

Design and Implementation of DuckDB Internals

①

Welcome & Setup

April 7, 2026

Torsten Grust
Universität Tübingen, Germany

1 | Welcome!



Welcome to this course which is all about digging deep into the internals of **tabular database management systems** (DBMSs).

Our tour through the DBMS kernel will touch on

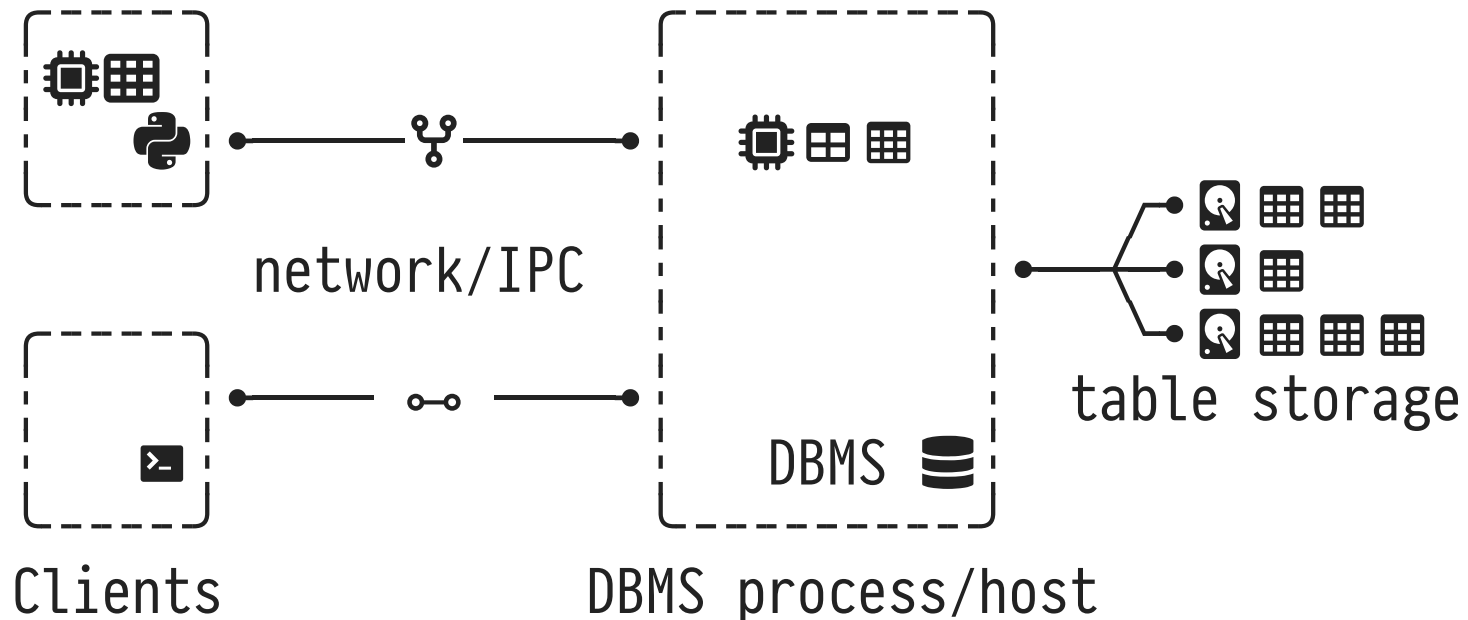
- the efficient representation of data on **secondary** (SSD/HDD 🗄️) and **primary storage** (RAM 🖨️),
- turning declarative SQL queries into efficient **flows of data**,
- a variety of interesting **data structures** 🗂️ for **sizable volumes** and associated **algorithms**,
- **modern CPUs** 🖨️ and how looping/branching code ↻↪ executes,
- various forms of **parallelism** 📄 on different levels (from single CPU instructions to threads), or
- ensuring **data integrity** under concurrent access or even if the host machine fails 🗨️ 🗑️.










Dissecting the Duck's Innards

This course will focus on **DuckDB** , a contemporary tabular DBMS built for high-speed SQL-based data analytics.

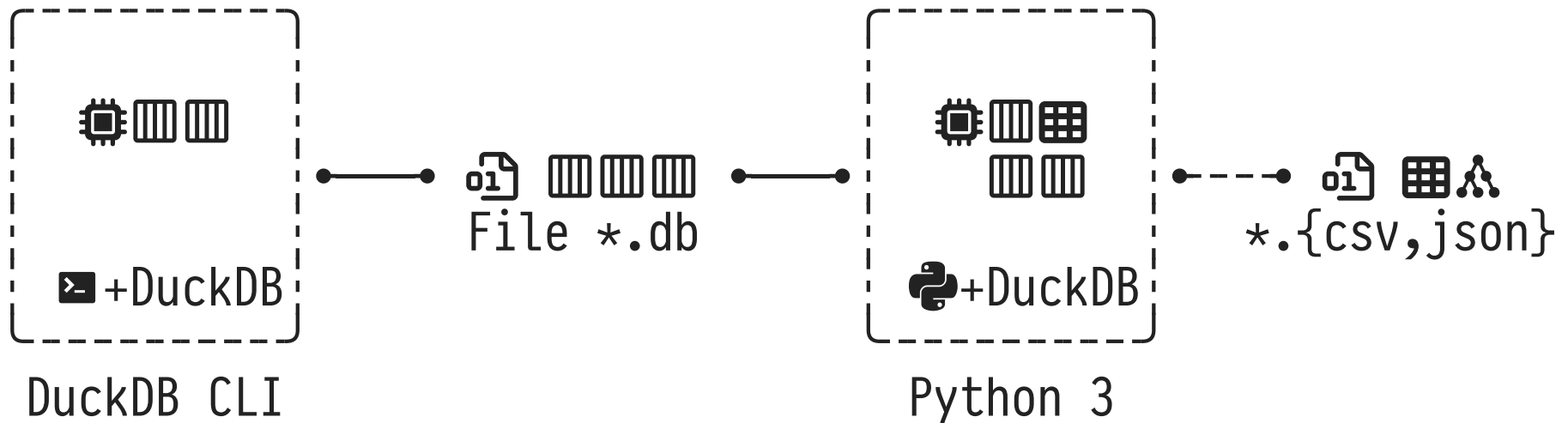
- DuckDB is young, first released in June 2019 (while most DBMSs—like PostgreSQL—originate from the 1980s).
- DuckDB development is moving fast. The system both adopts newest research results and builds on established DB wisdom.
- DuckDB is developed in the open (, MIT license). All code is available for inspection and tinkering.
- DuckDB provides a range of hooks to observe its internals.
- DuckDB comes with a CLI  and programming language APIs.
- DuckDB is easy to install and requires no maintenance.

Not DuckDB: Dedicated DBMS Process/Host Isolated From Users



- DBMS  controls family of disks holding table data , maps relevant table fragments into RAM buffer .
 - On-disk data organized in directories  or in raw blocks.
- Client processes  and DBMS process  isolated, connected via network  or inter-process communication .
- Data needs to be de-/serialized after/before wire transfer.
- Data structures in clients  are inaccessible by the DBMS.
- DBMS archetypes: PostgreSQL, MySQL, SQL Server™, Oracle®.

DuckDB: A Tabular DBMS Inside Your Own Process



- DuckDB kernel and client share a **single process** [] .
 - Python: `import duckdb`, C/C++: link with `libduckdb`.
- Table data resides in a **single database file** [] *.db, native DuckDB data formats [] in file and in RAM are similar.
- DuckDB sees in-process client data and can **directly read/write client data structures** [] using SQL (“zero copy”, replacement scans).
- DuckDB allocates sizable RAM buffers (but can use disk [] for temporary storage if required).

2 | This Course (*Dissecting the Duck's Innards*, short: *DiDi*)

- We will focus on DuckDB as a **tabular** SQL-based DBMS 🗃️. There are other kinds of DBMSs (for graphs, key/value pairs, vectors, ...), but we will not discuss those here.
- We will get our hands dirty using **DuckDB** 🐥 and its extensive **SQL** dialect. Lots of SQL will be read and written.
- Whenever possible we try to observe DuckDB under load or use hooks to inspect its operation while SQL queries are processed.
- We thus assume basic familiarity with the tabular data model and SQL, e.g., as discussed in *Tabular Database Systems (TaDa)*.
- We will draw data and queries from a variety of sources and have **fun** along the way! 😎

Torsten Grust?





Time Frame	Affiliation/Position
1989-1994	Diploma in Computer Science, TU Clausthal
1994-1999	Promotion (PhD), U Konstanz
2000	<i>Visiting Researcher</i> , IBM (USA)
2000-2004	Habilitation, U Konstanz
2004-2005	Professor Database Systems, TU Clausthal
2005-2008	Professor Database Systems, TU München
since 2008	Professor Database Systems, U Tübingen

- Web: <https://db.cs.uni-tuebingen.de/grust>
- Bluesky 🦋: [@tegggy.org](https://bsky.app/profile/tegggy.org)
- E-mail: torsten.grust@uni-tuebingen.de
- Feel free to reach out with criticism, bug reports, suggestions for improvement, pats on the back, or simply to say “Hi!” 🙌👋

Slides and Further Lecture Material



These **slides** (PDF), **code fragments** (SQL, Python, C), and **sample data** will be uploaded to a GitHub  repository:

github.com/DBatUTuebingen/DiDi 

- Slides point to relevant code files or extra material using tags like  #001:
 - Refers to a file named `001-*` on the GitHub repository (e.g., `001-sum-quantity.awk`).
- **NB.** Code and extra material provide essential content (e.g., details on SQL-based experiments or sample data).
 -  +  = : Only slides + code provide a complete picture.

Material

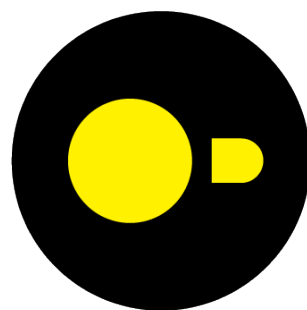
This course is *not* based on a single textbook. Rather, we build on

- a variety of scientific papers,
- textbook excerpts (few),
- the DuckDB  documentation at <https://duckdb.org/docs/>,
- Python/C/C++ code snippets (our own and from inside the ) ,
- blog posts from a range of authors,
- SQL references/standards,
- experience, and best practices.

There is a plethora of books on tabular DBMSs (both usage and internals), sample SQL snippets (experiments, benchmarks, and idioms), or performance tweaks. If we will use such sources, we will provide pointers.

Get Your Hands Dirty: Install DuckDB!

The tabular DBMS **DuckDB** will be the primary tool in this course:



DuckDB

<https://duckdb.org>, version 1.5 (March 2026: 1.5.0)

- Implements an extensive SQL dialect, is highly performant, open to contributions, and generally awesome.
- Straightforward to install and use on macOS 🍏, Windows 🪟, Linux 🐧 (x86 + ARM).

No DuckDB CLI (📄 SQL prompt/REPL) on iOS or Android.²

² Run the DuckDB CLI in the web browser: <https://shell.duckdb.org>. Suffices for quick SQL experiments.

3 | DuckDB? 🦆?

In case you were wondering:

DuckDB has been named after *Wilbur, the Duck*, which has been living as a pet with Hannes Mühleisen³—co-inventor of DuckDB with Mark Raasveldt—on Hannes' houseboat in Amsterdam.

Hannes (CEO) and Mark (CTO) run **DuckDB Labs**, a company that provides support and consultancy services around DuckDB. The labs are located in Amsterdam, The Netherlands.



Hannes and Wilbur (© Hannes Mühleisen)

³ Hannes originally is from the Stuttgart area. Back then his car had the license plate **S:QL1337**.

Design and Implementation of DuckDB Internals

②




The Query Performance Spectrum

April 7, 2026

Torsten Grust
Universität Tübingen, Germany


1 | DBMSs Exploit Modern Computer Architecture¹

The internals of DBMSs are carefully engineered to exploit the performance features of modern computer architecture:

- **CPUs**  (and their multi-threading capabilities),
- **main memory** (DRAM ) and its hierarchy of caches, and
- **secondary memory** (mass storage on SSDs or rotating disks ).

Since database queries typically process millions of rows, the effect of even the tiniest performance tweaks/tricks played in the innermost loops of DBMS routines multiply.

Goal: Understand the performance spectrum for a simple “query.”

quick one-liner
shell script  hand-written
C program

¹ This chapter adapts and expands on a discussion found in Thomas Neumann's lecture [“Foundations in Data Engineering” \(TUM\)](#) .

A Simple Benchmark Query

1. Read the CSV file for TPC-H table `lineitem` (scale factor `SF = 1`: 6+ million rows × 16 columns ≈ 720 MB of data) and
2. sum the `quantity` integer values in the 5th column:

`lineitem.csv`

```
1|155190|7706|1|17|21168.23|0.04|0.02|N|0|1996-03-13|...NL
1|67310|7311|2|36|45983.16|0.09|0.06|N|0|1996-04-12|...NL
1|63700|3701|3|8|13309.60|0.10|0.02|N|0|1996-01-29|...NL
⋮
[6+ million more rows]
```

- Real TPC-H benchmark data and queries are more complex but this suffices to demonstrate the effect of code optimizations.
- We will implement the query in awk, Python, C, and SQL.

2 : Performance Limits

What is the fastest query time we can hope for in principle?

- Torsten's current computer (🍏 MacBook Pro M2 Max, 2023):

Memory (📄 Primary/📄 Secondary)	Read Bandwidth	🕒 Query Time
(Ethernet)	2.5 GB/s	0.28s
📄 External USB-C SSD (2 TB)	800 MB/s	0.90s
📄 NVMe SSD (2 TB)	5 GB/s	0.14s
📄 DRAM (64 GB)	21 GB/s	0.03s

- **NB.**
 - Column **Query Time** based on I/O speed, ignores CPU cost (less significant for secondary mem, very significant for DRAM).
 - \Rightarrow We will *not* reach these limits. Let us try to get close.
- Understand how DuckDB achieves 0.002s for our query. 🐧☰

3 | Sum of Quantities ① — awk

- `awk`: interpreted text processing language popular on UNIX™.
 - Read input line by line, match each line against (regular) patterns in order, on a match invoke action `{...}` on line.

BEGIN	{ FS = " " }		delimiter in CSV is	#001
	sum = 0 }		match first line, reset sum	
	{ sum = sum + \$5 }		match any line, sum 5th column	
END	{ print sum }		match last line, output sum	

- Invoke the `awk` script, measure elapsed wall-clock time (s) 🕒:

```
$ time ./001-sum-quantity.awk lineitem.csv
153078795
      1.58 real          1.43 user          0.14 sys
```

Sum of Quantities ① — awk

- 🕒 Query time on Torsten's computer: $\approx 1.6\text{s}$:


Output of time	Measurement
real	elapsed wall-clock time 🕒 ($\approx \text{user} + \text{sys} + \Delta$)
user	time spent in application/library code
sys	time used by OS (system calls)

- The interpreted awk script cannot even keep up with secondary memory (SSD) read bandwidth:²
 - awk processes the CSV file with a throughput of 471 MB/s.
 - awk is **CPU-bound** for this query.
 - \Rightarrow The OS file system cache in DRAM does not help.

² Execution speed of awk variants vary greatly. We are using GNU awk ([gawk](#)) here. macOS awk is about 10 times slower for our benchmark query.

4 | Sum of Quantities ② — Python

- **Python**: established scripting/programming language, mainly follows an imperative paradigm.
 - Translates to bytecode, then interprets.

<pre>sum = 0</pre>		reset sum	 #002
<pre>with open(sys.argv[1], 'r') as file:</pre>		open file (reading)	
<pre> for line in file:</pre>		read line by line	
<pre> sum = sum + int(line.split(' ', 5)[4])</pre>		extract 5th col,	
		cast to int, add	
<pre>print(sum)</pre>		output	

- 🕒 Query time on Torsten's computer: ≈ 2.75 s.
 - Python processes the CSV file with a throughput of 275 MB/s.
 - Python is **CPU-bound** for this query.

5 | Sum of Quantities ③ — C

- Switch to compiled programming language **C**. Start out with a direct transcription of the Python code:
 - Read CSV file line by line using `getline(3)`.³
 - Use `strchr(3)` to search for delimiter '|' in line (4×).
 - Convert string (up to next '|') into integer using `atoi(3)`.

```
while (getline(&line, &linecap, file) > 0) {  
    char* delim = line;  
  
    for (column = 1; column < 5; column++) {  
        delim = strchr(delim, '|');  
        delim++;  
    }  
  
    sum = sum + atoi(delim);  
}
```

 #003

³ The (3) suffix in `getline(3)` refers to Section 3 of the UNIX™ manual which describes the function in the C Standard Library.

Sum of Quantities ③ — C

- 🕒 Query time on Torsten's computer: $\approx 0.5s$.
 - C processes the CSV file with a throughput of 1.5 GB/s.
 - Yet, C still is **CPU-bound** for this query. Getting closer to SSD read bandwidth, though.
- But where does time go?
 - **Profile** the running program, identify code portions consuming the most CPU time (UNIX™: [perf](#), macOS: [Instruments](#)).

C library routine	% of CPU time
(all other C code)	(16%)
getdelim (\equiv getline)	58%
atoi	14%
memchr (\equiv strchr)	12%

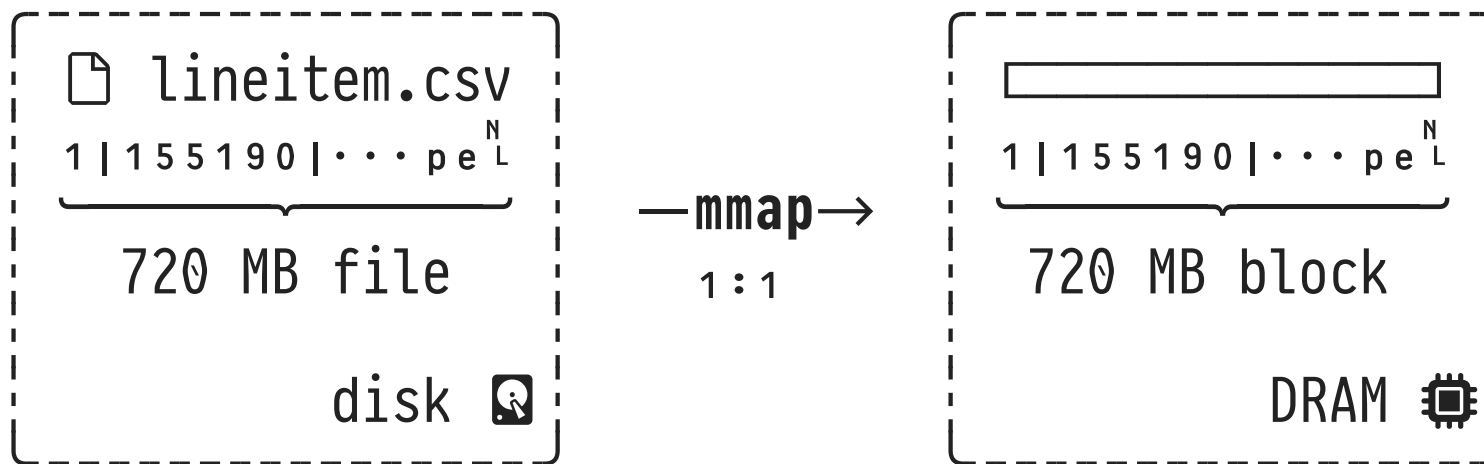
- Reading the CSV file line by line is too slow.
 - 💡 Read entire file at once, impose line structure ourselves.

6 : Sum of Quantities ④ — C with mmap(2)

Aim to *read the CSV file* into DRAM using a *single OS system call*:

- `mmap(2)` returns a pointer `data` to a contiguous block of memory that holds the contents of an entire disk file:

```
data = (char*)mmap(NULL, size, PROT_READ, MAP_SHARED, file, 0);
```



- If file size exceeds available DRAM, OS pages in file contents on demand.

Sum of Quantities ④ — C with mmap(2)

- Cannot use `strchr()` to find '|' (next column) ❶.
- No `getline()`: need to locate '\n' (❷) on our own now:

```
column = 1;

while (data < end) {
    switch (*data) {
        case '|': column++; break;
        case '\n': ...
    }
    data++;

    if (column < 5)
        continue;

    sum = sum + atoi(data);

    column = 1;
    while (*data++ != '\n');
}
```

start in column 1

 #004

scan memory block, byte by byte

❶ found |: next column begins
error: line has too few columns

proceed through memory block

reached column 5 already?
no, keep scanning

convert to int (up to |, '\n'), add

next line starts with column 1
❷ skip rest of line until '\n'

Sum of Quantities ④ — C with `mmap(2)`

- 🕒 Query time on Torsten's computer:
 - Once the OS caches the file in DRAM, `mmap()` directly maps the file system cache into the program's address space.

OS file system cache	Query time 🕒	Throughput
cold	1.6s	471 MB/s
warm	0.42s	1.8 GB/s


- NB.** The C program's profile has changed:

C code fragment/function	% of CPU time
<i>(all other C code)</i>	(25.5%)
<code>atoi</code>	21.4%
<code>while (*data++ ≠ '\n')</code>	👎 53.1%

- Search for `\n` 📄 dominates. 💡 Use `strchr(data, '\n')` instead:
 - 🕒 Query time: 0.27s (throughput 2.8 GB/s). 👍
 - How can `strchr()` be so efficient?


7 : Sum of Quantities ⑤ — C with mmap(2) + Block-Wise '\n' Search

Avoid byte-wise search for '\n'. Modern CPUs operate on 64-bit words.

-  Load **8 bytes (64 bits) at a time**, search for '\n' ('\n' = 0x0a) in this block. Advance pointer `data` in strides of 8 bytes.

```

/* HAS_NL(x): find '\n' in 64 bit-wide character block x */
#define HAS_ZERO(x) (((x) - 0x0101010101010101) &
                    ~(x) & 0x8080808080808080)
#define HAS_NL(x)   (HAS_ZERO(x ^ 0x0a0a0a0a0a0a0a0a))

HAS_NL(0x0a42440a6b637544) ≡ "Duck\nDB\n" reversed 
=      0x8000000800000000      (ARM CPUs: little endian)
      ↑      ↑
high bit set: found '\n' at offsets 4 and 7

```

- How do C macros `HAS_ZERO()` and `HAS_NL()` work?⁴

 #005

⁴ See the [Stanford Bit Twiddling Hacks](#)  (section “Determine if a word has a byte equal to n”) for a discussion of these C macros.

