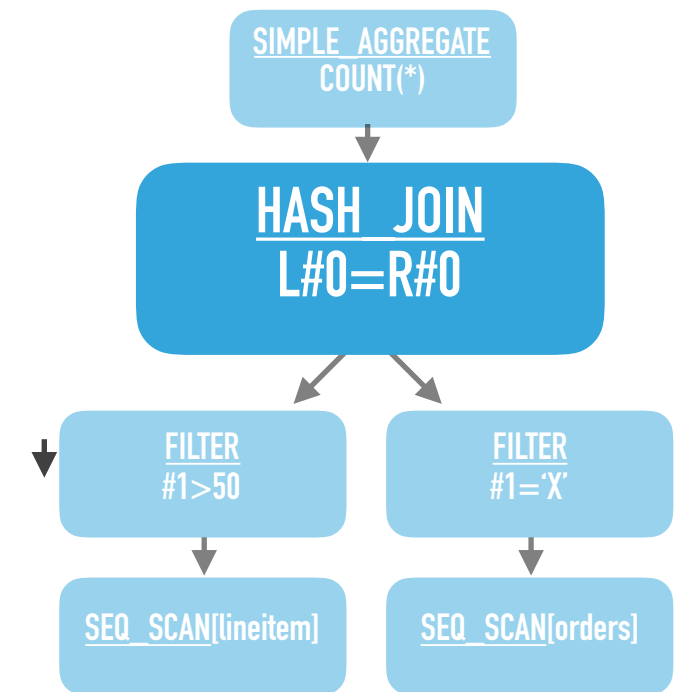
► Scan base table again:

► The scan is finished, return empty chunk

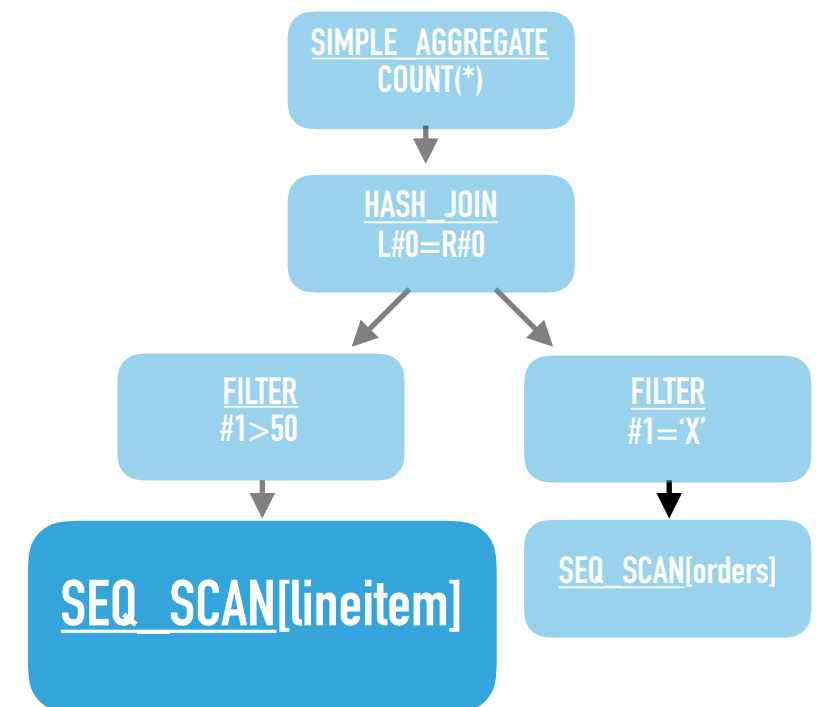


► Hash Join

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

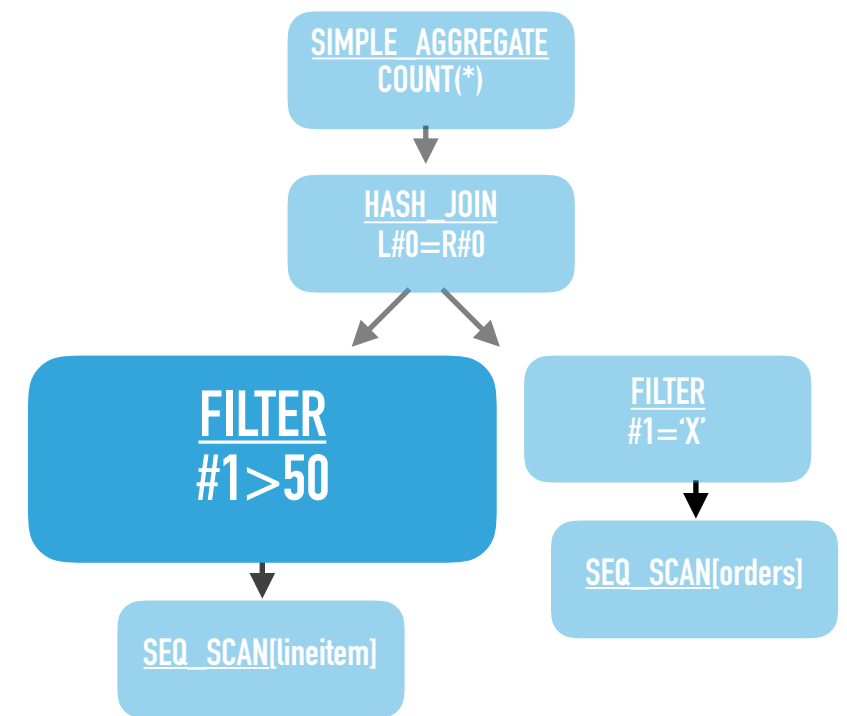


- ▶ **Sequential Scan**
- ▶ We arrive at scan on lineitem
- ▶ DataChunk with two columns
 - ▶ `l_orderkey` and `l_tax`



DataChunk			
	V1		V2
INTEGER	1	VARCHAR	80
	2		60
	2		20

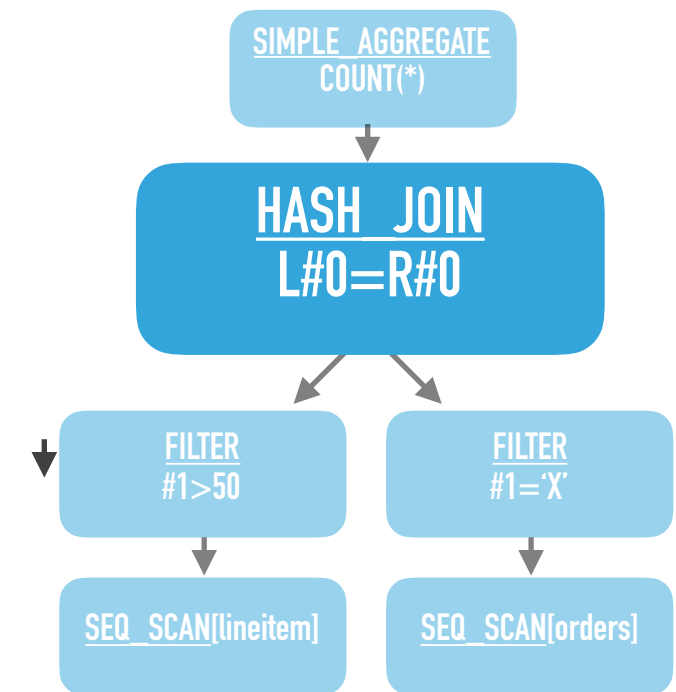
- ▶ **Filter**
- ▶ Performs the filter **#1>50**
- ▶ Again, add a selection vector



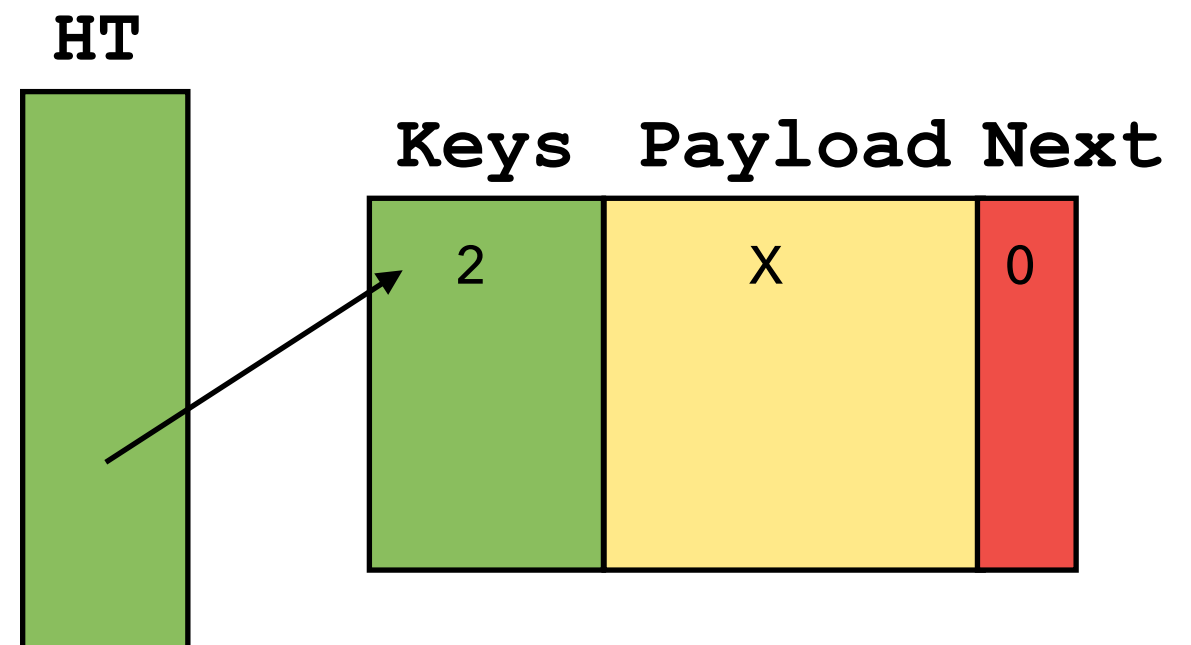
DataChunk				
	V1		V2	SEL
INTEGER	1	VARCHAR	80	0
	2		60	1
	2		20	