Sum of Quantities ⑤ — C with mmap(2) + Block-Wise % Search

```
while (data < end) {  
    :  
    sum = sum + atoi(data);  
  
    column = 1;  
  
    block = (uint64_t*)data;  
    while (!(nl = HAS_NL(*block)))  
        block++; /* advance by one 8-byte-block (64 bits) */  
  
    data = (char*)block;  
    if (nl & 0x000000000000000080ULL) { data += 1; continue; }  
    if (nl & 0x00000000000000008000ULL) { data += 2; continue; }  
    if (nl & 0x0000000000000000800000ULL) { data += 3; continue; }  
    if (nl & 0x00000000008000000000ULL) { data += 4; continue; }  
    if (nl & 0x00000000800000000000ULL) { data += 5; continue; }  
    if (nl & 0x00000080000000000000ULL) { data += 6; continue; }  
    if (nl & 0x00800000000000000000ULL) { data += 7; continue; }  
    data += 8;  
}
```

 #006

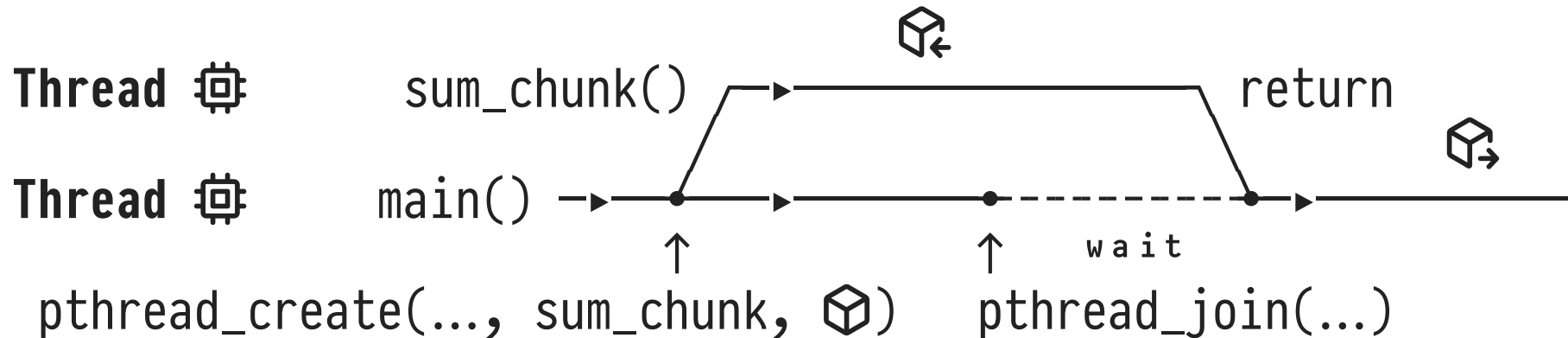
Sum of Quantities ⑤ — C with `mmap(2)` + Block-Wise `ℓ` Search

- 🕒 Query time on Torsten's computer (warm cache): ≈ 0.28 s.
 - C with `mmap()` and block-wise search for `ℓ` processes the CSV file with a throughput of 2.8 GB/s.
 - We match the performance of the built-in `strchr()` function.
- Our code definitely got more complex and fiddly.
 - Slowly getting an impression of how much careful performance engineering is required to build a DBMS kernel.

8 : Sum of Quantities ⑥ — C with mmap(2) + Threads

CPUs feature **multiple cores** that can execute code in parallel. The CPU (M2 Max) in Torsten's computer features $T = 12$ such cores ④.

- ? Split CSV file at \backslash boundaries into T **partitions (chunks)**.
- Spawn T **parallel threads**, each summing column 5 in one partition. Add thread-local partial sums to obtain overall sum.



- Threads use shared memory area ④ to exchange data (*e.g.*, thread ID, pointers to chunk start/end, thread-local sum).


Sum of Quantities ⑥ — C with mmap(2) + Threads

- C declaration of shared memory area :

```
struct chunk {
    pthread_t thread; /* thread ID */
    char      *data;  /* pointers to chunk start/end */
    char      *end;
    int       sum;    /* sum of partition */
};
typedef struct chunk chunk_t;
```



- Code for a thread (sums a chunk):

 #007

```
void *sum_chunk(void *arg)
{
    chunk_t *chunk = (chunk_t*) arg; /* this is the  */

    char *data = chunk->data; /* extract relevant arguments */
    char *end  = chunk->end;

    : /* sum chunk, just like sum-quantity-mmap.c did */

    chunk->sum = sum; /* put return value into  */
    return NULL;    /* NULL ≡  */
}
```

Sum of Quantities ⑥ — C with mmap(2) + Threads

- 🕒 Query time on Torsten's computer ($T = 12$ threads spawned, warm cache): ≈ 0.04 s.
 - Jointly, the threads process the CSV file with a throughput of 18.8 GB/s. This approaches DRAM read bandwidth.
- Profile shows that each `sum_chunk()` + `main()` use about the same chunk summing time. 👍
 - Creating chunk sizes (and thus thread-local work) of roughly equal size has worked well.
 - Wait time after `pthread_join(...)` expected to be small.

Sum of Quantities — Summary so Far

Query Implementation	Query Time 🕒	Throughput
awk	1.60s	471 MB/s
Python	2.75s	275 MB/s
C with getline	0.50s	1.5 GB/s
C with mmap	0.27s	2.8 GB/s
C with mmap + block-wise scan	0.28s	2.8 GB/s
C with mmap + 12 threads	0.04s	18.8 GB/s

- Implementation language and techniques matter **a lot**.
 - 50+ years after the inception of the relational model, database query optimization is a lively field of research.
- Even your laptop can read and process multiple GB/s.
 - Here we saturate everything (but DRAM).
 - Do we always need “big iron” or server clusters? (🐧: “No!”)

9 | Sum of Quantities ⑦ — SQL

Aggregate function `sum()` is all we need to formulate a SQL variant of the benchmark query over the CSV file:⁵

```
D .timer on
D SELECT sum(l_quantity)
   FROM read_csv('../databases/lineitem.csv',
                 header = false,
                 names = [ ..., 'l_quantity', ... ])
```

#008

sum(l_quantity)
153078795



used CPU time

Run Time (s): real 0.448 user 3.322090 sys 0.133819

- 🕒 Query time on Torsten's computer: ~0.45s.
 - Overall CPU time is >3s: DuckDB uses **parallel processing**.

⁵ Command `.timer on` instructs the DuckDB CLI to report query times for all subsequent SQL commands. The DuckDB documentation contains a [complete list of such commands](#) ↗.

Interlude: SQL EXPLAIN

SQL DBMSs typically implement an **EXPLAIN facility**⁶ that exposes details of the system's **query evaluation plans**:

- Supports query and performance debugging (profiling).
- **Shows order of query operations** explicitly, making effects of query optimization visible (e.g., projection pushdown).
- **Profiles query behavior during execution**, annotates plan with:
 - breakdown of how query operations use time,
 - # of rows processed,
 - size of intermediate results (in bytes), ...

```
D EXPLAIN
  query;
```

Do *not* run query. Show query plan and estimated row count.

```
D EXPLAIN ANALYZE
  query;
```

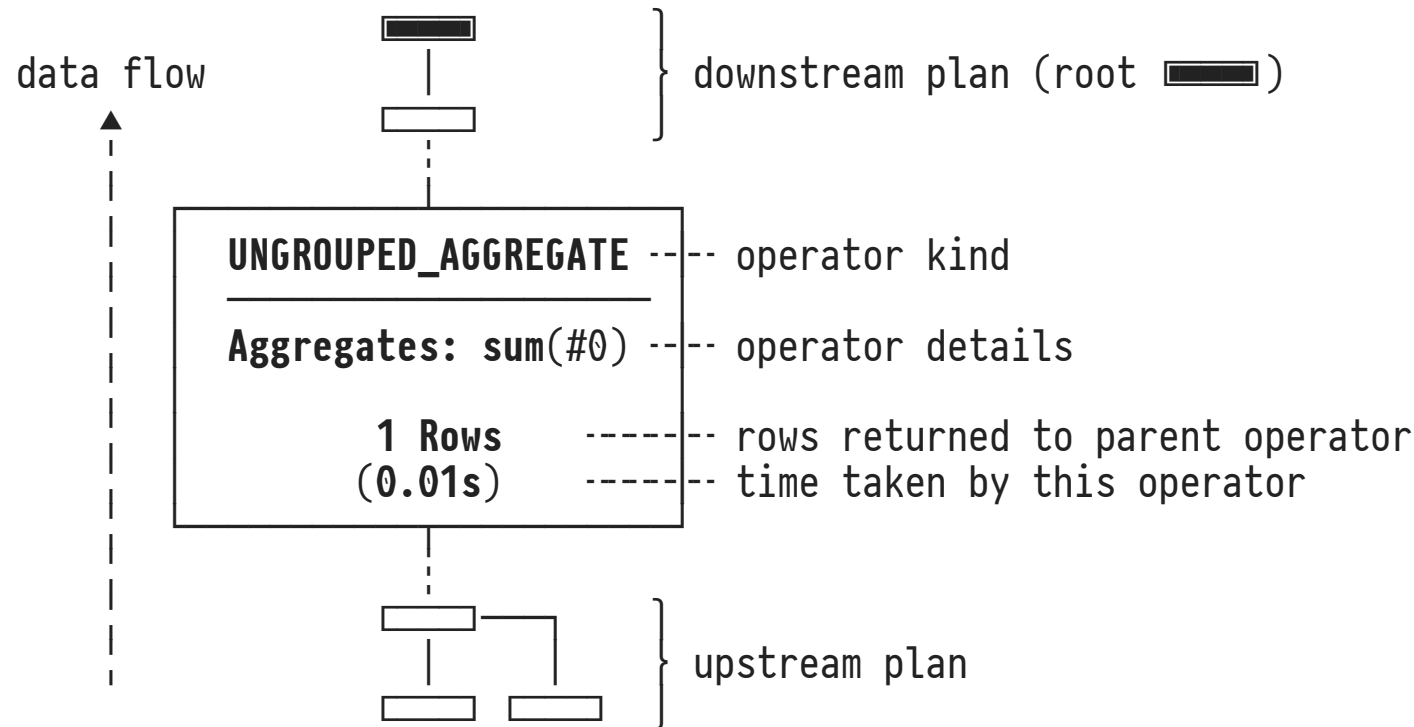
Actually run query. Measure times/rows, annotate plan.

⁶ DuckDB's EXPLAIN facility is extensive. Can see even more plan details via `EXPLAIN (FORMAT json) query`.

Interlude: SQL EXPLAIN

Query evaluation plans visualize bottom-up **data flow**:⁷

- Data sources (e.g., `TABLE_SCAN`) reside at the leaves.
- Query result is produced by the root (top-most operator).



⁷ U Tübingen has developed the [DuckDB Execution Plan Visualizer](#) which can render and inspect plans in the web browser (use `EXPLAIN (ANALYZE, FORMAT json) query` to produce plans that the visualizer can render).

Sum of Quantities 7 — SQL

Copy the CSV file into DuckDB's **native database storage format**:

```
D COPY lineitem FROM '../databases/lineitem.csv';
D SELECT sum(l_quantity)
FROM lineitem;
```

#008

sum(l_quantity)
153078795

```
Run Time (s): real 0.002 user 0.007288 sys 0.000294
```

A 🕒 query time of 0.002s for the benchmark indicates that

- the DBMS uses multiple cores (threads) to evaluate SQL queries,
- the query plan does *not* scan—or skip over—all 16 columns of table `lineitem` (projection pushdown focuses on `l_quantity`),
- column values are *not* stored as text and thus do *not* have to be parsed again and again, and that
- table `lineitem` does reside in DRAM (not in secondary storage).

Design and Implementation of DuckDB Internals

③

Managing Memory + Grouped Aggregation


April 7, 2026

Torsten Grust
Universität Tübingen, Germany

1 | Memory: Fast, But Tiny vs. Slow, But Large

- There are enormous **latency gaps** between accesses to the (L2) CPU cache, RAM, and secondary storage (SSD/HDD):

Memory	Actual Latency ⌘	Human Scale 🕒	Typical Size
CPU L2 cache	2.8 ns	1 s	$\frac{1}{4}$ –16 MB
RAM 🏠	≈ 100 ns	36 s	4–256 GB
SSD	50–150 μs	5–15 h	$\frac{1}{2}$ –4 TB
HDD 🗑️	1–10 ms	4–40 days	1–16 TB



- Fact: faster memory is significantly smaller. We will *not* be able to build cache-only DBMSs. Thus:
 - Keep **persistent database data** on secondary memory (disk 🗑️).
 - During **query processing**, try to keep all of the currently relevant (or: “hot”) data in RAM 🏠.
 - General: attempt to process all queries at the upper levels of the memory hierarchy (▲).

Memory Management in DuckDB

Most DuckDB internals are optimized for **in-memory operation**. By default, 80% of the host RAM¹ is used to process data and queries.

Memory 🧠 is a precious resource. DuckDB is built to use all RAM available but gradually moves to disk-based processing if required (*out-of-core processing*):

- If queries permit, **stream small data chunks** (packs of *vectors*) through the system. Avoid to materialize entire tables in RAM.
- Try to hold **intermediate data structures** (e.g., hash tables) in memory. **Spill to disk** if a query produces huge intermediates.
- In the remaining memory space (if any), cache disk-resident table data to speed up future accesses.

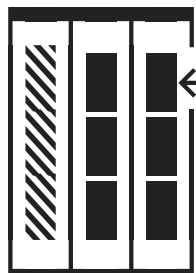
📄 #009

DB lingo: **buffering**.

¹ Compare these 80% to PostgreSQL's default memory buffer size of 128 MB ([shared_buffers](#)). DuckDB's configuration parameter [memory_limit](#) 🐭 controls the maximum RAM size available to the DBMS.

Streaming Query Execution (Pipelining)

Many SQL queries admit **streaming execution** in which **vectors** of only 2048 elements each flow from data source (tables) to result:



←vector

1 Route data chunk (■,■) through query plan

(▨: irrelevant col)

2 Memory requirement is ■+■+δ per thread at any instant (δ: intermediate)

- Simple filter (**WHERE**) + project (**SELECT**) operations.
- Ungrouped aggregations or grouped aggregations (**GROUP BY**) for which the number of groups is small.
- Returning N rows only (where N is small, **LIMIT N**).
- Reading data from one file and writing to another (e.g., reading from CSV and writing to Parquet, **COPY**).

Such queries can be evaluated over larger-than-memory tables even under very constraining **memory_limit** settings.

Spilling



Intermediate data structures larger than available RAM are common when processing complex SQL queries:

- Aggregation in which the grouping criteria create *many* groups (`GROUP BY`: hash table).
- Counting the *distinct* values in a column *C* with many distinct values (`count(DISTINCT C)`: hash table).
- Joining two tables that *both* are larger than memory (join builds hash table for smaller table).
- Sorting larger-than-memory input data (`ORDER BY`).
- Evaluating a complex window function over larger-than-memory input (`... OVER ...<frame>` with a sweeping frame).

DuckDB temporarily uses **out-of-core** memory on disk to hold (parts of) large intermediates (DB lingo: **spilling**).

DuckDB: Unified Memory Management

Use the *same memory area* of maximum size `memory_limit` to perform

1. **non-paged allocations:** fragmented (typically small) bits of data, e.g., as used in pointer-based data structures , and
 2. **paged allocations:** blocks of contiguous data of various sizes, typically holding (parts of) tabular data structures . Typical block sizes range from 32KB to 256KB (default).
- DuckDB prefers paged allocation to avoid memory fragmentation:

Non-paged allocation



After deallocation (▨: freed)



Paged allocation (single size)






Column-Wise Buffering of Base Table Data (👤→🔧)

As long as memory usage permits, DuckDB **buffers** data read from persistent database file `*.db` file in RAM:

- File I/O is performed in 256KB blocks (as opposed to byte-by-byte) which matches the memory manager's default page size.
- Only buffer table columns relevant for query processing (recall `Projections: ...` in operator `SEQ_SCAN`).
- If a page is no longer needed, add it to a LRU queue (DB lingo: **unpin**) as a candidate for replacement (DB lingo: **eviction**).
 - Buffered pages can be useful across queries: only evict once memory indeed is tight and needs to be reclaimed.
 - Buffered pages cache on-disk data (that could be re-read): *no need* to write evicted pages back to disk (🔧↯👤).

Unified Memory Management in DuckDB

DuckDB flexibly assigns the available primary memory  to all three kinds of allocations and thus is able to adapt to workloads/query specifics:



	Non-Paged Intermediate	Paged Intermediate	Buffered Base Table Data
Size	flexible	typical: 32...256KB (max: 2^{62} bytes)	256KB
Lifetime	query	query or session	across queries (database session)
Spilling			(not needed)

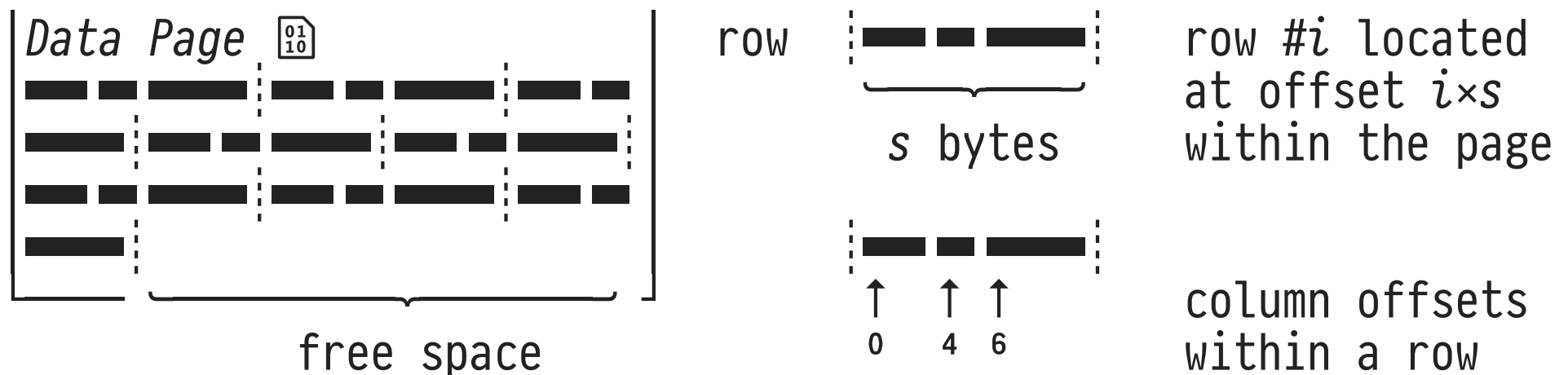
Unified Memory Management in DuckDB


- In contrast to many DBMSs that require a pre-configuration² of buffer size (# of pages) or memory size for intermediates.

² Since DuckDB shares the address space with its host process, it is important that memory consumption is low when the DBMS is idle. DuckDB thus does use not a fixed-size pre-configured buffer.

2 : Layout of Paged Intermediate Data Structures

- DuckDB uses **column-based** storage  for base tables, but relies on a **row-major layout**  for paged intermediate data:
 - Co-locating the columns of a row helps row comparison during sorting, hashing, or joins (column access locality).
- **Fixed-size rows** with **columns of known width** admit efficient row/column access via offset calculations:

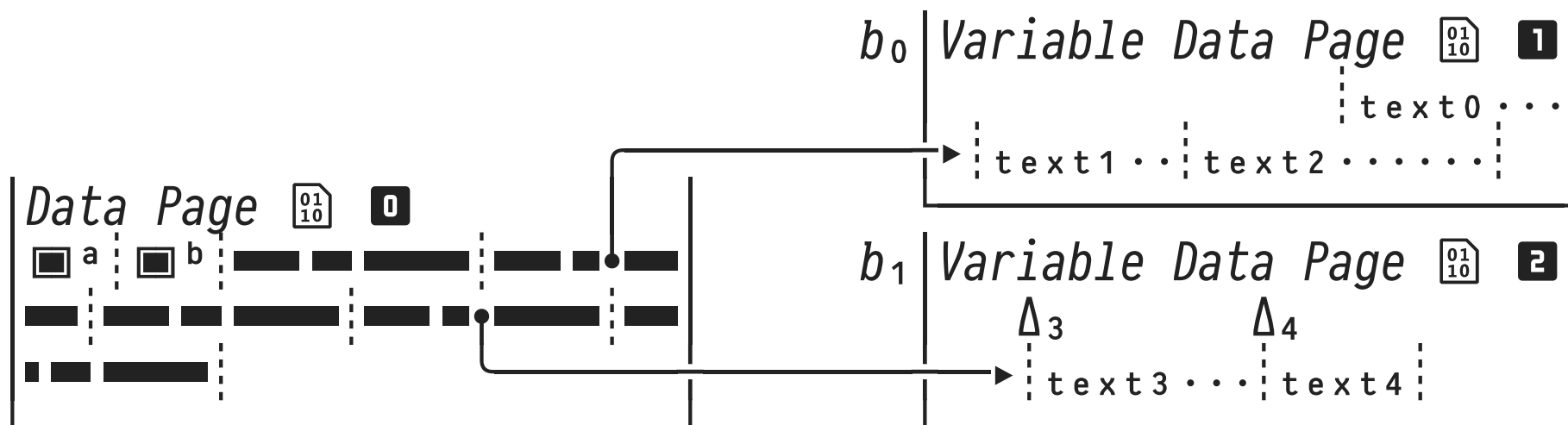


- Row size s and column offsets are known at query plan generation time: store these only once per *Data Page* .

Spilling Pages that Hold Pointer Targets

When a *Variable Data Page* $\boxed{\begin{smallmatrix} 01 \\ 10 \end{smallmatrix}}$ is spilled and later loaded back, its base address b_0 may change and pointers $\bullet \rightarrow$ become invalid. 

1. A string on a *Var Page* $\boxed{\begin{smallmatrix} 01 \\ 10 \end{smallmatrix}}$ with base address b lives at $b+\Delta$.
2. Store b as meta data \blacksquare at the start of the *Data Page* $\boxed{\begin{smallmatrix} 01 \\ 10 \end{smallmatrix}}$.
3. If the *Var Page* $\boxed{\begin{smallmatrix} 01 \\ 10 \end{smallmatrix}}$ is loaded back at base address B , on the next dereference (lazily) replace pointers to $b+\Delta$ by $B+\Delta$.



\blacksquare^a : \langle data page $\mathbf{0}$, row offset 0, var page $\mathbf{1}$ @ b_0 \rangle } meta
 \blacksquare^b : \langle data page $\mathbf{0}$, row offset 3, var page $\mathbf{2}$ @ b_1 \rangle } data

3 | Hash-Based Grouped Aggregation

Grouping and aggregation are core operations in OLAP workloads:

```
SELECT l_orderkey, count(*)
FROM   lineitem           -- TPC-H (sf = 100): 600,000,000 rows
GROUP BY l_orderkey      -- creates 150,000,000(!) groups
```







- DuckDB implements grouped aggregation in terms of plan operator `HASH_GROUP_BY`, i.e., through **hashing**:⁴
 1. Row with `l_orderkey` value o hashes to key k .
Update entry `hash_table[k]`: $(o, count) \rightarrow (o, count+1)$.
 2. Row with `l_orderkey` value $o' \neq o$ also hashes to key k .⁵
Collides with entry for o at `hash_table[k]`. DuckDB uses linear probing to find entry $(o', count)$ —hopefully nearby.

⁵ If DuckDB statically knows that no collisions will occur (e.g., when the grouping key has few distinct values like `lineitem.l_returnflag`), the efficient alternative operator `PERFECT_HASH_GROUP_BY` is used.



⁴ Other tabular DBMSs also implement grouped aggregation through sorting (here: on column `l_orderkey`). Profitable if the query also contains an `ORDER BY l_orderkey` clause (DB lingo: *interesting order*).

DuckDB: Robust External Grouped Aggregation

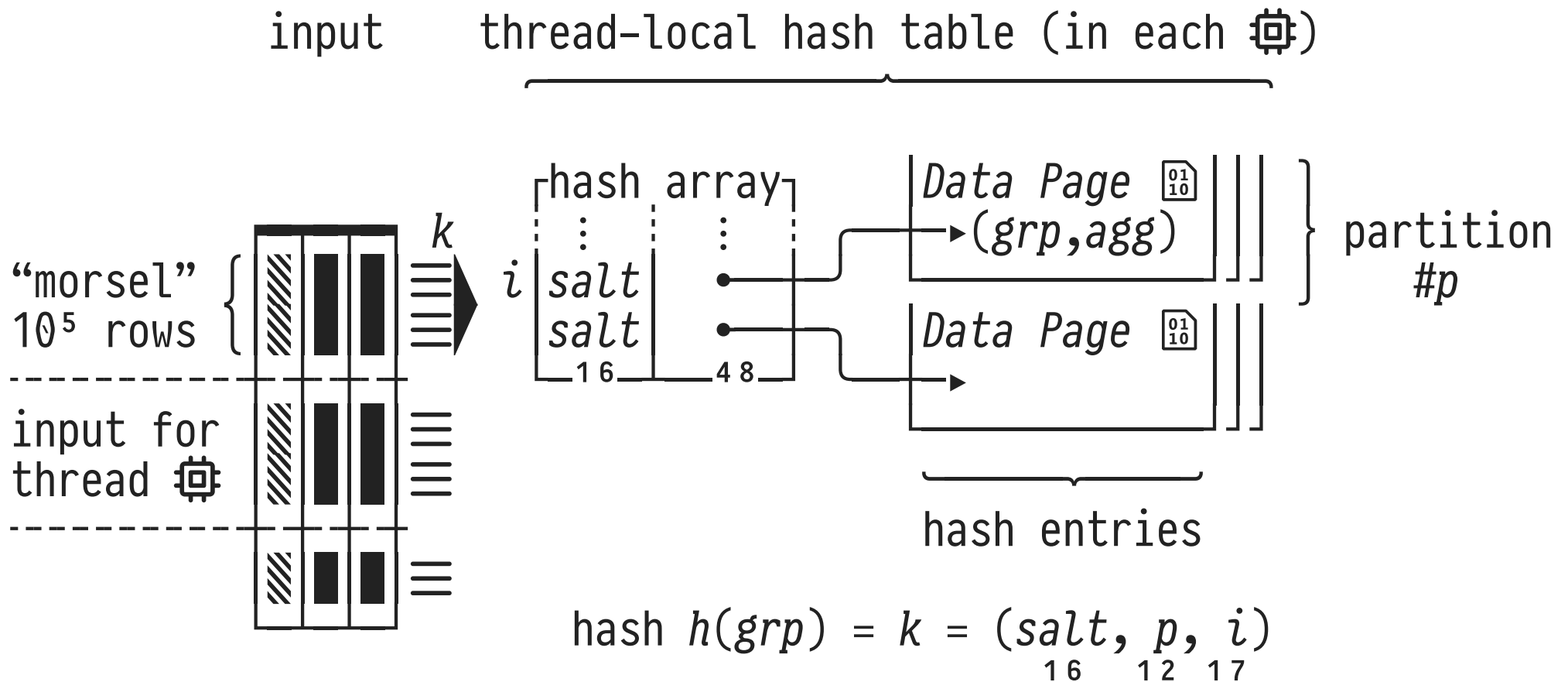
Designing the implementation of `HASH_GROUP_BY`:



- Aim to use all cores    of modern **multi-threaded** CPUs.
- If we create **many groups**, the hash table entries may exceed available memory  and thus need to be **spilled** to disk. 
 -  Maintain hash table entries as a *paged intermediate data structure*, let the memory manager handle the spilling.

DuckDB performs external grouped aggregation in two phases:

- Phase **①**: Thread-local pre-aggregation (one  per morsel).
- Phase **②**: Partition-wise aggregation (one  per partition).

External Aggregation (Phase ①): Thread-Local Aggregation



- Relevant columns of input table (grouping keys grp + arguments of aggregate agg) are split into **morsels** of $\approx 100,000$ rows each.
- Each thread  reads a morsel. Typical: more morsels than .

External Aggregation (Phase ①): Thread-Local Aggregation

In each thread :

└ For each row in morsel:

└ Compute hash key $k = h(grp)$. Split the bits of k :

└ Lower 17 bits i : index in hash array (131,072 slots)

└ Middle ≤ 12 bits p : partition # for hash entry

└ Upper 16 bits $salt$: used to **optimize** collision handling

└ Hash array slot i occupied?

└ **Are salt values equal?**

└ **Yes** Follow $\bullet \rightarrow$ to entry (grp', agg) on data page .

└ **Is $grp' = grp$?**


└ **Yes** No collision. Update value of agg . Done.

└ **No** Collision. Linear probing to find proper slot.

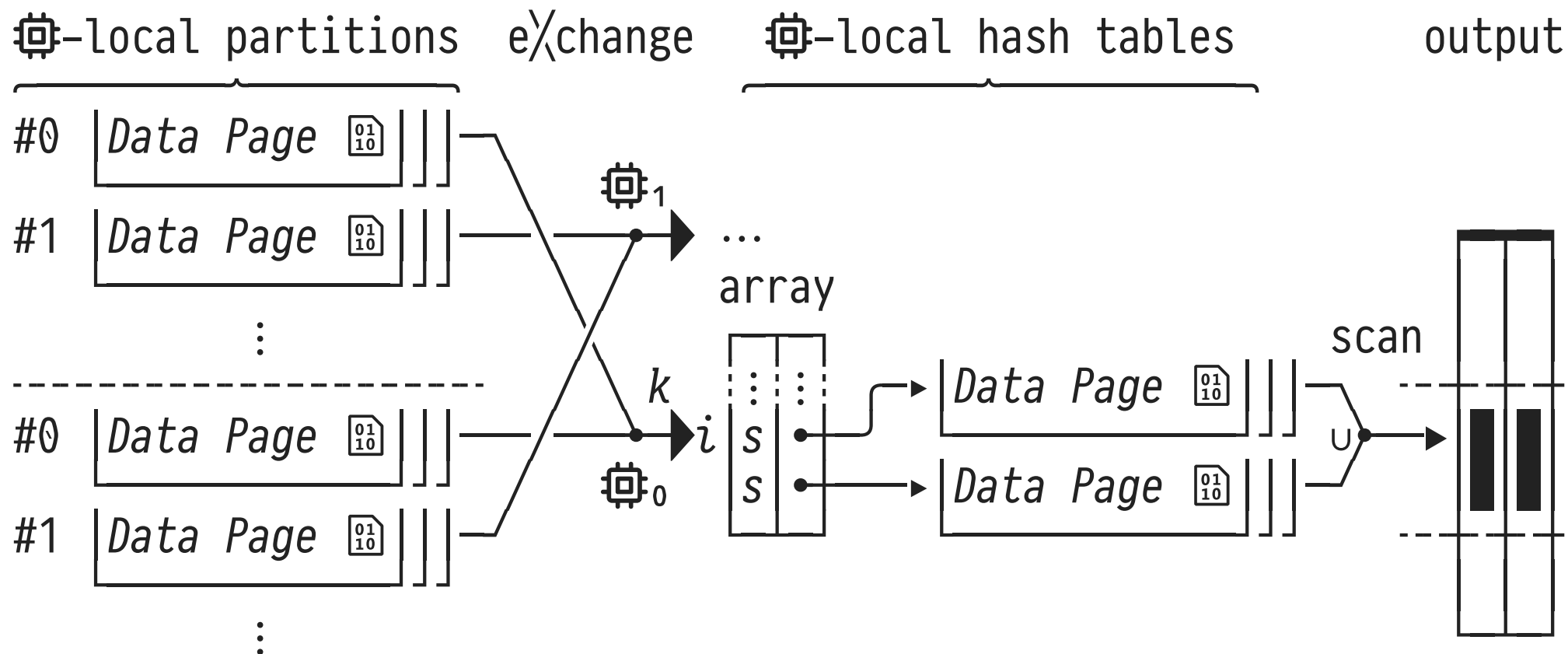
└ **No** Collision. Linear probing.

└ **No** Add entry $e = (grp, agg_0)$ to data page  in partition # p .

└ Place $(salt, \bullet \rightarrow$ to $e)$ in array slot i .

- When hash slots are $\frac{2}{3}$ full and collisions become frequent:
 1. Unpin data pages —the memory manager may evict these.
 2. Empty hash array, continue (reset OK because of Phase ②).

External Aggregation (Phase ②): Partition-Wise Aggregation



- Run one thread ⚙ per partition $\#$. Choose bits for p such that the hash table for a fully aggregated partition fits memory (create more but smaller partitions).
- When a ⚙ ends, immediately scan its output data pages Ⓜ . These become a morsel to be fed to the downstream plan.

Design and Implementation of DuckDB Internals

④

Sorting Large Tables

April 7, 2026

Torsten Grust
Universität Tübingen, Germany

1 | Sorting is a Core Operation

The tabular model is based on unordered bags of rows, but the **sorting** of (intermediate) results of SQL queries is material:

1. **ORDER BY** clause,¹
2. ordered aggregates (`list(e_1 ORDER BY e_2)`),
3. window functions (`row_number() OVER (ORDER BY e)`),
4. joins over time series data (**ASOF JOIN**), or
5. to ensure deterministic **LIMIT** and **OFFSET** behavior.

```

SELECT *           -- wide payload (here: reorder all columns)
FROM   lineitem
ORDER BY l_shipdate DESC NULLS FIRST , l_quantity;
--                                                 ↑
--           defaults: ASC  NULLS LAST  lexicographic ordering

```

¹ In DuckDB, can adapt ordering defaults via `SET default_order = {ASC,DESC}` and `SET default_null_order = {NULLS_FIRST,NULLS_LAST,NULLS_FIRST_ON_ASC_LAST_ON_DESC,NULLS_LAST_ON_ASC_FIRST_ON_DESC}`.

Sorting (Large) Tables in DuckDB

To this day,² DuckDB's internal sorting routines are on the workbench 🛠️ and still undergo active development.

Key aspects:

- **Efficient row comparison** that respects column types, directions (↑↓), handles **NULLs**, allows wide keys (lexicographic ordering).
- Adaptation to **pre-sorted input data**.
- Sorting strategy can process **tables larger than main memory**, spilling to disk 🗄️ if required.
 - Effectively, sorting needs to materialize its entire input table (last input row could be the first in sort order).
- **Balanced use of all available CPU cores** 🏗️ during sorting.

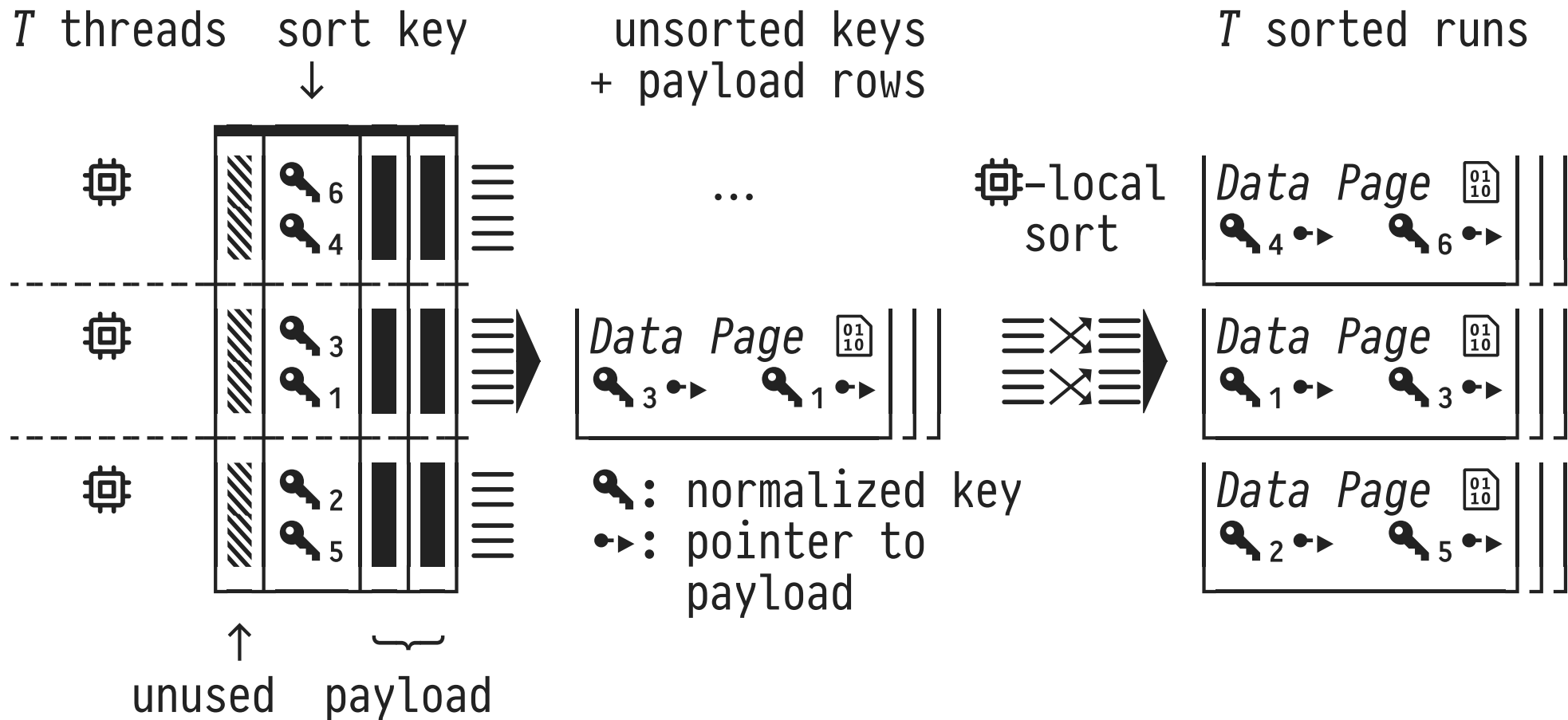
² The last major rewrite of DuckDB's sorting code has shipped with version 1.4.0 in September 2025.

2 : DuckDB: Two-Phase Merge Sort

DuckDB implements a variant of a **two-phase merge sort** strategy that utilizes all available CPU cores ☒ (say T) in both phases:

- **Phase ① (thread-local sorting):**
 - Each ☒ : read $1/T$ fraction of input, gather relevant columns to form rows on data pages $\begin{bmatrix} 01 \\ 10 \end{bmatrix}$, **normalize sort keys** 🔑 ,
 - then generate a **sorted run** of keys.
- **Phase ② (T -way merge of the sorted runs):**
 - Precompute boundaries of T non-overlapping **segments** in each run that can be **merged independently** by the ☒ s.
 - Each ☒ immediately outputs its merged run segments for consumption by the downstream query plan.

Phase ①: Thread-Local Sorting



Phase ①: 🔑 Key Normalization

- Sorting compares keys to decide row order \bowtie . These comparisons are frequent, but can be complex and costly:

```
⋮
ORDER BY  $c_1$  ASC NULLS LAST,  $c_2$  DESC NULLS FIRST
```


- Dispatch on type of c_i to select proper $<$ operator, respect sort order ($\uparrow\downarrow$), implement lexicographic ordering (column c_1 dominates), proper **NULL** treatment.
- **💡 Key normalization:** Map³ values in columns c_i to a byte sequence such that a *single* $<$ comparison on two sequences can decide row order.

 #014

³ DuckDB's key normalization is available at the user level in terms of scalar SQL function `create_sort_key()`.

Phase ①: Key Normalization for Fixed-Size Keys

CPUs compare fixed-size types (e.g., 32/64-bit ints) faster than variable-length byte sequences.

-  If sequence length is known to be n bytes, store it in groups of $\lceil n/8 \rceil$ 64-bit ints. Directly compare these integers:⁴

```

struct FixedSortKey {
    uint64_t      part0;    // assume  $8 < n \leq 16$ : group size = 2
    uint64_t      part1;
    sort_key_ptr_t payload; // (only if query has payload)
};

bool LessThan(const FixedSortKey &k1, const FixedSortKey &k2) {
    return k1.part0 < k2.part0 ||
           (k1.part0 == k2.part0 && k1.part1 < k2.part1);
}

```

⁴ If SQL queries have no payload (`SELECT c_1, c_2 FROM ... ORDER BY c_1, c_2`), DuckDB only keeps the normalized keys in fields `part0, ...`: normalization can be inverted before output is generated.

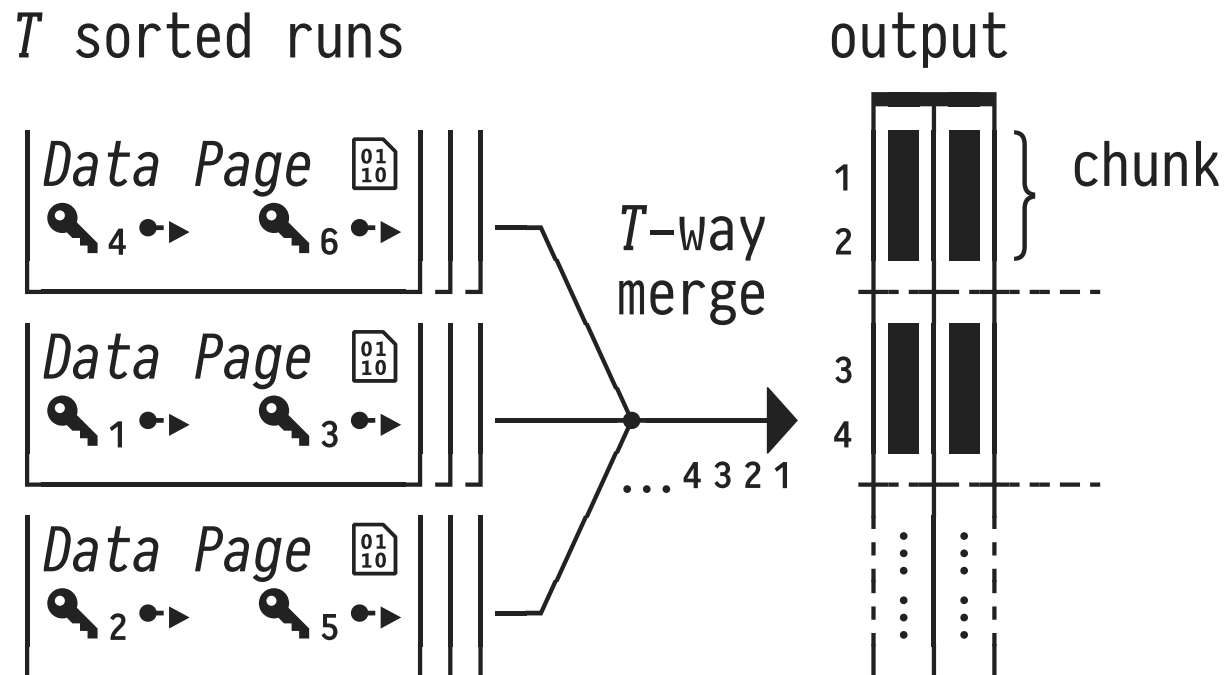
Phase ①: Sorting Algorithms

Typed keys (that can be compared via `<`) fit off-the-shelf sort implementations that assume the C++ `std::iterator` interface.

- DuckDB implements `std::iterator` for chains of 256KB data pages:
 $\boxed{\begin{smallmatrix} 01 \\ 10 \end{smallmatrix}} \rightarrow \boxed{\begin{smallmatrix} 01 \\ 10 \end{smallmatrix}} \rightarrow \dots \rightarrow \boxed{\begin{smallmatrix} 01 \\ 10 \end{smallmatrix}}$, maps element indices to (page#, offset) pairs.
 - (A single 256KB data page would only be able to hold about 10,000 24-byte `FixedSortKey` structs. DuckDB aims to generate runs much longer than that.)
- DuckDB thus is able to adapt and combine three existing sorting algorithms:
 1. *Vergesort* (detects runs of already sorted data).
 2. *Ska Sort* (radix sort on the first 64 bits of the key).
 3. *Pattern-defeating QuickSort* (fallback, if the first 64 bits do not already order the data).