► Hash Join

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

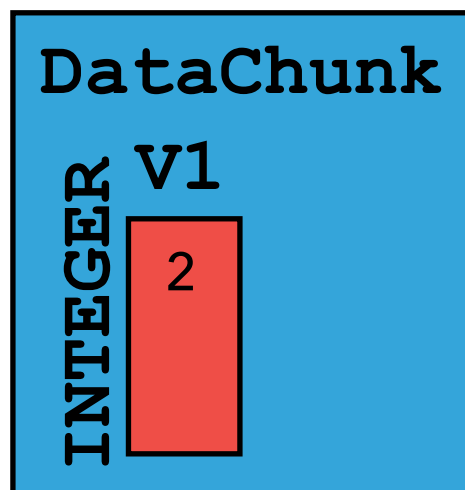
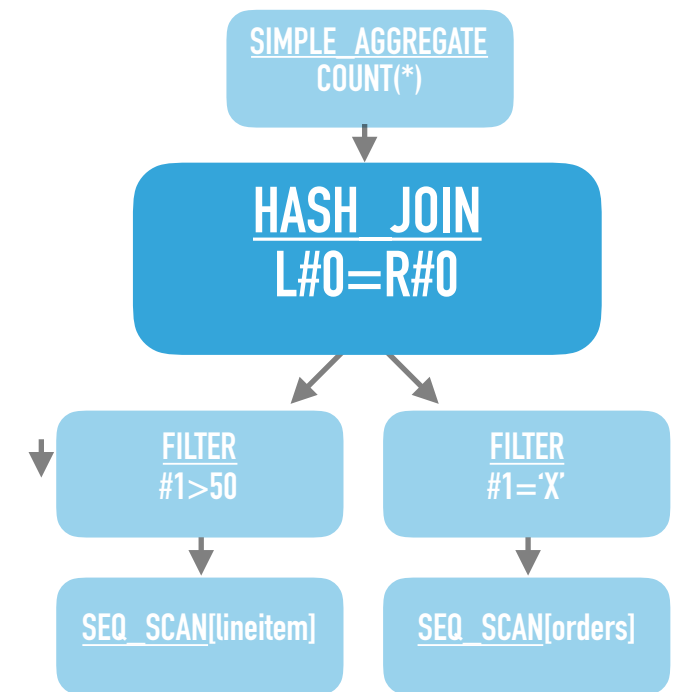