Phase ②: T -Way Merge

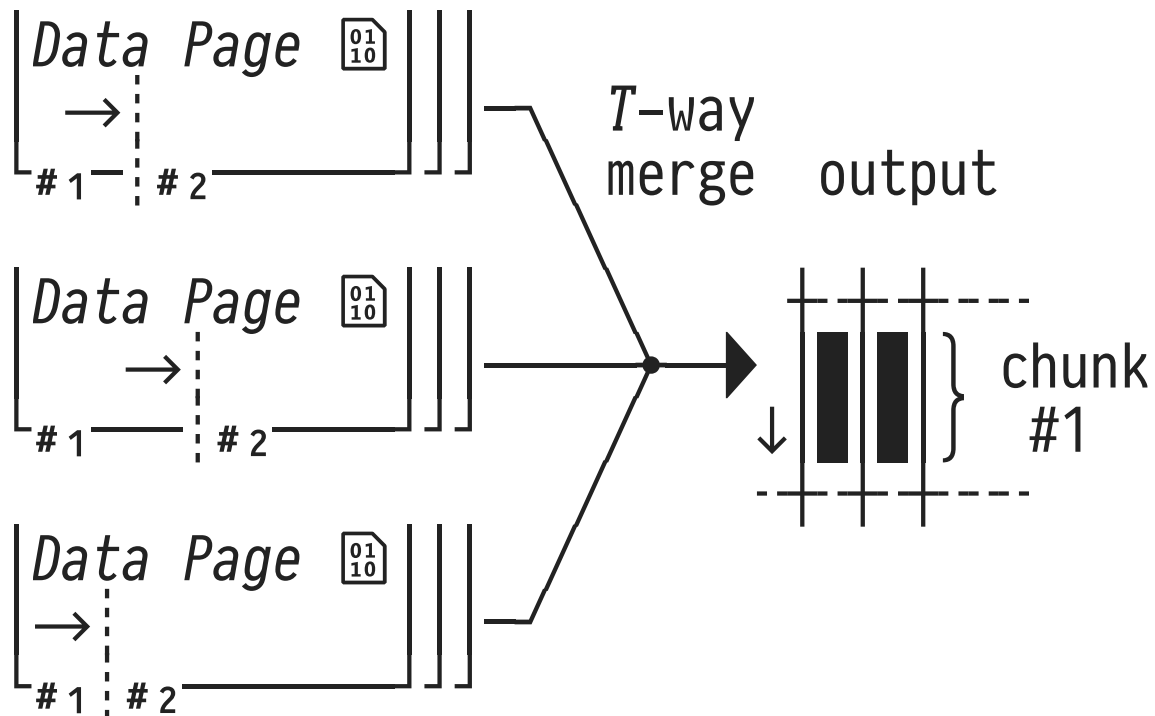
- Consume the T sorted runs produced by Phase ①:
 - Merge pairs ($\text{key}_k, \bullet \rightarrow$) using $<$ on key key_k ,
 - dereference pointers $\bullet \rightarrow$ to access payload,
 - emit sorted output columns chunk-by-chunk:



- How can all CPU cores  contribute equally during this merge?

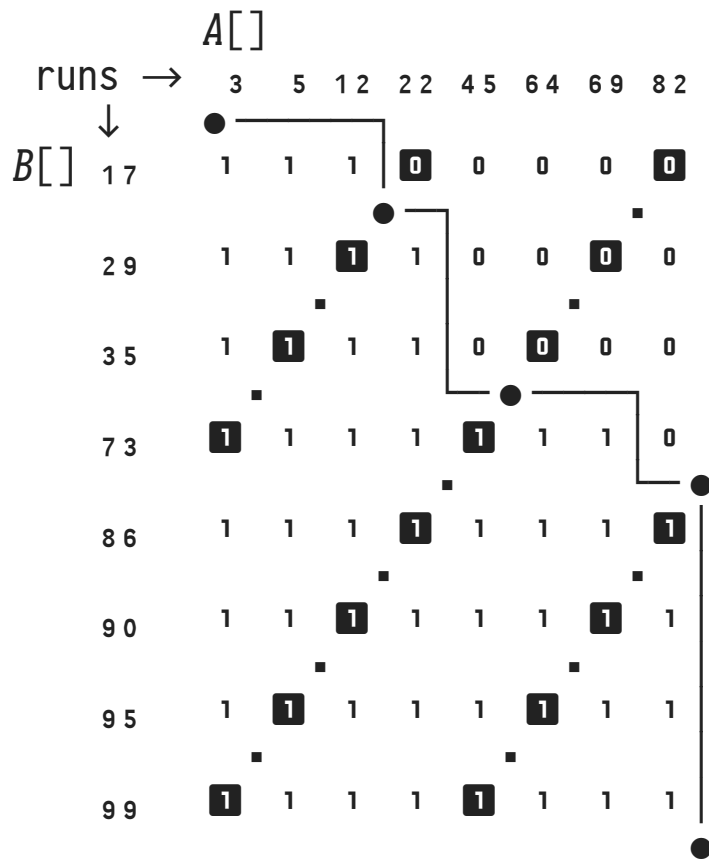
Phase ②: Parallelizing T -Way Merge


T sorted runs



- Runs contribute rows in **segment #1** (up to \vdots) to output chunk #1.
- Run segments $\#i$ hold the rows for chunk $\#i$. No overlap between segments.
- 💡 Precompute segment boundaries \vdots . Threads 🔧 work on one segment.
- Once run segments $\#i$ are merged, the thread 🔧 *immediately* emits chunk $\#i$. The downstream plan can assemble the sorted output based on indices $\#i$ (DuckDB lingo: **batch indices**).

Phase ②: Precomputing Segment Boundaries (*Merge Path*)



- Merge matrix $M[i,j]$: 1 if $A[i] < B[j]$, 0 otherwise.
- Merge path  marks the 1/0 boundary.
- The n th point of the merge path lies on the n th cross diagonal \dots of M .
- Partitioning the path into k segments of equal size can be done by considering $k-1$ evenly-spaced diagonals.
- Point \bullet marks the only transition $1 \rightarrow 0$ on the diagonals. Can be found by binary search on A and B . (Merge matrix and merge path *need not* be constructed.)
- Points \bullet define segments in A and B ($\vdots \equiv$ boundaries):

A:	3	5	12	22	45	64	69	82	
B:	17			29	35	73			86 90 95 99
			#1		#2		#3		#4

- DuckDB generalizes the *Merge Path* idea from two runs to T runs. (Profiling shows that boundary computation uses $\leq 2\%$ of the overall execution time.)

Design and Implementation of DuckDB Internals

⑤

The ART of Indexing

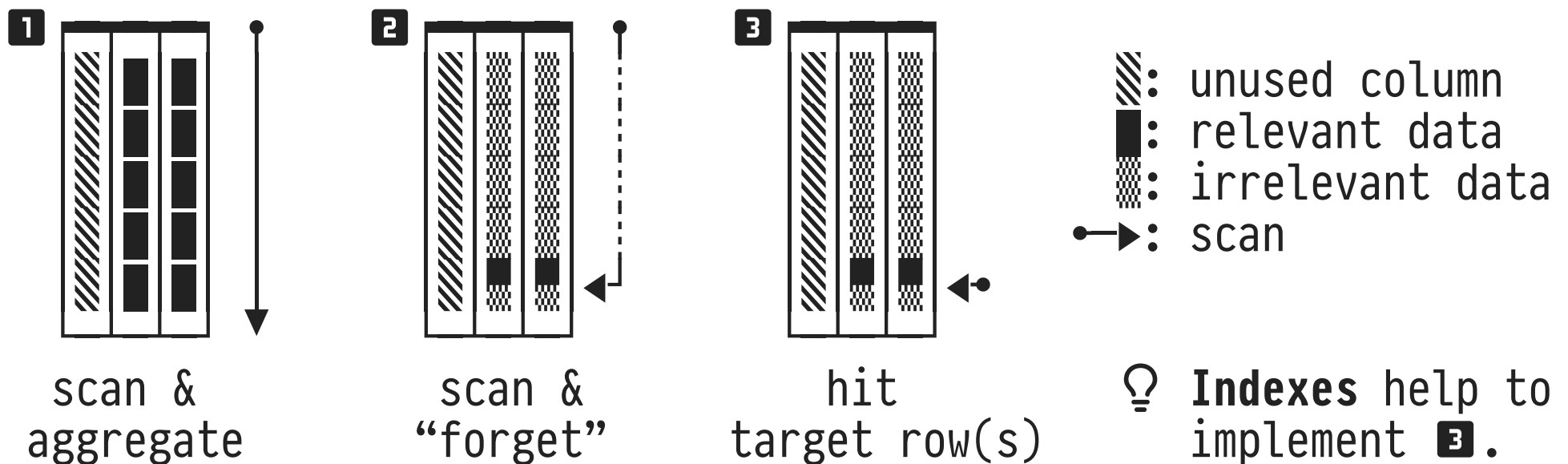
April 7, 2026

Torsten Grust
Universität Tübingen, Germany

1 | Scanning All Rows vs. Narrowing In On Few Rows




The internals of DuckDB have been designed to efficiently support analytical SQL queries: “*online analytical processing*” (OLAP).

- Typically, these queries read (and aggregate) *all rows* of the input tables **1**.
- Yet, scanning all rows wastes I/O and memory bandwidth **2** if a query focuses on *small row subsets (or even a single row)*:



Indexes in DuckDB

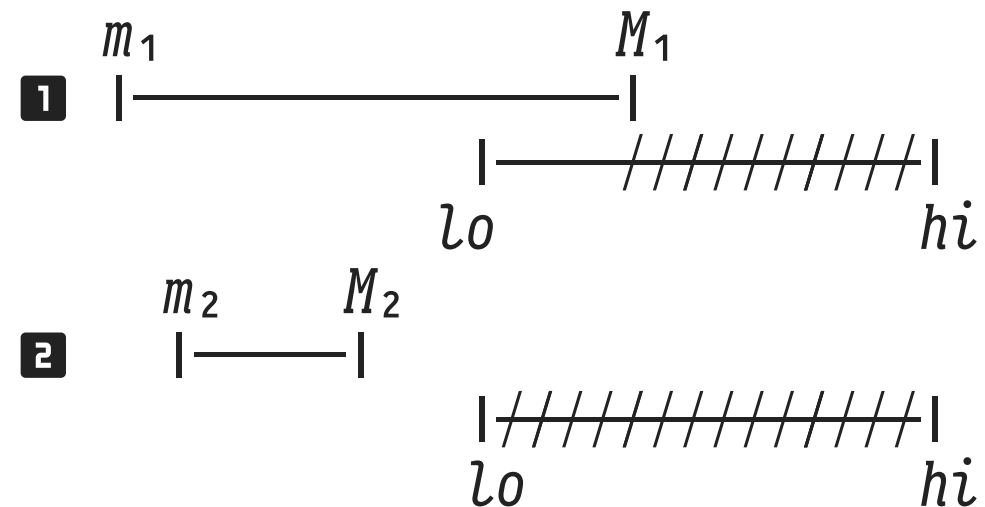
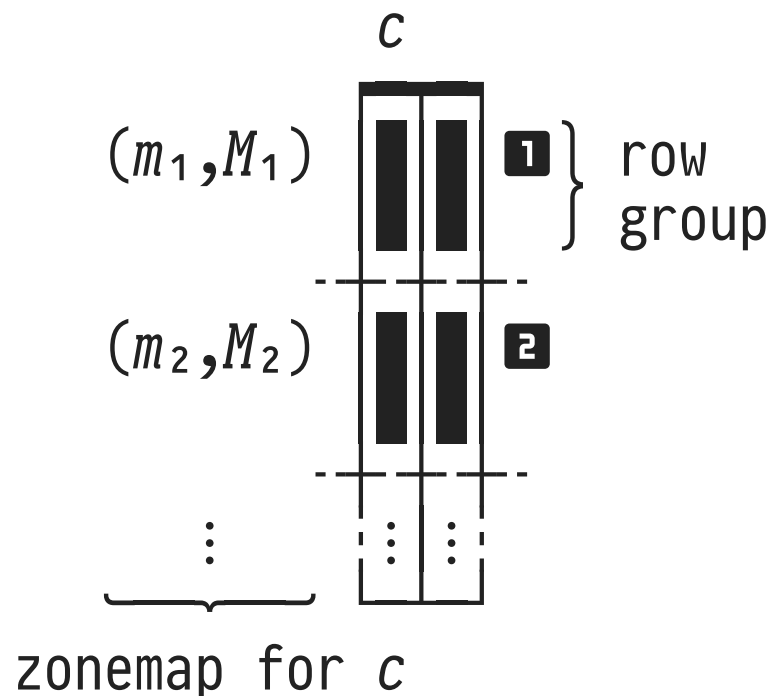
DuckDB implements two kinds of **indexes**:

1. **Zonemaps** (also: *min-max indexes*) are integral part of the columnar table storage 
 - Automatically created, always present for all columns.
 - Used when predicates are pushed down into **Sequential Scan**.
2. **Adaptive radix trees (ART)** are tree-shaped data structures maintained outside/in addition to table storage +
 - Consume working memory space.
 - Require maintenance on table updates.
 - Created either
 - manually via SQL's DDL statement **CREATE INDEX** or
 - automatically when constraints **UNIQUE**, **PRIMARY KEY**, and **FOREIGN KEY** are declared for a table.

2 : Zonemaps

Zonemaps are baked into DuckDB's columnar table storage format:

- Each column c is divided into *row groups* of 120K (122880) rows.
- Each row group holds a (min, Max) entry that—quite roughly—characterizes the contained values.
- Operator **Sequential Scan** can safely **skip row groups** for which a predicate $lo \leq c \leq hi$ will always fail (///):

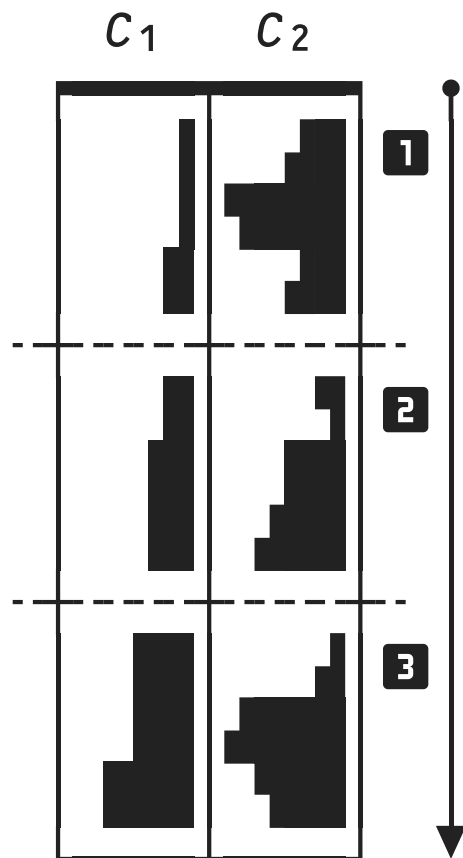


- Need to scan row group **1** 
- May skip row group **2** 


Zonemaps and Column Ordering

Column ordering affects the effectiveness of zonemaps.


 #015



c_1 : In each row group, the zonemap entries (min, Max) for c_1 have a small span $|—|$.

⇒ A large number of row groups will never satisfy $lo \leq c_1 \leq hi$. Can skip. 

c_2 : All spans in unordered column c_2 are wide $|—————|$ and cover the active domain of c_2 .

⇒ Most (all) row groups will potentially contain hits for $lo \leq c_2 \leq hi$.
No skipping. 

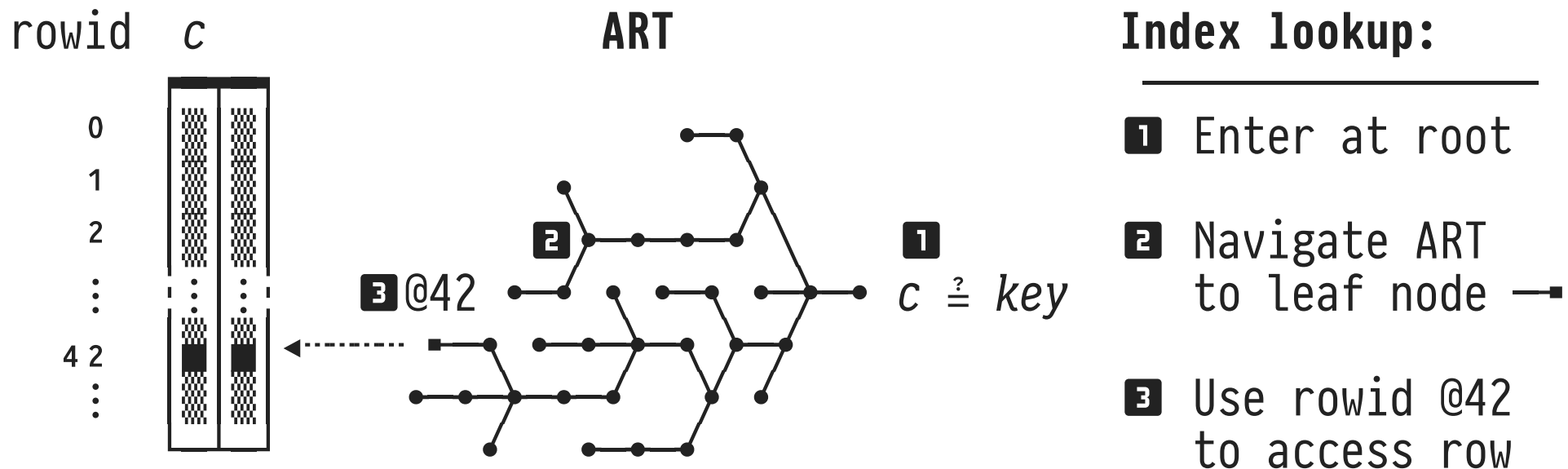
💡 Keeping base tables ordered may help predicate evaluation.

3 | Adaptive Radix Tree (ART) Indexes

ART indexes are **ordered search tree structures** built to evaluate simple **equality and range predicates** (assume an ART on column c):

$c = key$ $c \text{ IN } (key_1, key_2, \dots)$ $c < key$ $c > key$

- ARTs and tables are separate. These extra data structures
 - use **row IDs** to refer to rows in their associated table,
 - must be created, then maintained under table updates, and
 - occupy space in working memory.



Creating ART Indexes in DuckDB

1. Manual:

```
CREATE [UNIQUE] INDEX index ON table (column [, column, ...]);
DROP INDEX [IF EXISTS] index;
```

- **NB.** More indexes aid query evaluation but incur maintenance and space¹ overhead. A tradeoff in physical database design.

2. Implicit (supports constraint enforcement). DDL statement **1** creates on a unique ART index on *t(c)* behind the scenes:

```
1 CREATE TABLE t (... , c ... PRIMARY KEY, ...); -- also: UNIQUE
2 CREATE TABLE s (... , f ... REFERENCES t(c), ...);
```

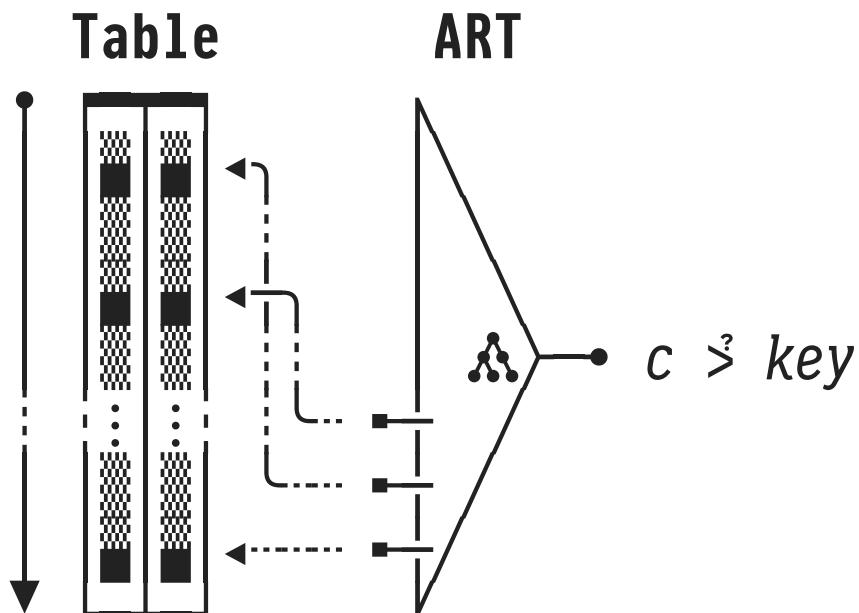
- DDL statements like `INSERT INTO s VALUES (... , key, ...)` or `UPDATE t SET c = key` lead to lookups $c \stackrel{?}{=} key$ in that index.

¹ As of DuckDB 1.4, memory occupied by ART indexes is under control of the buffer manager but cannot be evicted. Work to rectify this is underway as I type this.

Use the Index? Ignore the Index?

An index lookup predicate p may yield *multiple rows* in table t :

- $p \equiv c \stackrel{?}{=} key$: If c is not unique in t : a single leaf node \rightarrow holds multiple rowids $\cdots \rightarrow \cdots \rightarrow \cdots \rightarrow$.
- $p \equiv c \stackrel{?}{>} key$: **1** perform lookup $c \stackrel{?}{=} key$, **2** access rows using the rowids in all leaves $\rightarrow \rightarrow \rightarrow$ following the found leaf:



- Lookup finds one leaf (or adjacent leaves $\rightarrow \rightarrow \rightarrow$). #016

- Yet the rowids $\cdots \rightarrow$ may point all over table t . May need to “jump around” to collect rows. Violates memory locality.

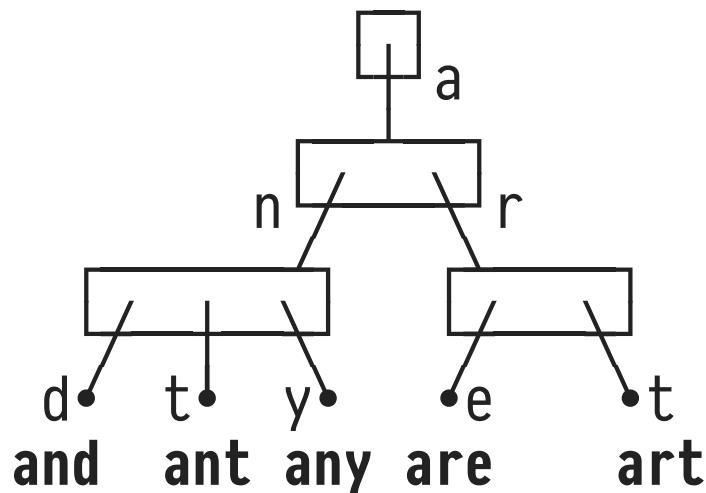
💡 Use index only if **selectivity** $sel(p)$ of predicate p is low:

$$0 \leq sel(p) \stackrel{\text{def}}{=} (\text{SELECT count(*) FILTER (p) FROM t}) / \text{count(*) FROM t} \leq 1$$

4 : The Internals of Adaptive Radix Trees

An ART on $t(c)$ uses the **bit-wise representation of values** in column c to organize itself. Values are *not* hashed or compared.

- Divide values into groups of s bits each (s is the *span*).
- During tree traversal, the next s bits of *key* determine the child node to descend into.




- Example: $c :: \text{text}$, span $s = 8$ bits.
- Values have $k = 24$ bits.

	s bits		$\text{ant} < \text{any}$
ant	$=_2$ 01000001	01001110	01001110
any	$=_2$ 01000001	01001110	01011001

common prefix \equiv shared path in ART

ART Internals: Mapping Values to Bit Sequences

Map values to bit sequences whose **lexicographic order** properly reflects value sort order:

- **Unsigned ints:** use standard binary representation (little-endian machines: reverse byte order so that MSB comes first).
- **Signed ints** (in two's complement): flip the sign bit, then treat like unsigned ints.
- **IEEE 754 floats:** **1** always flip the sign bit, **2** if the sign bit was originally set, now flip all bits.  #017
- **Text strings:** map UTF-8 byte sequence to binary², end strings with $\backslash\text{u}_L^N$ (values must not be prefixes of other values).
- **Composite values** (v_1, \dots, v_n): map the individual fields v_i , then *concatenate* the resulting bit sequences.
- **NULL:** increase bit length k (e.g., by one bit/one byte) to accommodate the additional value.

² DuckDB uses function `ucol_getSortKey()` of the C/C++ library ICU to perform this mapping.

ART Internals: What is a Good Span?

- The height of ARTs depend on the value bit length k (not the number n of entries).
 - An ART for k -bit values has $\lceil k/s \rceil$ levels of inner nodes.
 - ART lookup and insert operations have complexity $O(k)$.
- ARTs vs. binary search trees (BSTs):
 - Height: if $n > 2^{k/s}$, ARTs have smaller height than BSTs ($\log_2(n)$).
 - Lookup: for values of k bits (k large), comparison ($<$) is $O(k)$. Complexity of BST lookup thus is $O(k \cdot \log_2(n))$.

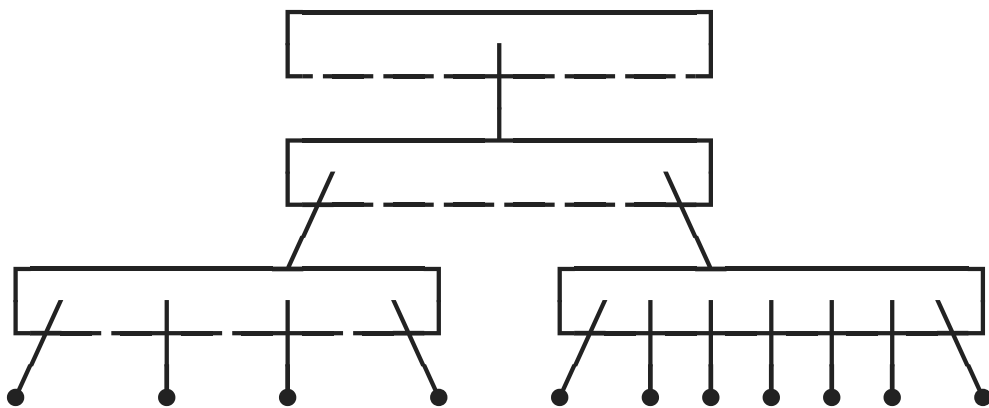
💡 Optimize ART **performance**: Work with a large span s ! 👍

🤖 But how about the **space usage** of such ARTs? 🗨️

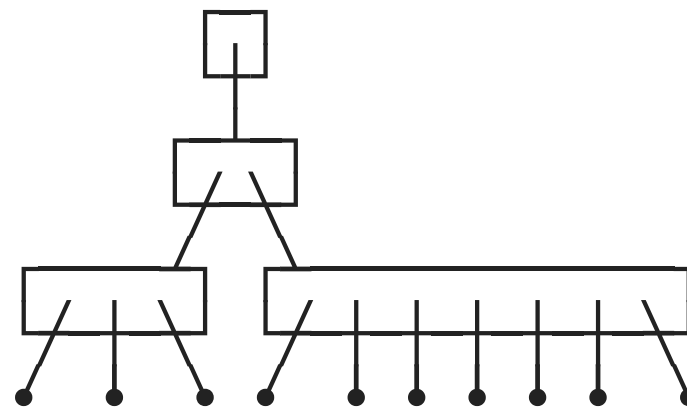
ART Internals: Designing Inner Nodes

- A simple representation of inner ART nodes:
 - Array of 2^s child pointers. (Nodes would have **fixed size**.)
 - Use s -bit chunk of *key* to index array + access child node.
 - But: Most child pointers will be NULL (\perp —), space usage will be excessive if s is large (grows exponentially):

Node size fixed



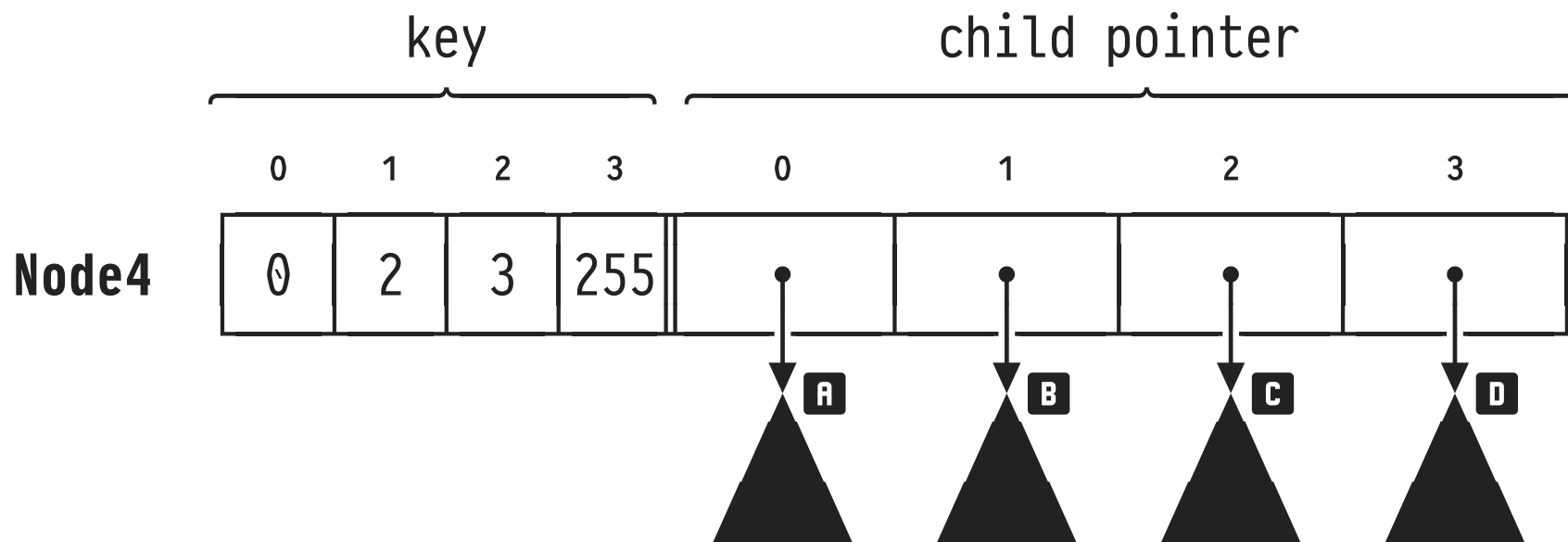
Node size adaptive



💡 Use a fixed large s but **adapt node size** (use variable fanout).

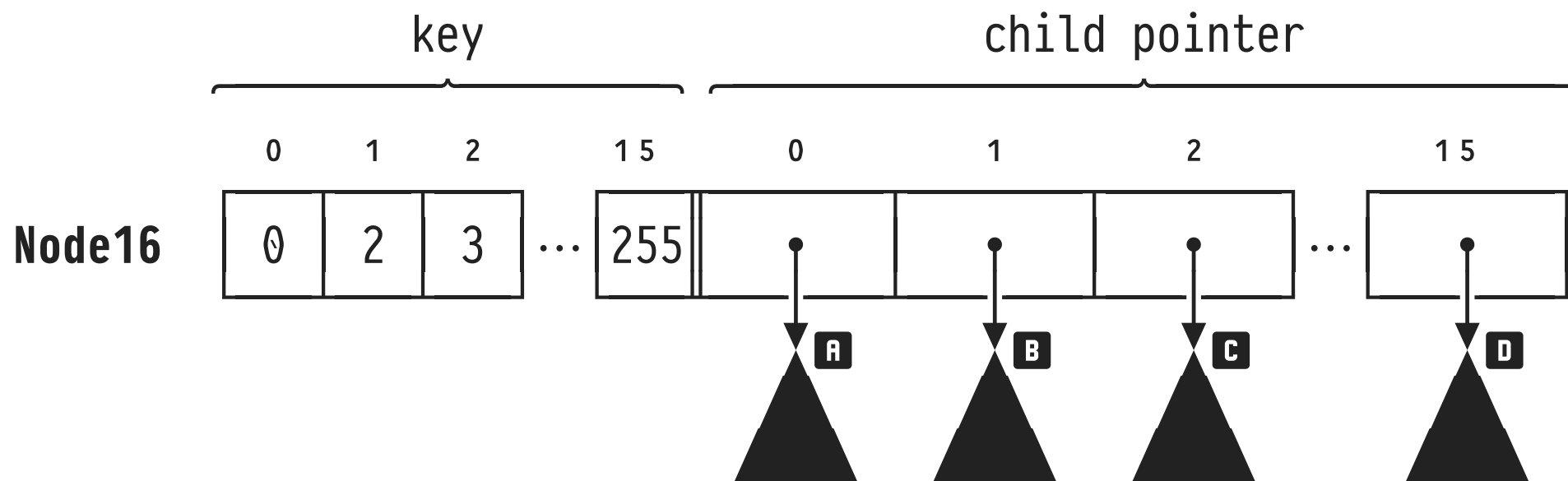
ART Internals: Four Inner Node Types ①

- DuckDB: $s = 8$ (can cut values into bytes, admits large fanout).
- Choose inner node type based on the number of non-NULL child pointers: 4, 16, 48, or 256 ($= 2^s$).
 - On value insertion/deletion, switch to larger/smaller node type when node overflows/underflows.



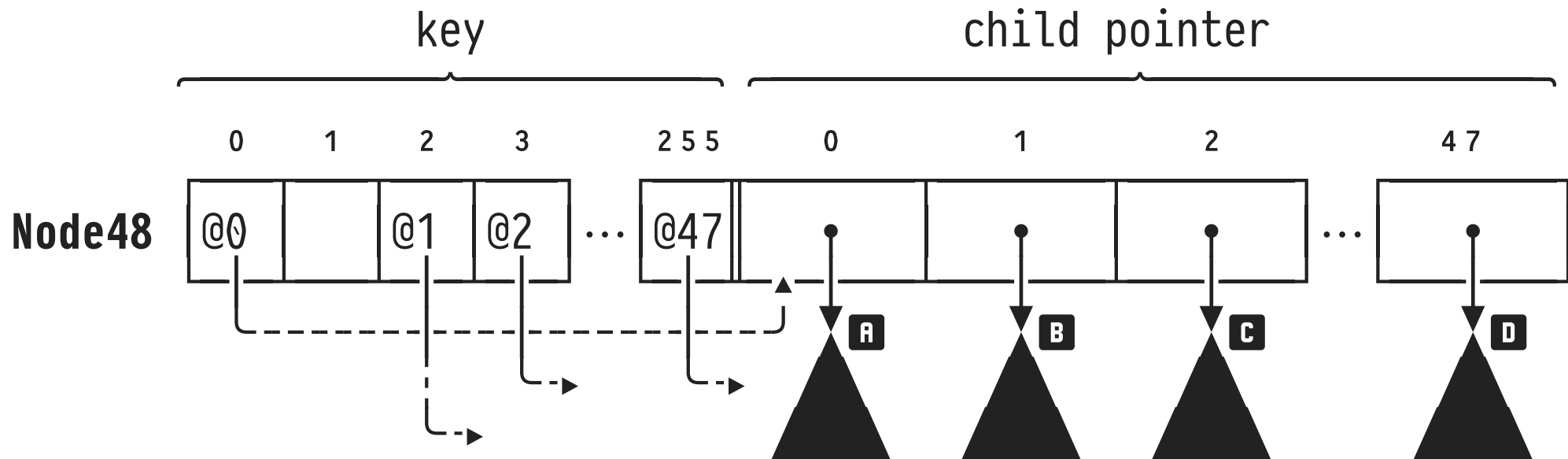
- **Node4:** Array of 4 keys is ordered (search left to right).
- Value/pointer pairs stored at corresponding indices 0...3.

ART Internals: Four Inner Node Types ②



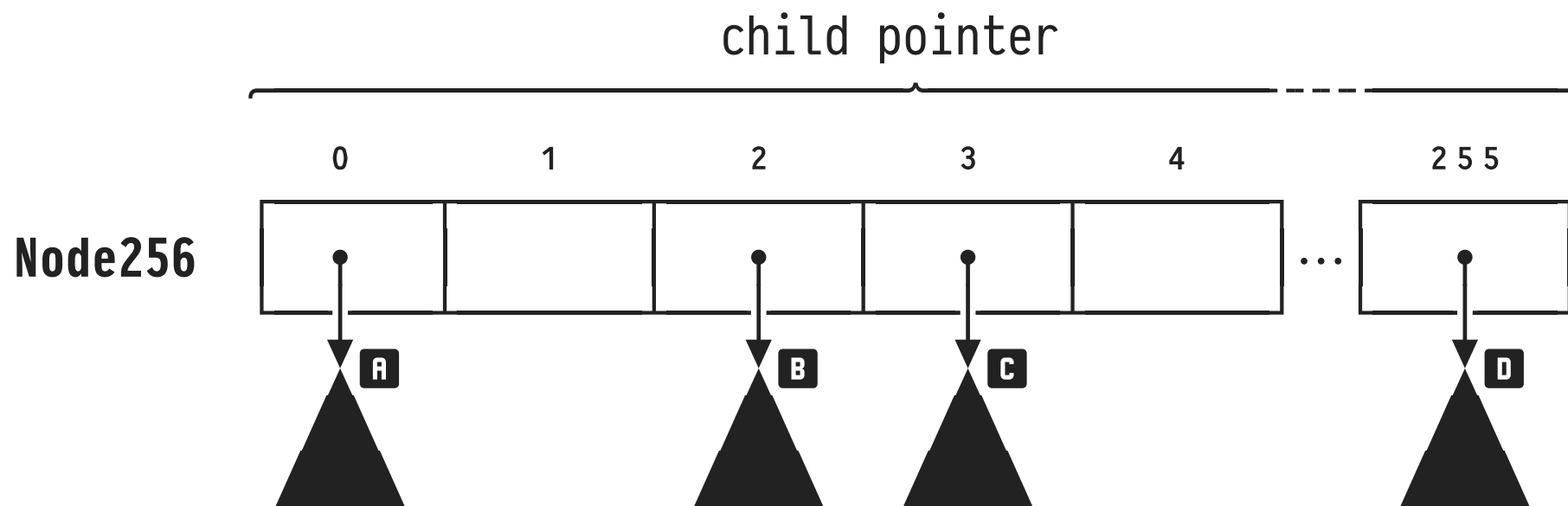
- **Node16:** stores 5...16 child/pointer pairs.
- Locate current *key* byte by binary search or parallel SIMD comparisons.

ART Internals: Four Inner Node Types ③



- **Node48:** With 17...48 values in a node, searching the value array becomes expensive. Thus: do *not* store values. Instead:
 - Use current byte of *key* as index into array of 256 indices.
 - `key[<key byte>]` holds the index `@i ∈ {0...47}` of the corresponding child pointer. Access `pointer[@i]`.
- Saves space compared to an array of 256 8-byte pointers:
 $640B = 256B + 48 \times 8B < 256 \times 8B = 2048B.$

ART Internals: Four Inner Node Types ③

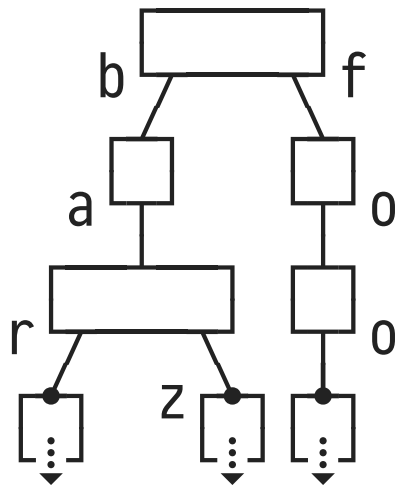


- **Node256:** If an inner node needs to hold 49...256 entries, invest $256 \times 8B = 2kB$ to hold an array of 256 child pointers.
 - Simply return `pointer[<key byte>]`.

ART Internals: Leaf Nodes

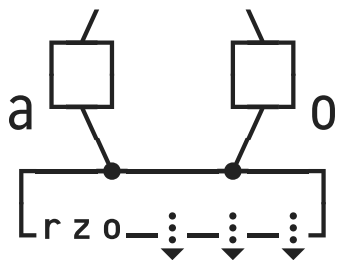
In an ART for $t(c)$, **leaf nodes** hold rowids that point to the rows in table t in which column c holds the search value key .

- Possible ART leaf designs:



(This is an ART for values **bar**, **baz**, **foo**.)

} **Single-entry leaves**, each hold rowids (↓) for one value val .



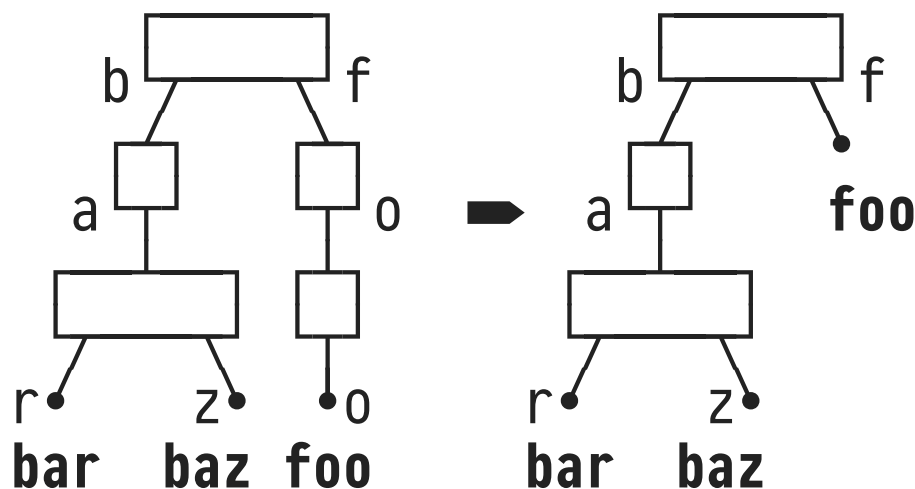
} **Multi-entry leaf**, holds rowids for multiple values. DuckDB has Leaf7/Leaf15/Leaf256 much like Node4...256 (hold rowids, not pointers).

ART Optimizations: Lazy Expansion + Path Compression

To efficiently use space and save tree traversal effort during lookups, try to remove inner ART nodes—and thus **reduce tree height**—if possible.

Most effective if values are long (bit length k is large).

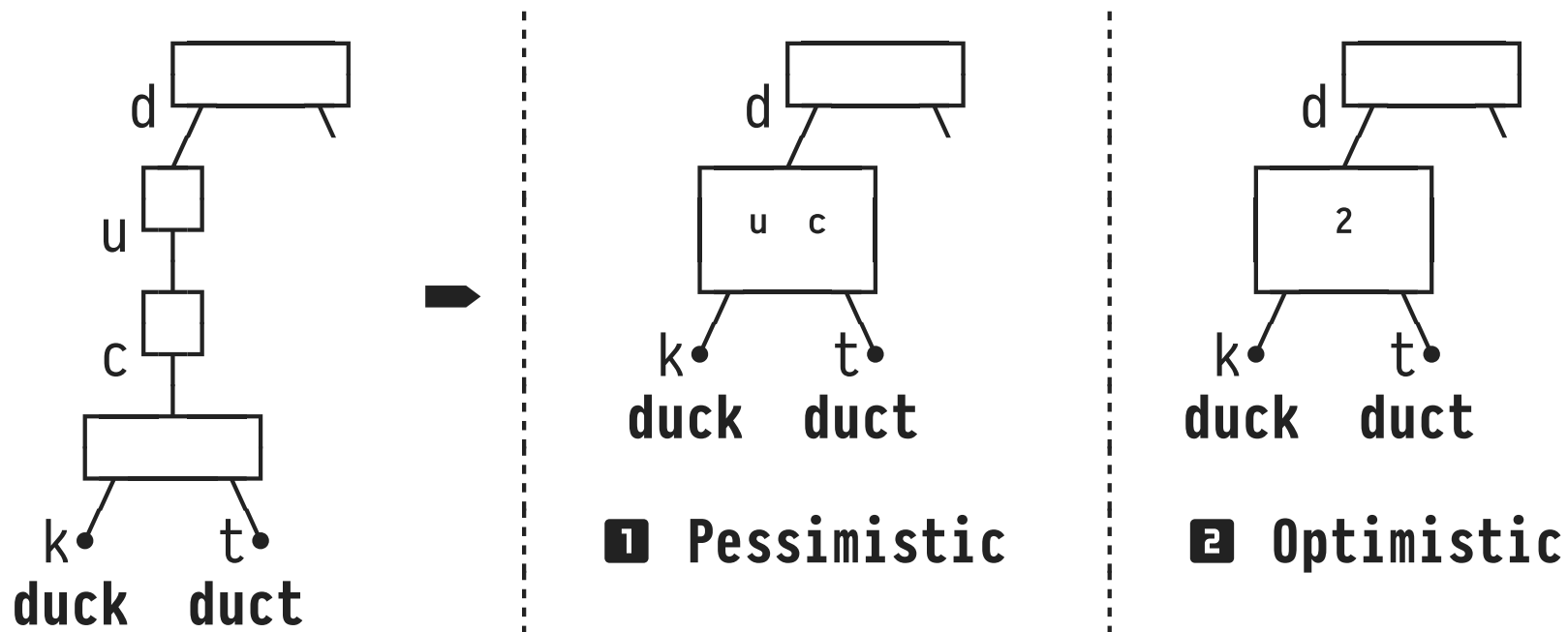
1. **Lazy expansion:** create inner nodes only if they are required to distinguish at least two leaf nodes.