DataChunk				
	V1		V2	SEL
INTEGER	1	VARCHAR	80	0
	2		60	1
	2		20	

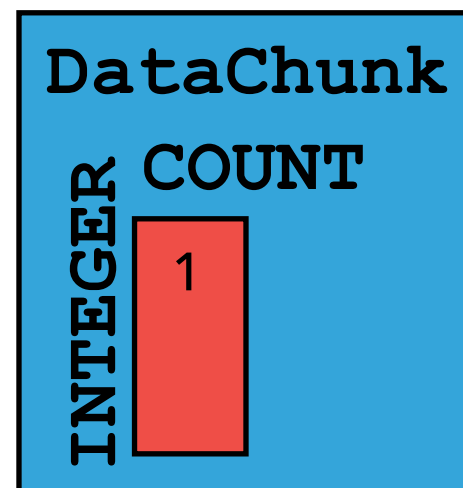
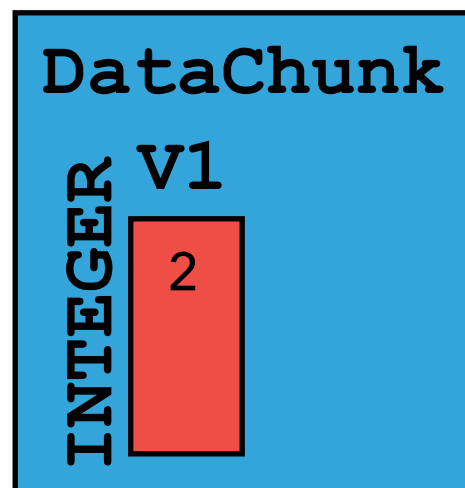
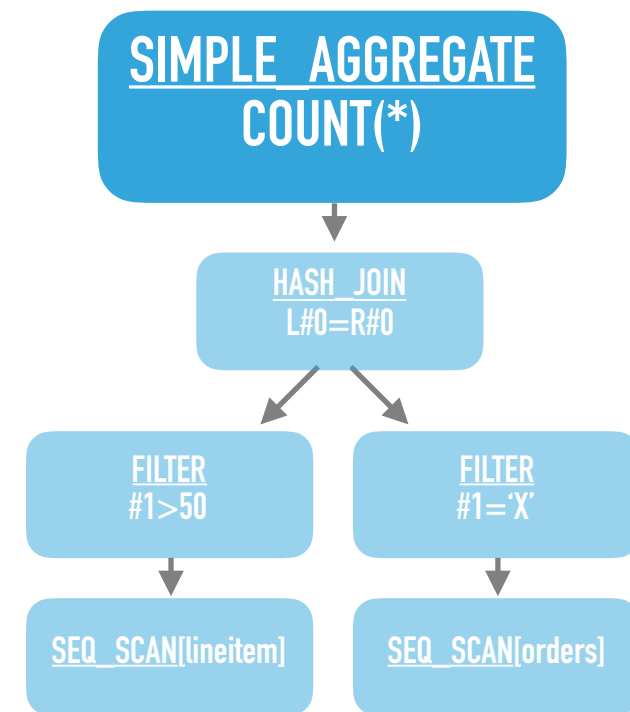