- Saves two inner nodes on the path to leaf **foo**.
- Expand path if another leaf with prefix **f** is added.
- Path does not spell **foo**: store at leaf or in table.

ART Optimizations: Lazy Expansion + Path Compression

2. **Path compression:** remove inner nodes that have a single child.



- **1** Inner node holds byte sequence (here: **uc**) of preceding single-child nodes. Compare this sequence against **key**.
- **2** Inner holds count (**2**) of removed single-child nodes. Skip that many bytes in **key**. At a leaf, compare its value to **key** to ensure that traversal indeed reached the proper leaf.

ART: Index Lookup (Lazy Expansion + Pessimistic Path Compression)

<pre> search(<i>node</i>, <i>key</i>, <i>depth</i>): 1 if <i>node</i> = NULL 2 return NULL 3 if <i>isLeaf</i>(<i>node</i>) 4 if <i>leafMatches</i>(<i>node</i>, <i>key</i>) 5 return <i>node</i> 6 return NULL 7 <i>s</i> ← <i>byteSequence</i>(<i>node</i>) 8 if <i>s</i> ≠ <i>key</i>[<i>depth</i>...<i>depth</i>+<i>len</i>(<i>s</i>)] 9 return NULL 10 <i>depth</i> ← <i>depth</i> + <i>len</i>(<i>s</i>) 11 <i>child</i> ← <i>findChild</i>(<i>node</i>, <i>key</i>[<i>depth</i>]) 12 return search(<i>child</i>, <i>key</i>, <i>depth</i>+1) </pre>	<pre> search failed (<i>findChild</i>() returned NULL) reached leaf level? reached the proper leaf? yes, success no, failure byte sequence in inner node does partial key match byte sequence? skip to next key byte key byte selects child descend into subtree </pre>
--	---

- Invoke with **search**(*node*: <ART root>, *key*: 'duck', *depth*: 0).
- Returns leaf node that holds rowid(s) or NULL.

5 : Making the Most of Indexes

DuckDB is an OLAP DBMS, optimized to scan entire tables. Basic index support is present—but there certainly is unused potential.

- The following predicates *could* be evaluated by [Index Scan](#):³

Predicate	Index Constellation	Index used by 🐧?
$c_1 \text{ } \textcircled{<} \text{ } val_1$	$t(c_1)$	👍 (if selective)
$c_1 \text{ } \textcircled{<} \text{ } val_1 \text{ AND } c_2 \text{ } \textcircled{<} \text{ } val_2$	$t(c_1, c_2)$ (composite)	👎
$c_1 \text{ } \textcircled{<} \text{ } val_1 \text{ OR } c_2 \text{ } \textcircled{<} \text{ } val_2$	$t(c_1)$ and/or $t(c_2)$	👎
$c_1 \text{ } \textcircled{<} \text{ } val_1 \text{ OR } c_2 \text{ } \textcircled{<} \text{ } val_2$	$t(c_1)$ and $t(c_2)$	👎
$c_1 \text{ } \text{LIKE 'patt_rn%'}$	$t(c_1)$	👎

Index support in DuckDB (as of version 1.4)

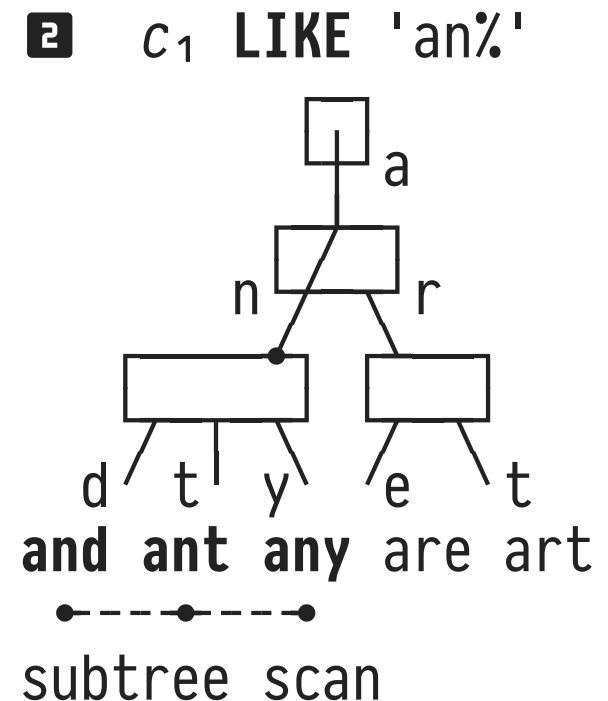
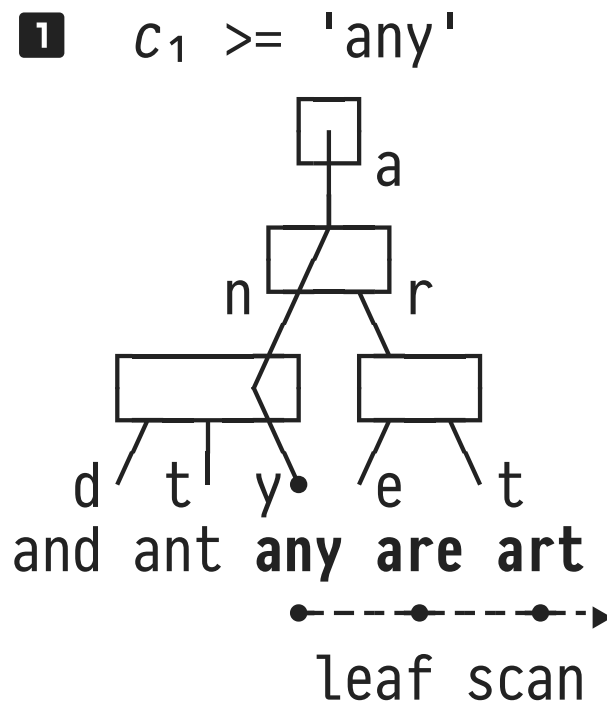
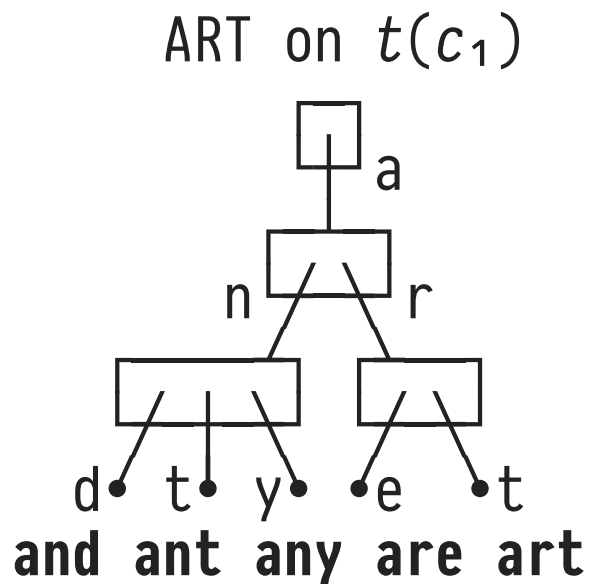
- (PostgreSQL🐘 is optimized to operate as an OLTP DBMS and can use indexes in all of these scenarios.)

³ In this table, $\textcircled{<}$ represents a SQL comparison operator: $<$, \leq , $=$, \geq , $>$.

Index Traversals

Index traversals (beyond equality-based index lookups for a key val) can support **range predicates**:  #019

- $c_1 \geq val_1$ or $c_1 \geq val_1$ **AND** $c_1 \leq val_2$: perform lookup \rightarrow for val_1 , then scan leaves left to right \dashrightarrow (until $> val_2$). **1**
- c_1 **LIKE** 'prefix%': perform lookup \rightarrow for $prefix$, hit inner node. Scan leaves in subtree \dashrightarrow . **2**



Design and Implementation of DuckDB Internals

⑥

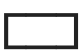




Query Execution Plans and Pipelining

April 7, 2026

Torsten Grust
Universität Tübingen, Germany

1 | Query Execution \equiv Plan Execution

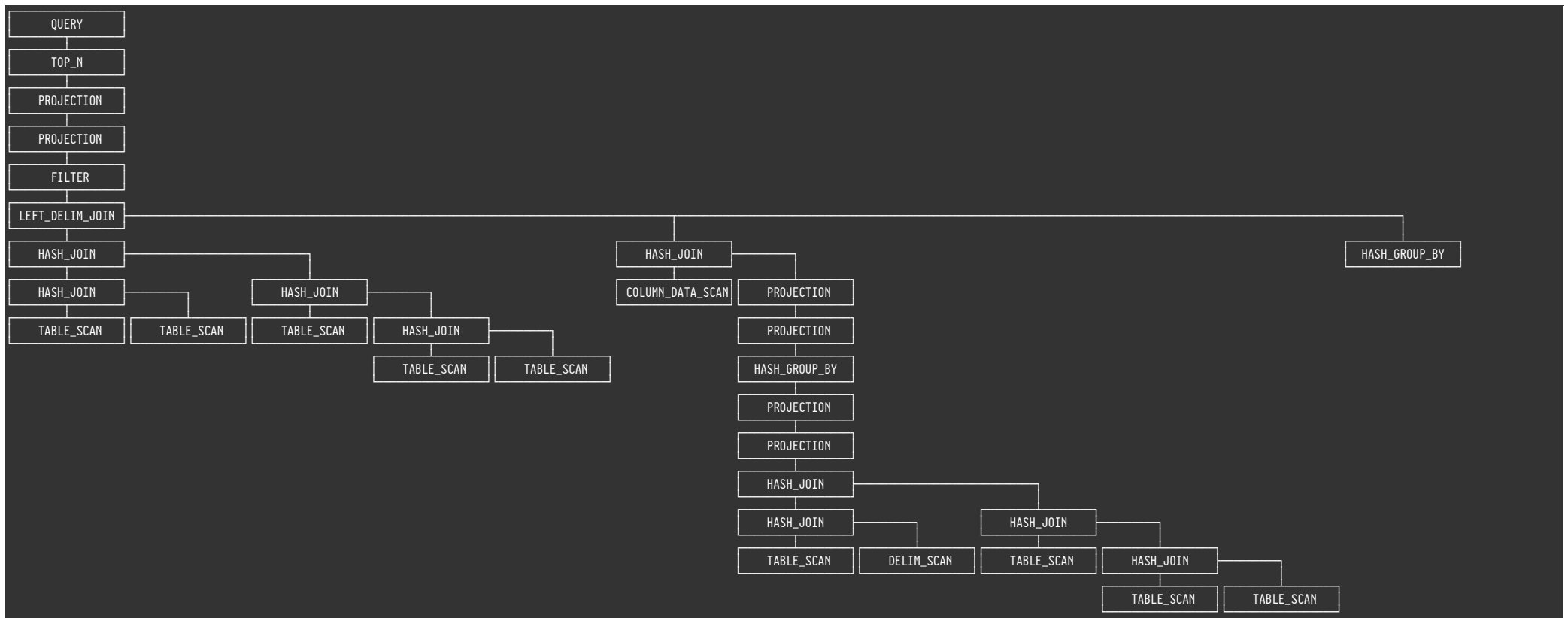
DuckDB translates SQL text into an **execution plan**, a tree-shaped¹  data flow network:

- Nodes  in this network are (physical) **operators**, each of which implement one computation step. Example operators:
 - **TABLE_SCAN**: reads rows from a table (leaf/source operator),
 - **PROJECTION**, **FILTER**: eval (Boolean) expressions, discard rows,
 - **HASH_GROUP_BY**: hash-based grouped aggregation.
- Edges are directed upwards/leftwards ( \rightarrow  \equiv  \leftarrow ) and route rows from child operator(s) to parent operator.
 - DuckDB: **data chunks** of 2048 rows flow jointly along.
- Rows returned by the root operator represent the query result.

¹ This is mostly true. Particular SQL features—e.g., correlated subqueries or common table expressions (CTEs)—can lead to DAG-shaped execution plans.

Execution Plans can be Complex

The execution plans of complex queries are intricate (TPC-H Q2):





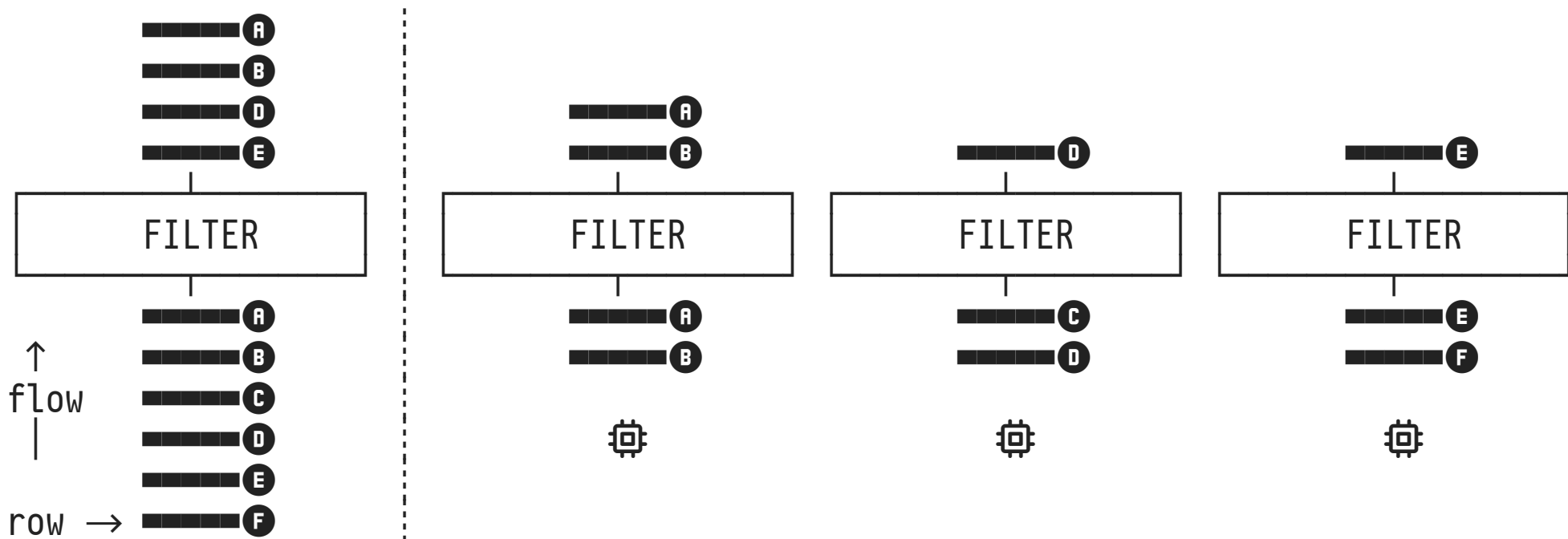
- **Q:** How to orchestrate the data flow such that DuckDB
 - does *not* need to materialize intermediate results and
 - and that all CPU cores 🏠 can equally contribute in parallel?

📄 #020

2 : Trivially Parallel Operators 1

Operators like **FILTER** (discard rows based on predicate) or **PROJECTION** (evaluate scalar expressions) are inherently parallel:

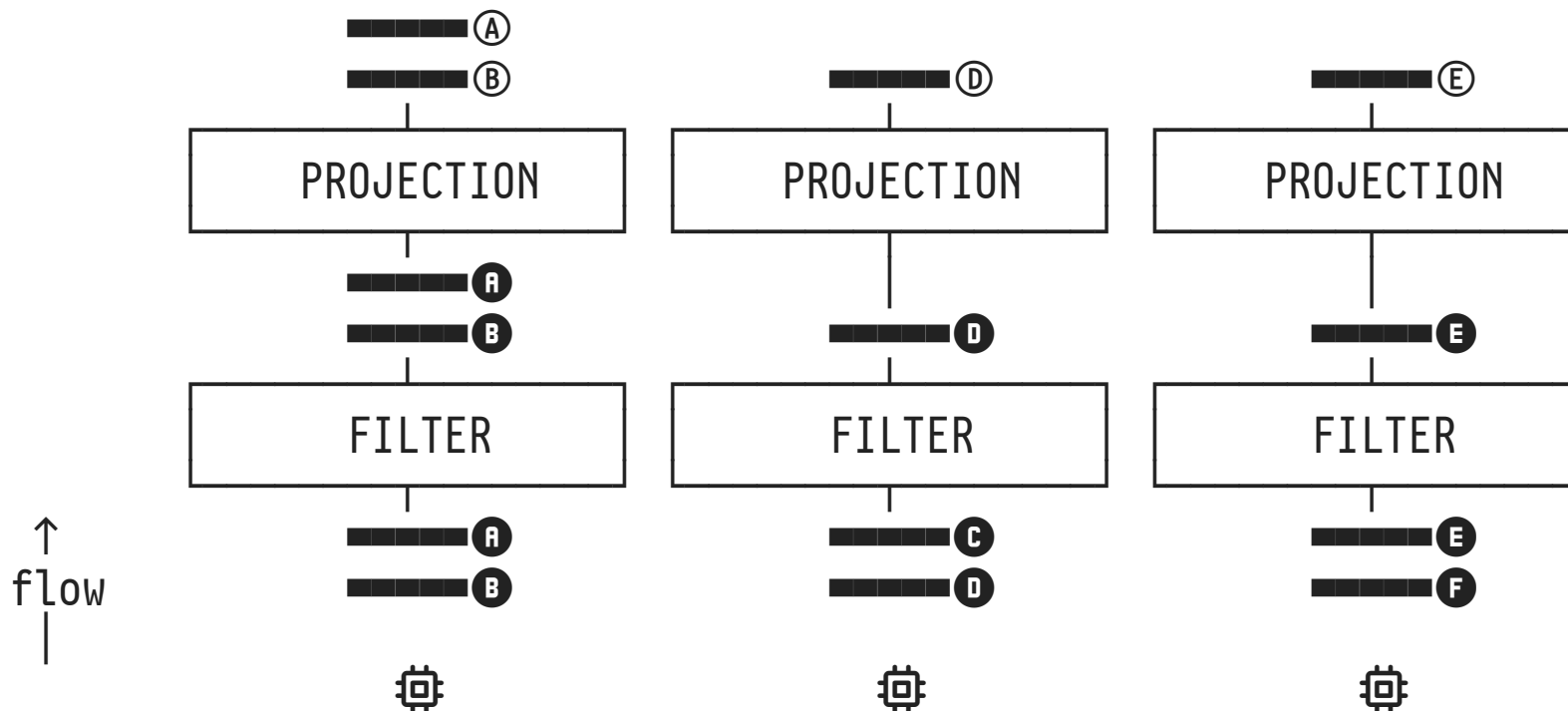
- Each thread  operates on a partition of rows. Overall result is the (disjoint) union of all -local results.
- No inter-thread dependencies, operator implementation itself does need *not* to be aware of //ism.



Trivially Parallel Operators

Sequences of “PROJECTION-like” operators are assembled into **pipelines**:



- Each thread/core    runs its own instance of the pipeline.
- Data is passed between operators in **chunks** (2048 rows).²



² Rows that pass **FILTER** are buffered before they are pushed down the pipeline (do not pass “single-row chunks”—DuckDB aims to pass more than 64 rows). If *all* rows pass, simply pass on the input rows.

Pipeline-Breaking Operators (Sinks)

“HASH_GROUP_BY-like” (DB lingo: **sinks**) operator phases:


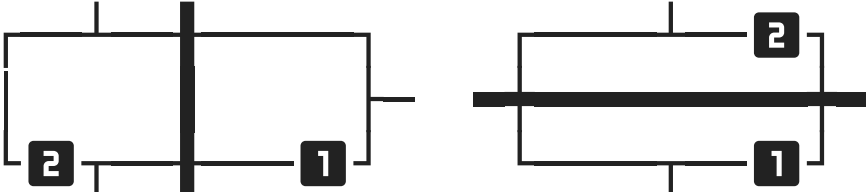
1. **Sink:** Each thread  “sinks” all incoming data into its *local state* (e.g., `HASH_GROUP_BY` Phase ①, recall Chapter 03).
2. **Combine:** Once a thread has read its input, combine local state with the operator's *global state* (`HASH_GROUP_BY`, Phase ②).
3. **Finalize:** *Combine* done, a single thread  post-processes the global state which is then pushed to the next pipeline.

Sinks may only be placed at pipeline ends (**pipeline breakers**).

- Examples of sinks in DuckDB:
 - `HASH_GROUP_BY` (builds aggregate hash table)
 - Build side of `HASH_JOIN` (hash table for rhs join input)
 - `ORDER_BY` (local state: sorted run)
 - `UNGROUPED_AGGREGATE` (local state: partial aggregate)

4 : Assembling Execution Plans from Pipelines

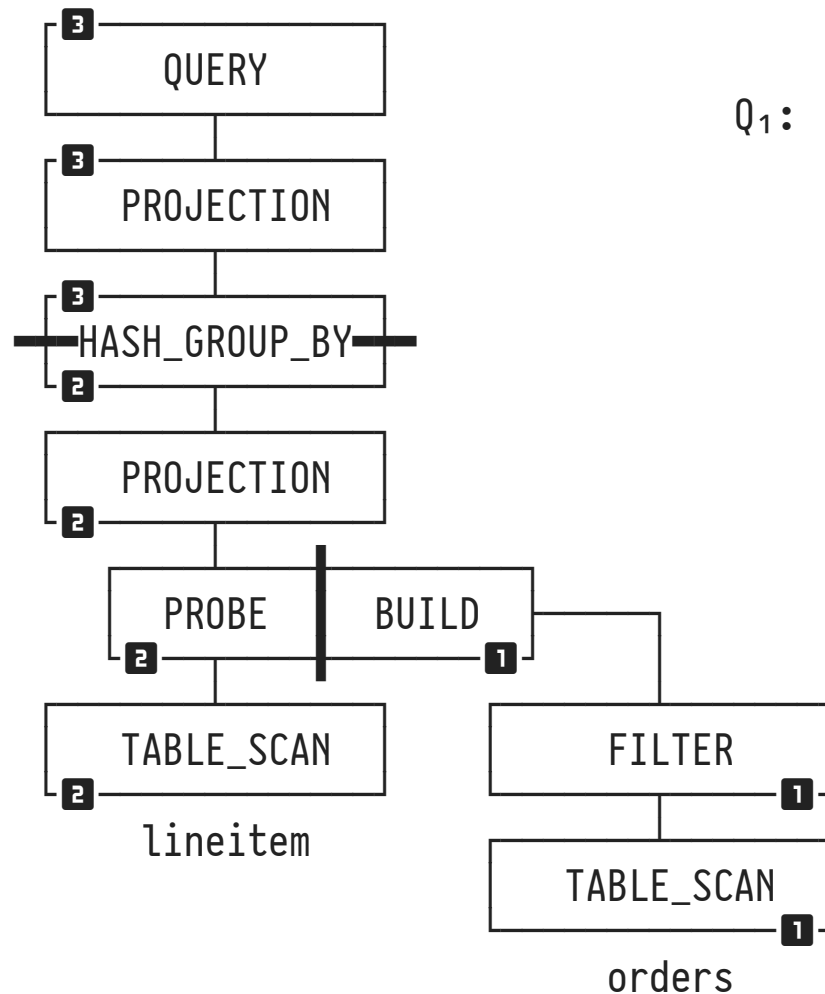
Execution plans for non-trivial SQL queries are assembled by **stitching multiple pipelines together**.

- In the plans below, operators  are tagged by **1**, **2**, ... to assign them to a pipeline. Pipeline breakers are marked using **†**.
- Operators like those on the right act like a sink for pipeline **1**. Their output can then be read by the subsequent pipeline **2**.
 
- Such operators introduce **pipeline dependencies: $1 < 2$** (pipeline **1** must finalize before **2** can read its output).
- Root operator QUERY acts as a sink that
 - collects all incoming rows and
 - constructs the final query result (to be shipped to the caller or displayed by the CLI).

Assembling Execution Plans from Pipelines

Sample query Q_1 over TPC-H data of medium complexity:³

 #021

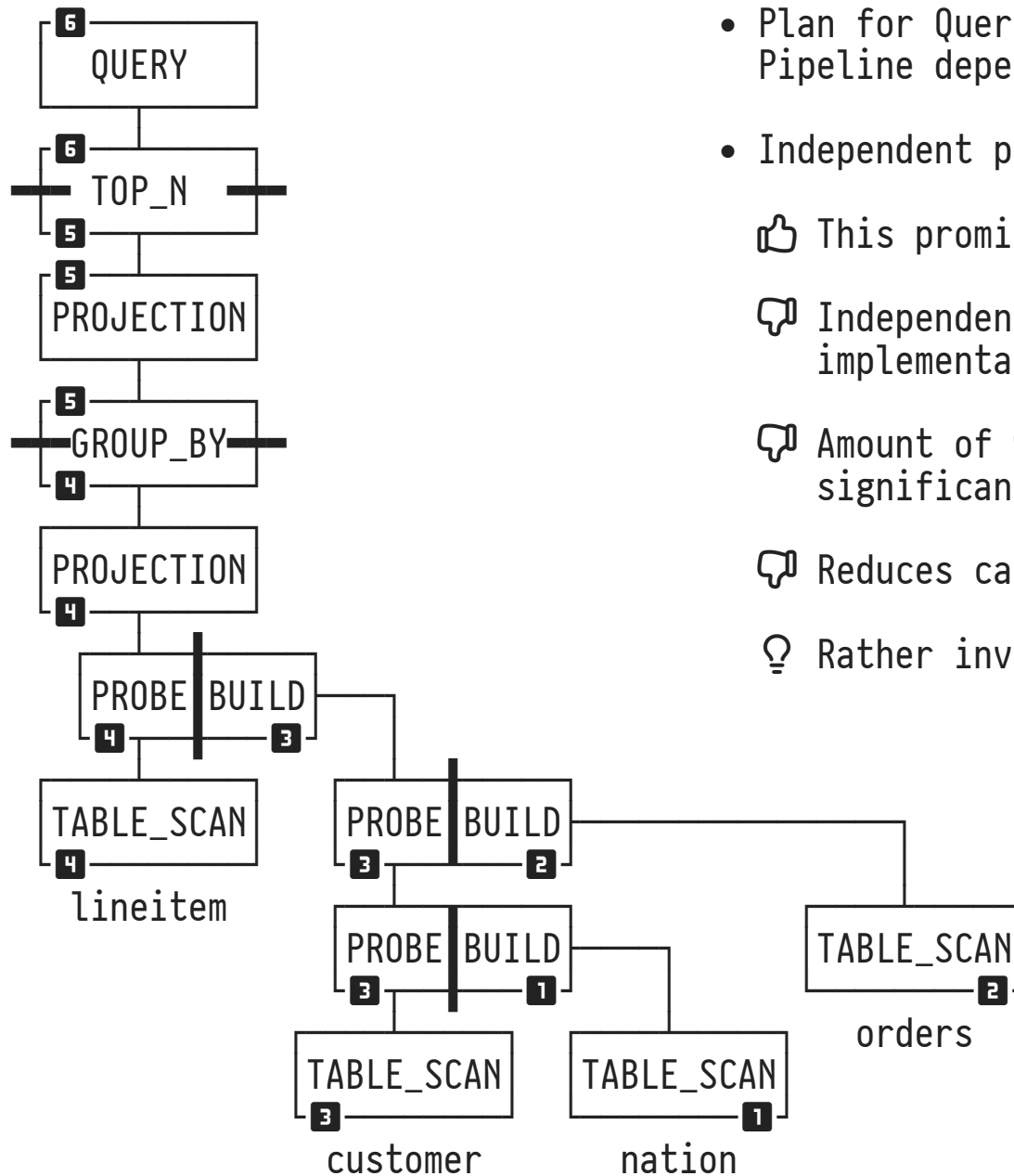


Q_1 : `SELECT DISTINCT o.o_orderkey AS violation
FROM lineitem AS l, orders AS o
WHERE l.l_orderkey = o.o_orderkey
AND o.o_orderstatus IN ('0', 'F')
AND l.l_linestatus <> o.o_orderstatus;`

- `PROBE | BUILD` \equiv `HASH_JOIN`
- For `HASH_JOIN`, DuckDB chooses the smaller input to be on the `BUILD` side.
- Pipeline dependencies:
`1 < 2 < 3`

³ Checks for violations of a TPC-H constraint: `o_orderstatus` is set to '0' ('F') iff all lineitems of this order have `l_linestatus` set to '0' ('F'). See the [TPC-H benchmark specification](#) , §4.2.3.

Assembling Execution Plans from Pipelines

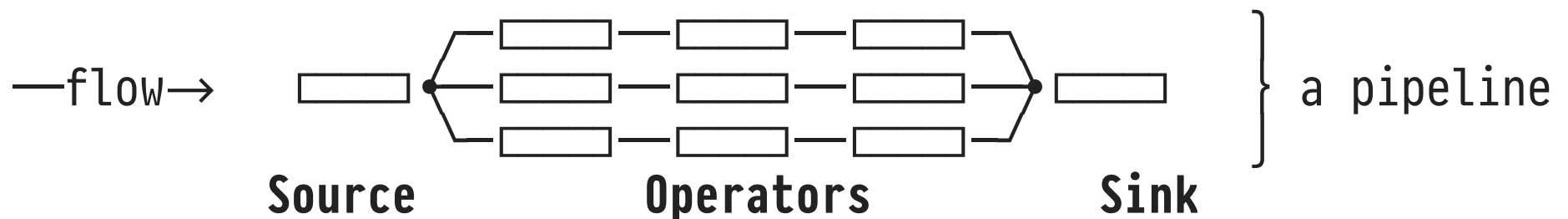


- Plan for Query Q₃ (TPC-H Query 10)
Pipeline dependencies: (1||2)<3<4<5<6.
- Independent pipelines 1||2 can be run in parallel:
 - 👍 This promises reduced latency.
 - 🗨 Independent pipelines are rather rare. Does the implementation of the required logic pay off?
 - 🗨 Amount of work in pipelines 1 and 2 may differ significantly.
 - 🗨 Reduces cache locality.
 - 💡 Rather invest more cores in a single pipeline.

5 : Parallelism in a Pipeline

A pipeline in a DuckDB execution consists of three segments.

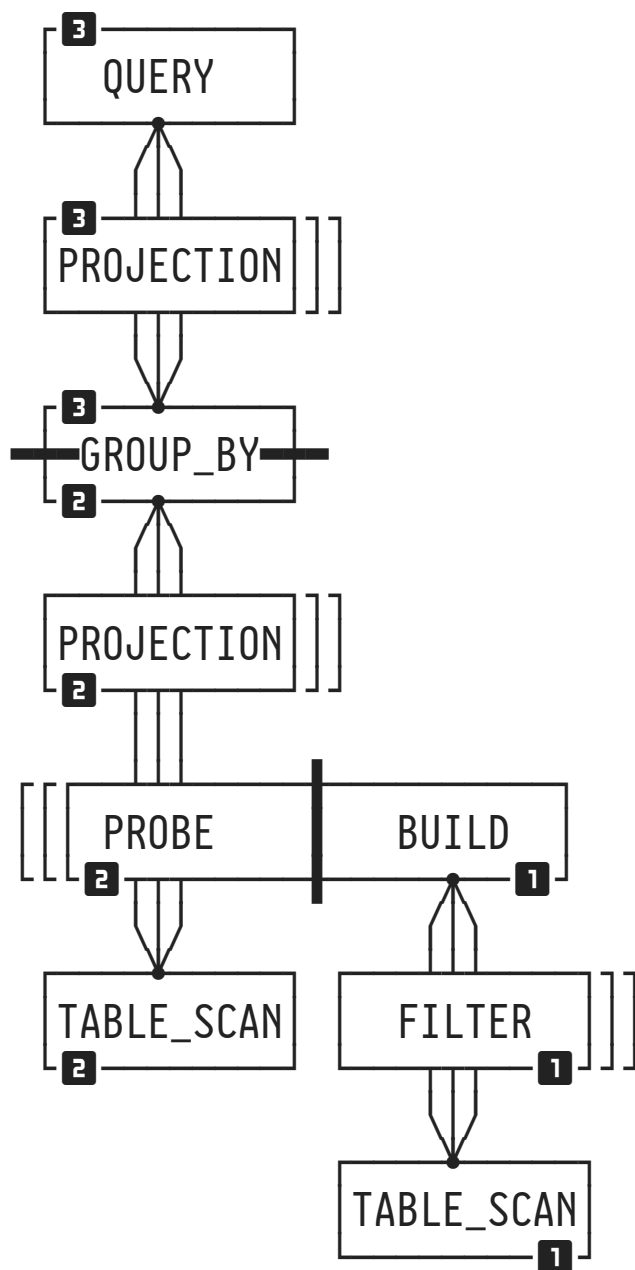
- **Sources** and **sinks** are parallelism-aware, **operators** need not be:



- Sources \leftarrow know how to partition rows. Examples:
 - `COLUMN_DATA_SCAN` (literal data), `RANGE` (row generator)
 - `TABLE_SCAN` (read base table), `CTE_SCAN` (read CTE)
 - `READ_CSV`, `PARQUET_SCAN` (read external file)
 - Some operators are sinks *and* sources: $\rightarrow \square \leftarrow$
(`HASH_GROUP_BY` serves rows from its aggregate hash table).


Parallelism in a Pipeline


01 #022



- Pipelines and parallelism in the plan for Q₁ (TPC-H constraint check, see above).
- Pipeline **sources** ←:
 - 1: TABLE_SCAN (orders)
 - 2: TABLE_SCAN (lineitem)
 - 3: Phase 2 of HASH_GROUP_BY
- Pipeline **sinks** →:
 - 1: BUILD phase of HASH_JOIN
 - 2: Phase 1 of HASH_GROUP_BY
 - 3: QUERY root node
- PROBE side of HASH_JOIN is an operator: multiple threads can peek into the hash table in parallel.

6 : Pipeline Execution

In DuckDB, **pipeline execution** is driven by a core loop (see next slide) that is run by each thread . This loop

- maintains operator state (-local as well as global),
- receives data chunks from operators (or the source) and **pushes** these chunks towards the downstream operator (or the sink),
- contains control flow that detects whether
 - 1 the source is exhausted or
 - 2 an operator outputs an empty chunk.

The pipeline driver can

- *cache* rows in case a **FILTER** or **PROBE** returns few results (avoid the overhead of passing on tiny or single-row chunks),
- *split* the data flow and pass chunks to multiple consumers (e.g., to share the output of a **TABLE_SCAN**).

Sketch of DuckDB's Pipeline Driver (Run by each Thread 🛠️)⁴

```

var operator[] pipeline[0...P]           pipeline to execute (operator count P+1)
var state[]    states[0...P]             🛠️-local + global operator state
var data_chunk[] intermediates[0...P-1]  rows output by source/operators (intermediate results)

source ← pipeline[0]                     pipeline source and sink
sink   ← pipeline[P]                     (pipeline[1...P-1] are operators)

for i in 0...P:                           initial operator states
| states[i] ← pipeline[i].GetOperatorState(...)
for i in 0...P-1:                          allocate intermediate chunks
| intermediates[i] ← data_chunk.Initialize(pipeline[i].return_type)

while true:
| intermediates[0] ← source.GetData(states[0], ...)  get next data chunk from 🛠️'s morsel
| if (intermediates[0].size = 0)                    no more rows? ❶
| | return FINISHED                                ↳ yes, bail out

reached_sink ← true
for i in 1...P-1:                                  execute operators in pipeline order
| intermediates[i] ← pipeline[i].Execute(intermediates[i-1], states[i])
| if (intermediates[i].size = 0)                    operator reduced data chunk to size 0? ❷
| | reached_sink ← false                            ↳ yes, restart at source with next chunk
| | break

if reached_sink:                                  did non-empty data chunk reach the sink?
| sink.Sink(intermediates[P-1], states[P], ...)    ↳ yes, collect rows at the sink

```


⁴ The pseudo code for the pipeline driver is modelled after [DuckDB GitHub issue 1583](#) 🐭.

Pipeline Driver Manages Control Flow, Operators are Simple

Since control flow is handled by the pipeline driver, source/operator/sink **implementations** turn out to be simple:

1. Build phase of `HASH_JOIN` (sink \rightarrow):

```
HashJoinBuild.Sink(input, state):  
| ht ← state.local_hash_table  
| ht.BuildHashTable(input)
```

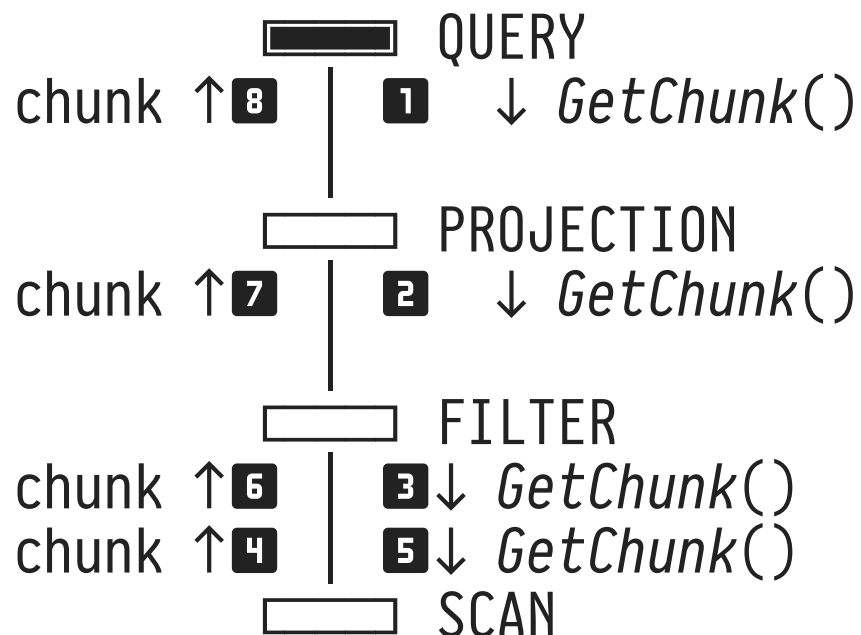
The *Combine* phase of `HASH_JOIN` merges the -local hash tables into a global hash table that can be probed against.


2. Probe phase of `HASH_JOIN` (operator):



```
HashJoinProbe.Execute(input, state):  
| ht ← state.global_hash_table  
| joined ← ht.ProbeHashTable(input)  
| return joined
```

Not DuckDB: Pull-Based Execution (Volcano Iterator Model)

Prior to release 0.3.0 (October 2021), DuckDB implemented **pull-based** operators⁵ without a pipeline driver:



1. Start at plan root . Call *GetChunk()* on child operator to request next data chunk.
2. Operators forward *GetChunk()* call down the plan until leaves can return a chunk.
3. Operators process obtained chunk. Call *GetChunk()* again if all of chunk was filtered, otherwise pass chunk upwards.
4. Operators return FINISHED if no more chunks to deliver.

⁵ This pull-based plan execution strategy is known as the [Volcano iterator model](#) . It is still widely implemented in today's DBMSs. Example: PostgreSQL  (operators return a *single row* on each “*GetChunk()*” call).

Not DuckDB: Pull-Based Execution (Volcano Iterator Model)

Pseudo code for `HASH_JOIN` in the pull-based Volcano model:

- Control flow handled inside operator.
- Probe + build phases are entangled, hard to parallelize. 🗨️

```

HashJoin.GetChunk(probe, build):
  state ← this.GetOperatorState(...)
  ht ← state.hash_table

  if not state.build_done:
    while true:
      chunk ← build.GetChunk(...)
      if chunk.size = 0:
        break
      ht.BuildHashTable(chunk)
      state.build_done ← true

  do:
    chunk ← probe.GetChunk(...)
    if chunk.size = 0:
      return FINISHED
    joined ← ht.ProbeHashTable(chunk)
  while joined.size = 0
  return joined

```

- are we in the build or probe phase?
- 1 build phase
 - pull next chunk from build side
 - no more rows?
 - ↳ yes, hash table complete—now probe
 - insert chunk into hash table
 - 2 probe phase
 - pull next chunk from probe side
 - no more rows?
 - ↳ yes, join complete
 - probe chunk against hash table
 - retry if no rows joined

Design and Implementation of DuckDB Internals

⑦

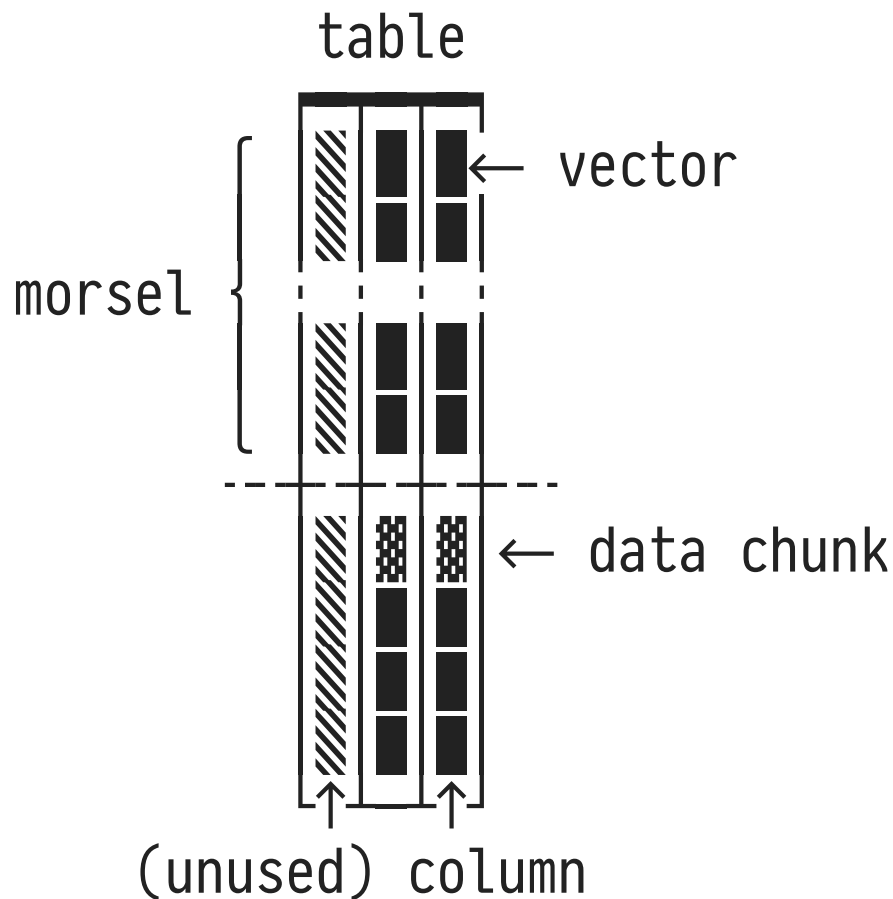
Vectorized Query Execution






April 7, 2026

Torsten Grust
Universität Tübingen, Germany

1 : DuckDB Terminology Recap

Before we proceed, let us make sure that we agree on the data items processed by DuckDB's query engine:



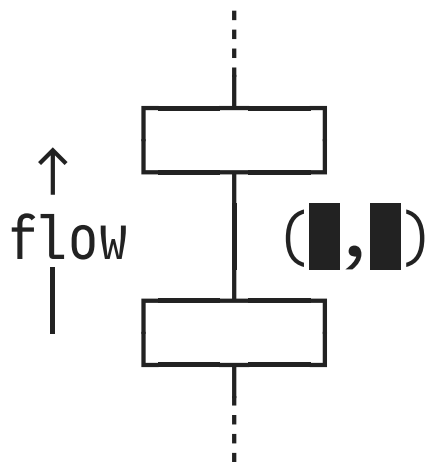
- **Tables** are organized in **columns** some of which may be unused  in any given plan operator.
- Columns are split into **vectors** , each containing 2048 values.¹
- A **data chunk** (, ) groups vectors of multiple columns, representing a horizontal slice of a table.
- Threads  consume input data **morsel**-by-**morsel**, each consisting of 60 data chunks (= 122,880 rows).

¹ DuckDB's vector size can be configured at compile time ([STANDARD_VECTOR_SIZE](#)).