► Hash Join

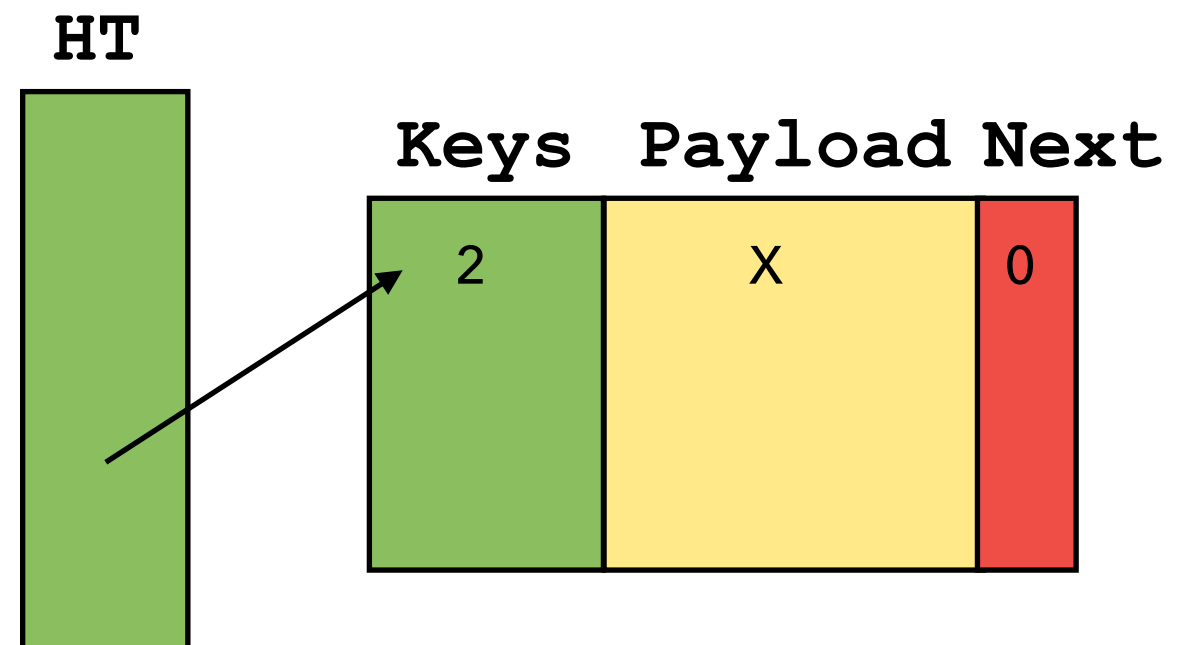
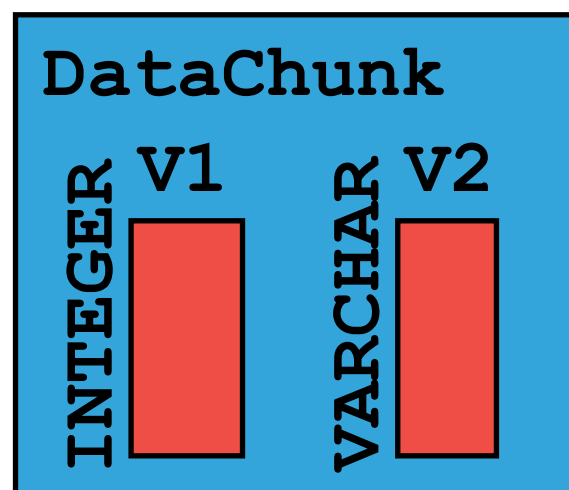
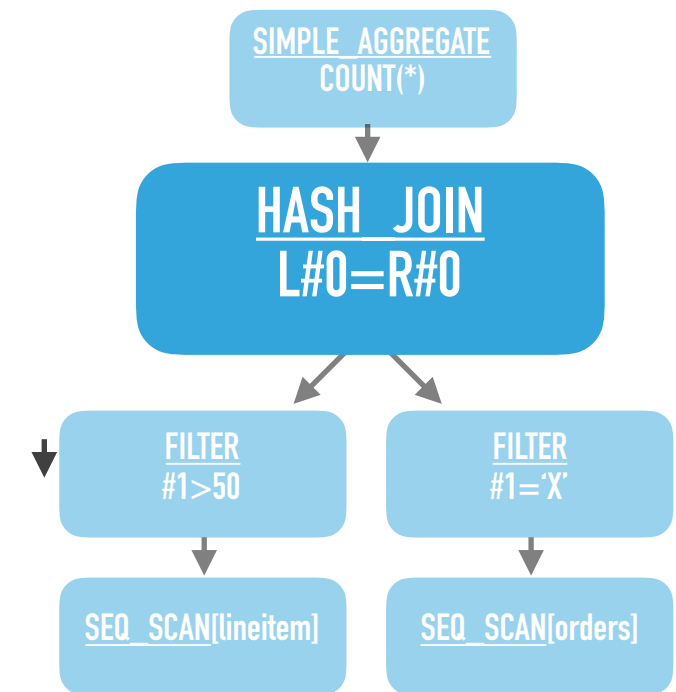
- We get one hit on our join!
- The hash join now produces the result
- We return this to the aggregate



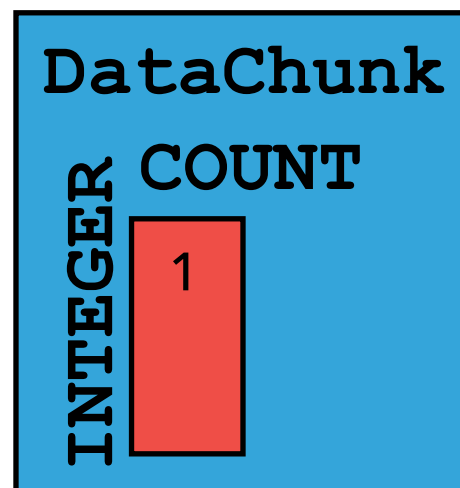
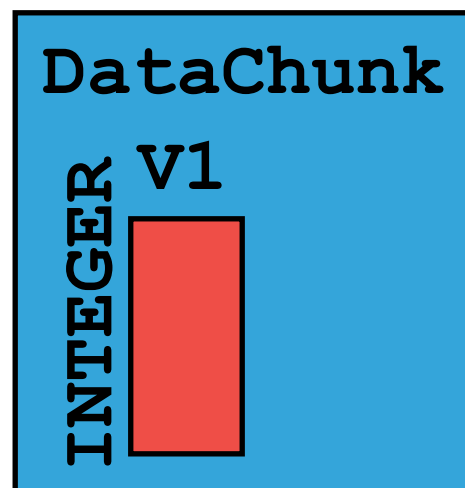
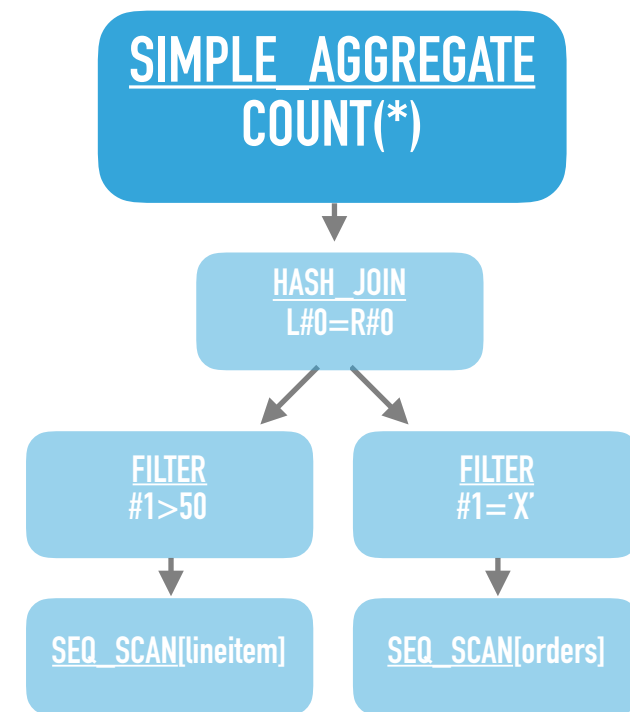
- ▶ The aggregate takes our input chunk
- ▶ Updates the aggregate
- ▶ Then fetches from the child again



- ▶ 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!



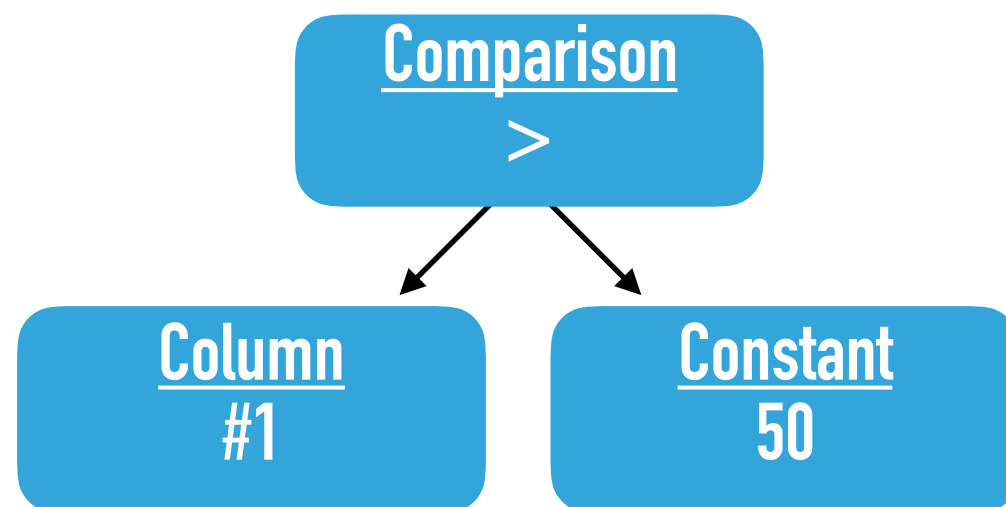
- ▶ Aggregate gets an empty chunk
- ▶ Returns the final result of our query



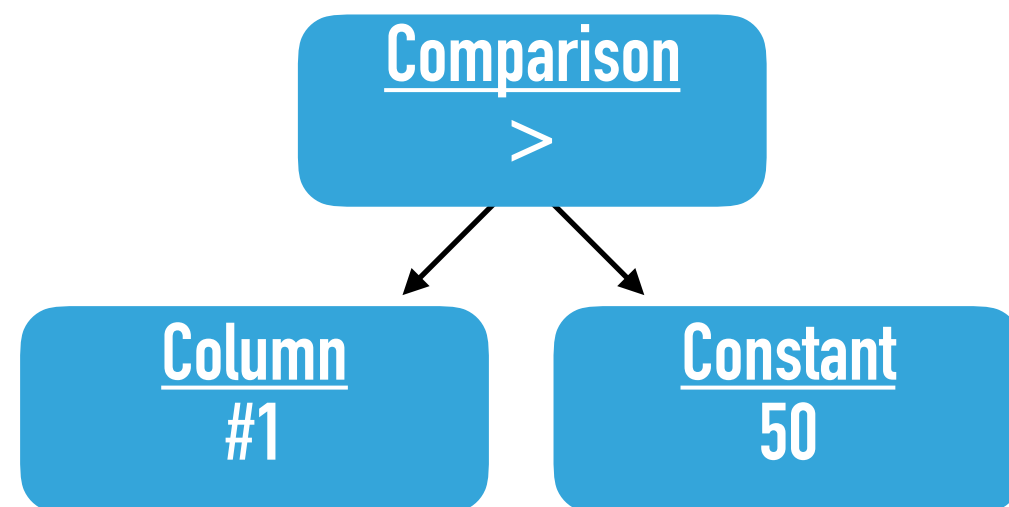
Expression Execution

- ▶ **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

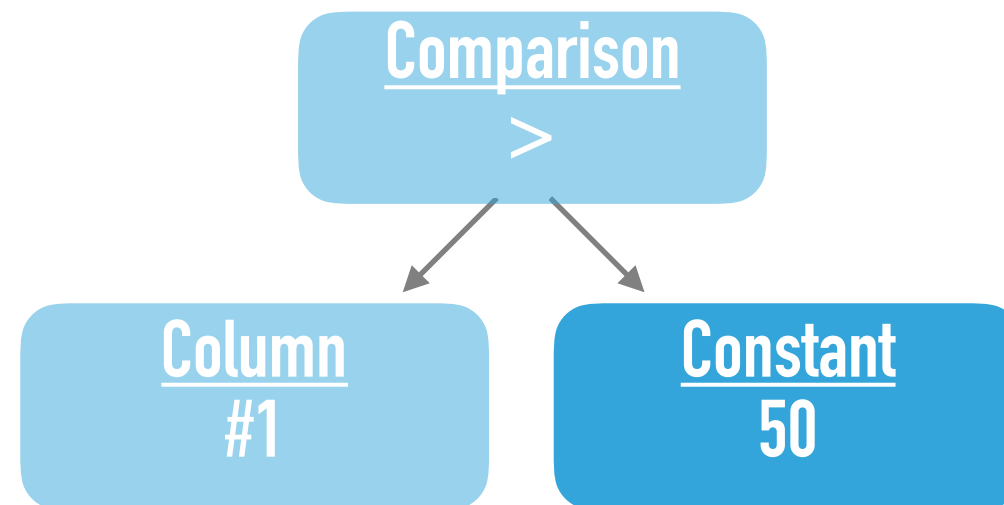
80
60
20

C1

80
60
20

Column reference
#1 fetches second
column from input

FILTER
#1 > 50



DataChunk			
INTEGER	V1	VARCHAR	V2
	1		80
	2		60
	2		20

C1

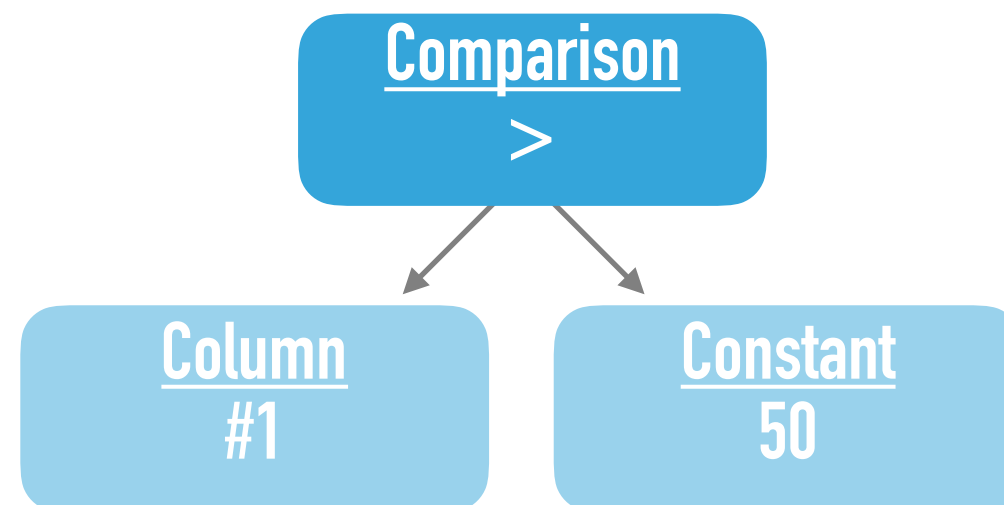
80
60
20

C2

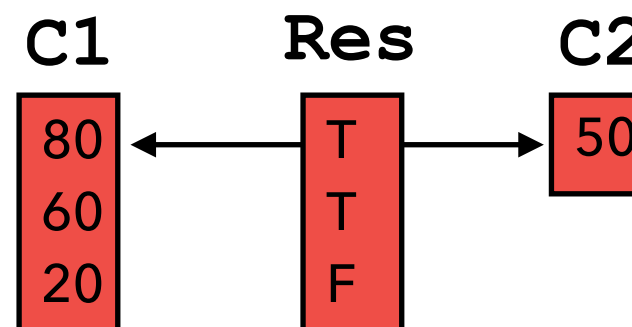
50

Constant is a
single value

FILTER
#1 > 50

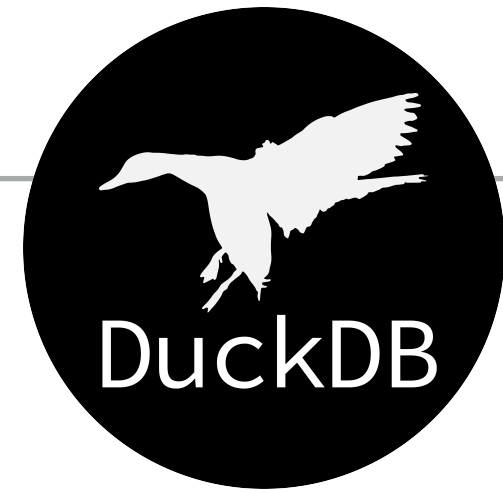


DataChunk			
	V1		V2
INTEGER	1	VARCHAR	80
	2		60
	2		20



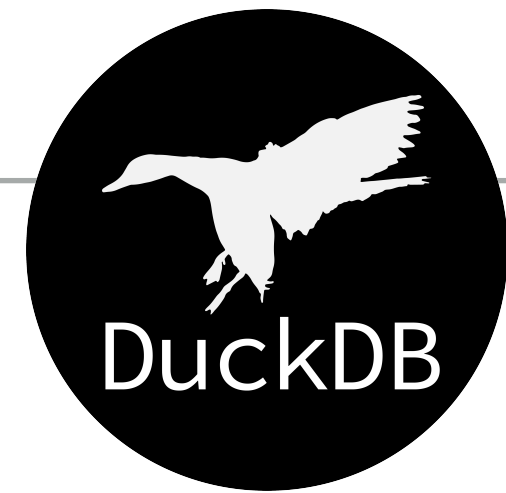
Comparison
runs and returns
matching tuples

Hands On

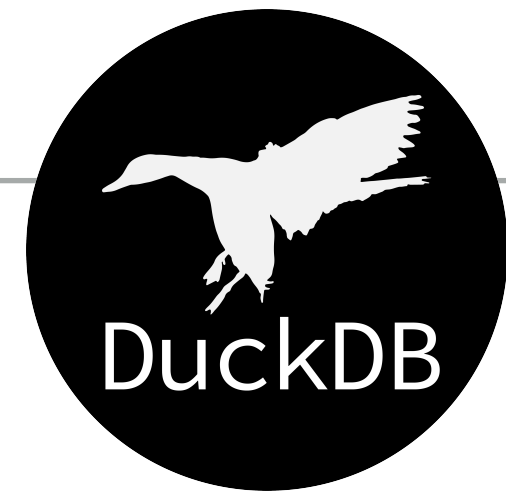


- ▶ **Assignment:** Implement a function in DuckDB
- ▶ Open issues for functions from other systems:
- ▶ <https://github.com/cwida/duckdb/issues/193>
- ▶ Implement one of those
 - ▶ For those that are successful, submit a PR!

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

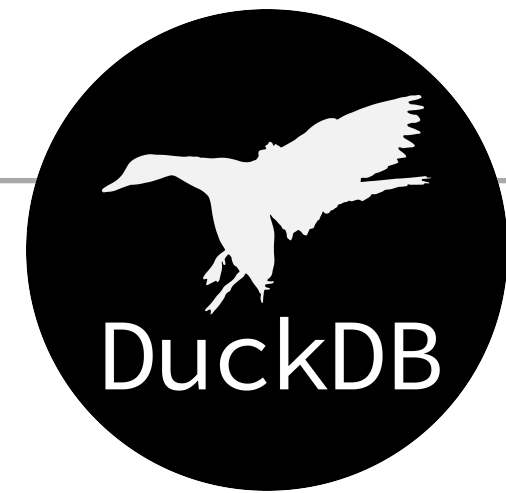


- ▶ Tests are in the `test` directory
 - ▶ We use the Catch framework for tests
- ▶ Tests look like this:

```
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]

▶ Remove spaces on right side of 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!