2 : Passing Data Chunks Between Operators

During plan execution, operators process **one data chunk at a time**, then pass that chunk downstream:

Memory Location	Latency ⌘	Size (🍏 M2 Max)	... Can Hold
CPU Registers		$n \times 64$ bits	cell value(s)
L1 Cache	< 1 ns	(per 🏠) 192 KB	vector
L2 Cache	2.8 ns	(shared) 32 MB	data chunk
RAM 🏠	≈ 100 ns	16–96 GB	column(s)



- The passed **data chunks** will fit in the L2 cache. Operators can access intermediate results with low latency. 📄 #023
- Passing **entire columns**: many intermediates will only fit into RAM. High latency. 👍
- Passing **single rows**: frequent context switches between operators. CPU overhead. 👎

3 : DuckDB Vectors: Logical vs. Physical

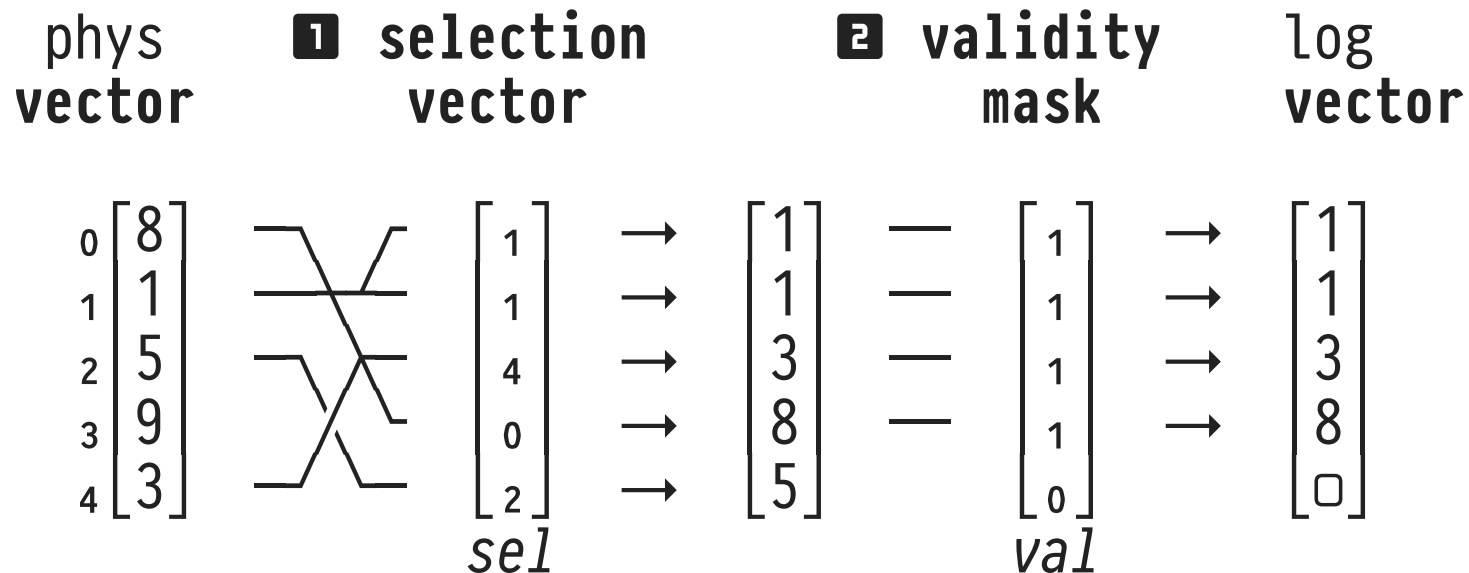
To form a horizontal table slice, data chunks (■,...,■) contain a **vector** for each relevant column.

- Vector ■ (**logical**): sequence of 2048 values of a *single type*.
- DuckDB exploits regularity in value sequences to implement a family of compressed **physical vector representations**:

FLAT		CONSTANT		DICTIONARY		SEQUENCE	
$\begin{bmatrix} 1 \\ 9 \\ 0 \\ 4 \\ 2 \end{bmatrix}$	$\begin{bmatrix} 1 \\ 9 \\ 0 \\ 4 \\ 2 \end{bmatrix}$	$\begin{bmatrix} 42 \\ 42 \\ 42 \\ 42 \\ 42 \end{bmatrix}$	$\begin{bmatrix} 42 \\ 42 \\ 42 \\ 42 \\ 42 \end{bmatrix}$	$\begin{matrix} 0 \\ 1 \end{matrix} \begin{bmatrix} DE \\ NL \end{bmatrix}$	$\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 1 \end{bmatrix}$	$\begin{bmatrix} DE \\ NL \\ DE \\ DE \\ NL \end{bmatrix}$	$\begin{bmatrix} 0 \\ 10 \\ 20 \\ 30 \\ 40 \end{bmatrix}$
		<i>const</i>		<i>dict</i>	<i>sel</i>		<i>base</i>
		$\begin{bmatrix} 5 \\ 5 \\ 5 \\ 5 \\ 5 \end{bmatrix}$	$\begin{bmatrix} 42 \\ 42 \\ 42 \\ 42 \\ 42 \end{bmatrix}$				$\begin{bmatrix} 10 \\ 10 \\ 10 \\ 10 \\ 10 \end{bmatrix}$
		<i>len</i>					<i>inc</i>
phys	log	phys	log	phys	log	phys	log

DuckDB Vectors: Selection Vectors and Validity Masks

Any vector is associated with a **selection vector** and **validity mask** used to permute/filter its value sequence and encode **NULLs** (\square):



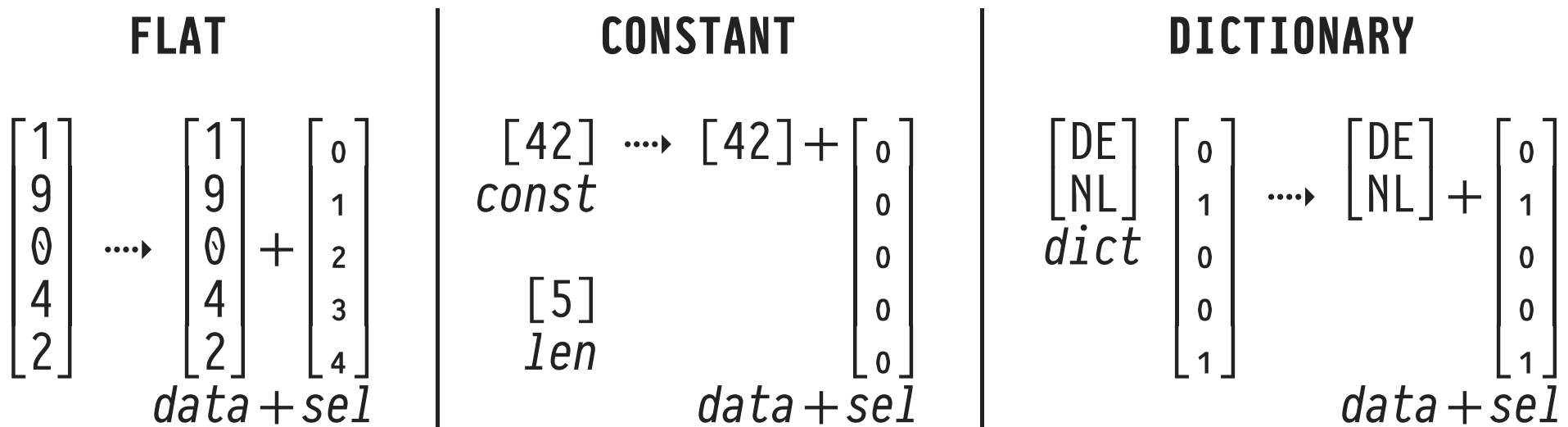
- Selection vectors contain vector indices $\in \{0, \dots, 2047\}$. Can help to avoid copying—potentially sizable—values after a sorting or filtering operation.
- Validity masks contain Booleans (0/1).

DuckDB Vectors: Unified Representation

DuckDB aims to push vectors in compressed physical representation from operator to operator. However, *handling all representations everywhere* would lead to combinatorial code explosion. 🗨️

- 💡 Convert to (\dots) and then process a **unified representation** comprising a *data vector+selection vector* pair.

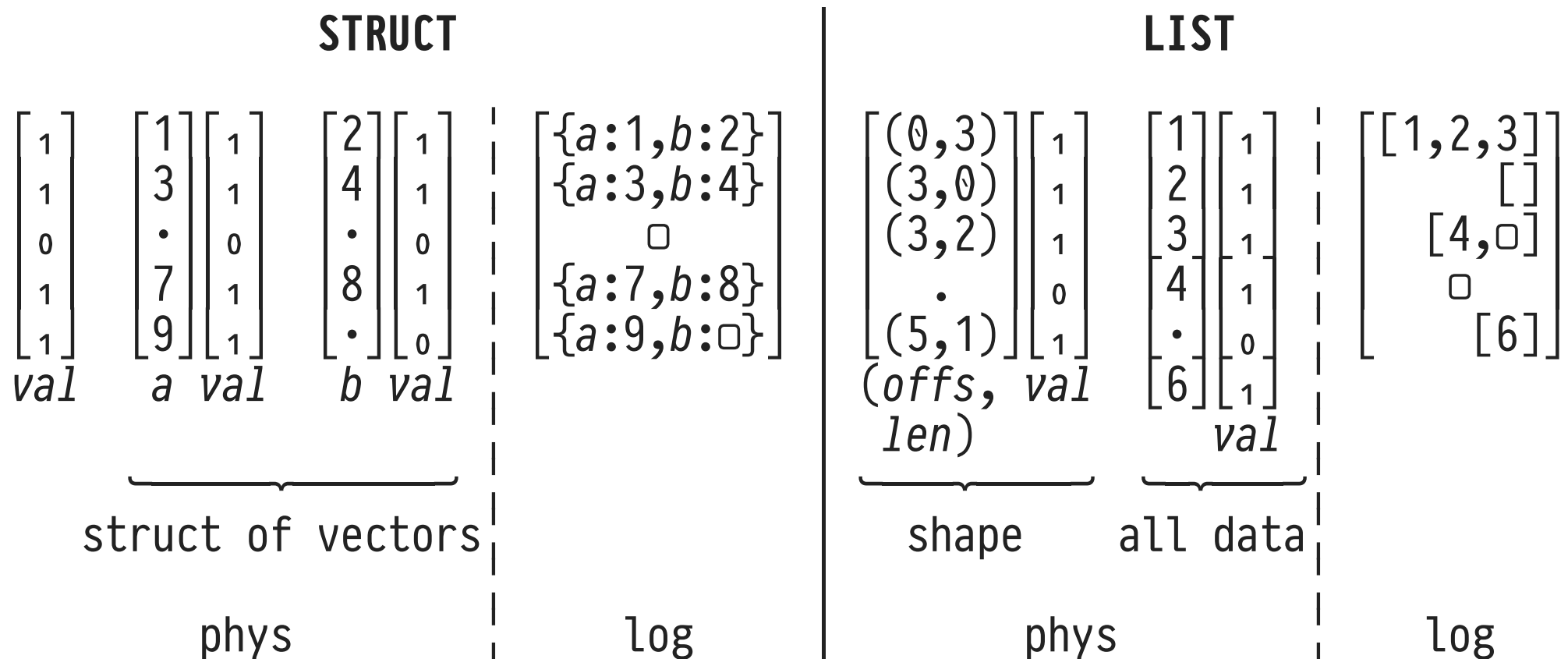
NB. This conversion does *not* involve value copying:²



² SEQUENCE vectors are first expanded into FLAT vectors, then converted into the unified representation.

DuckDB Vectors of Complex Values (Structs {...}, Lists [...])

- **Vectors of structs** are represented as structs of vectors.
- **Vectors of lists** store list shapes and data values separately, combine data of all lists into a single data vector.



- Nested types: apply representation schemes recursively.

4 : Super-Specific Vector Code

Q: How does DuckDB perform a binary operation on vectors `l` and `r`?

That depends on

- **1** the kind of operation (e.g., arithmetics: `+`, `*`, ...),
- **2** value types (`int`, `float`, ...),
- vector representations of **3** `l` and **4** `r` (`FLAT`, `CONST`, ...).

💡 Branch on **1**...**4** *before* we enter the tight loop that scans the values in `l` and `r`:

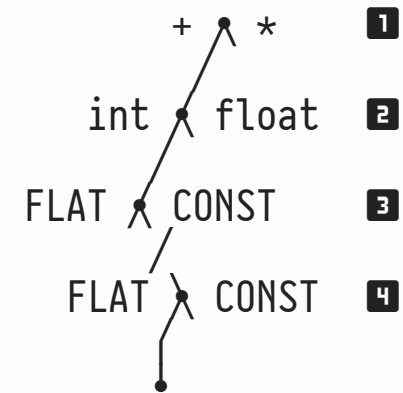
- Avoids repeated tests (with identical outcome),
- leads to *tight inner loops* that feature less branches.

Super-Specific Vector Code

```

switch (1 op)
  case +: switch (2 type)
    case int: switch (3 l.repr)
      case FLAT: switch (4 r.repr)
        case FLAT:
          res.repr = FLAT;
          for (i = 0; i < count; i++)
            [ res.data[i] = l.data[i] + r.data[i];
              break;
        case CONSTANT:
          res.repr = FLAT;
          for (i = 0; i < count; i++)
            [ res.data[i] = l.data[i] + r.const;
              break;
        ...
      case CONSTANT: switch (4 r.repr)
        case FLAT:
          res.repr = FLAT;
          for (i = 0; i < count; i++)
            [ res.data[i] = l.const + r.data[i];
              break;
        case CONSTANT:
          res.repr = CONSTANT;
          res.const = l.const + r.const;
          break;
        ...
    case float: ...
  case *: ...
  ...

```



Super-specific code section:
Add the constant vector l of integers to flat vector r



Super-Specific Vector Code

Previous slide does not handle **selection vector/validity** mask (each leaf of the decision tree λ requires code much like this):

```
for (i = 0; i < count; i++)      /* code for leaf 1: +, 2: int, 3: FLAT, 4: FLAT */
| lindex = l.sel[i];
| rindex = r.sel[i];
| if (l.val[lindex] && r.val[rindex])
| | res.data[i] = l.data[lindex] + r.data[rindex];
| | res.val[i] = 1;
| else
| | res.val[i] = 0;
```

DuckDB aims to control the resulting C++ code volume:


1. Convert vectors l , r to **unified representation**, process these by generic loops (cuts lower levels of λ).
2. Use **C++ templates** to generate code at *DBMS compile time*.³

³ A different breed of DBMS—for example [Umbra](#)  —generates and compiles C++/IR/machine code at *query translation time*. These query-specific code fragments are then linked with the DBMS and invoked at runtime.

Deep Dive Into DuckDB C++ Code: Vector Operations ①

```
SELECT i = 42          -- i: SEQUENCE (base=1, inc=1), 42: CONST
FROM   generate_series(1, 100) AS _(i)
```

Trace how DuckDB vectorizes the evaluation of expression `i = 42`:

ExpressionExecutor::Execute(, ..., count, result)

```
void ExpressionExecutor::Execute(
  const BoundComparisonExpression &expr, ..., idx_t count, Vector &result) {
  ...
  Execute(*expr.left, ..., count, left);    ↪ left  ≡ [1|1]s    (SEQUENCE)
  Execute(*expr.right, ..., count, right);  ↪ right ≡ [42]c    (CONSTANT)

  switch (expr.GetExpressionType() ①) {
  case ExpressionType::COMPARE_EQUAL:
    VectorOperations::Equals(left, right, result, count);
    break;
  case ExpressionType::COMPARE_NOTEQUAL:
    VectorOperations::NotEquals(left, right, result, count);
    break;
  ...
  }
}
```

Deep Dive Into DuckDB C++ Code: Vector Operations ②

VectorOperations::Equals([1|1]^s, [42]^c, result, count)

```
void VectorOperations::Equals(
    Vector &left, Vector &right, Vector &result, idx_t count) {
    ComparisonExecutor::Execute<duckdb::Equals>(left, right, result, count); ①
}
```

```
struct ComparisonExecutor {
private:
    template <class T, class OP>
    static inline void TemplatedExecute(
        Vector &left, Vector &right, Vector &result, idx_t count) {
        BinaryExecutor::ExecuteSwitch<T, T, bool, OP>(left, right, result, count); ③
    }
public:
    template <class OP>
    static inline void Execute(
        Vector &left, Vector &right, Vector &result, idx_t count) {
        ...
        switch (left.GetType().InternalType() ②) {
        ...
        case PhysicalType::INT64:
            TemplatedExecute<int64_t, OP>(left, right, result, count); ④
            break;
        ...
        }
    }
}
```

Deep Dive Into DuckDB C++ Code: Vector Operations ③

```
BinaryExecutor::ExecuteSwitch<int64_t, int64_t, bool, ..., duckdb::Equals, ...> (
    [1|1]s, [42]c, result, count, ...)
```

```
struct BinaryExecutor {
:
    template <class LEFT_TYPE, class RIGHT_TYPE, class RESULT_TYPE, ..., class OP, ...>
    static void ExecuteSwitch(
        Vector &left, Vector &right, Vector &result, idx_t count, ...) {
        auto left_vector_type = left.GetVectorType();           ↗ SEQUENCE_VECTOR
        auto right_vector_type = right.GetVectorType();        ↗ CONSTANT_VECTOR
        if (left_vector_type == VectorType::CONSTANT_VECTOR ③ &&
            right_vector_type == VectorType::CONSTANT_VECTOR ④) {
            ExecuteConstant<LEFT_TYPE, RIGHT_TYPE, RESULT_TYPE, ..., OP, ...>(
                left, right, result, ...);
        }
        :
        else {
            ExecuteGeneric<LEFT_TYPE, RIGHT_TYPE, RESULT_TYPE, ..., OP, ...>(
                left, right, result, count, ...);
        }
    }
:
}
```

DuckDB 1.4 source: `src/include/duckdb/common/vector_operations/binary_executor.hpp`

Deep Dive Into DuckDB C++ Code: Vector Operations ④

No specific code section for vector operation **SEQUENCE = CONSTANT**:
first transform **left** and **right** into **unified representation**.

```
BinaryExecutor::ExecuteGeneric<int64_t, int64_t, bool, ..., duckdb::Equals, ...> (
    [1|1]s, [42]c, result, count, ...)
```

```
struct BinaryExecutor {
:
    template <class LEFT_TYPE, class RIGHT_TYPE, class RESULT_TYPE, ..., class OP, ...>
    static void ExecuteGeneric(
        Vector &left, Vector &right, Vector &result, idx_t count, ...) {
        UnifiedVectorFormat ldata, rdata;

        left.ToUnifiedFormat(count, ldata);    ↪ ldata ≡ [1,2,3,...,100]u (UnifiedVector)
        right.ToUnifiedFormat(count, rdata);   ↪ rdata ≡ [42,42,...,42]u (UnifiedVector)

        result.SetVectorType(VectorType::FLAT_VECTOR);
        auto result_data = FlatVector::GetData<RESULT_TYPE>(result);
        ExecuteGenericLoop<LEFT_TYPE, RIGHT_TYPE, RESULT_TYPE, ..., OP, ...> (
            UnifiedVectorFormat::GetData<LEFT_TYPE>(ldata),
            UnifiedVectorFormat::GetData<RIGHT_TYPE>(rdata),
            result_data, ..., count, ...);
    }
:
}
```

Deep Dive Into DuckDB C++ Code: Vector Operations ⑤

Both argument vectors [...] are now in unified representation:

```
BinaryExecutor::ExecuteGenericLoop<int64_t, int64_t, bool, ..., duckdb::Equals, ...> (
    [1,2,3,...,100]u, [42,42,...,42]u, result, ..., count, ...)
```

```
struct BinaryExecutor {
:
:
template <class LEFT_TYPE, class RIGHT_TYPE, class RESULT_TYPE, ..., class OP, ...>
static void ExecuteGenericLoop(
    const LEFT_TYPE *__restrict ldata, const RIGHT_TYPE *__restrict rdata,
    RESULT_TYPE *__restrict result_data, ..., idx_t count, ...) {
:
    for (idx_t i = 0; i < count; i++) {
        auto lentry = ldata[i];
        auto rentry = rdata[i];
        result_data[i] = Operation<..., OP, LEFT_TYPE, RIGHT_TYPE, RESULT_TYPE> (
            ..., lentry, rentry, ...);
    }
}
:
}
```

DuckDB 1.4 source: `src/include/duckdb/common/vector_operations/binary_executor.hpp`

Deep Dive Into DuckDB C++ Code: Vector Operations ⑥

```
BinaryExecutor::ExecuteGenericLoop<int64_t, int64_t, bool, ..., duckdb::Equals, ...> (
    [1,2,3,...,100]u, [42,42,...,42]u, result, ..., count, ...)
```

```
struct BinaryExecutor {
:
:
template <class LEFT_TYPE, class RIGHT_TYPE, class RESULT_TYPE, ..., class OP, ...>
static void ExecuteGenericLoop(
    const LEFT_TYPE *__restrict ldata, const RIGHT_TYPE *__restrict rdata,
    RESULT_TYPE *__restrict result_data, ..., idx_t count, ...) {
    :
    for (idx_t i = 0; i < count; i++) {
        auto lentry = ldata[i];
        auto rentry = rdata[i];
        result_data[i] = lentry == rentry;
    }
}
:
}
```

DuckDB 1.4 source: `src/include/duckdb/common/vector_operations/binary_executor.hpp`

- `ldata/rdata`: plain C++ arrays of `LEFT_TYPE/RIGHT_TYPE` elements.
- `const ... *__restrict`: *Hey C++, we never write to these arrays and they do not overlap in memory. Go on, optimize!*

5 : Compiling Tight Loops

Tight loops over vectors form the core of DuckDB's query engine.

- Routine to subtract two flat integer vectors (⚠️ simplified):

```
#define STANDARD_VECTOR_SIZE 2048
```

```
#023
```

```
void PROJECT_sub_int_col_int_col(  
    int* col1, int* col2, int *res, int* sel)  
{  
    int i;  
  
    if (sel) {  
        // skip discussion of selection vector for now  
    }  
  
    for (i = 0; i < STANDARD_VECTOR_SIZE; i += 1) {  
        res[i] = col1[i] - col2[i];  
    }  
}
```

Assembly Code for Tight Loops

Use `clang` (options `-O2 -fno-vectorize -fno-unroll-loops`) to generate code for x86-64 CPUs.⁴

- Register assignment:

`col1: %rdi, col2: %rsi, res: %rdx, i: %rax`

```
PROJECT_sub_int_col_int_col:
    xorl %eax, %eax           # %rax ←32 0
loop:
    movl (%rdi,%rax,4), %ecx  # %ecx ←32 mem[%rdi + %rax × 4]
    subl (%rsi,%rax,4), %ecx  # %ecx ←32 %ecx -32 mem[%rsi + %rax × 4]
    movl %ecx, (%rdx,%rax,4)  # mem[%rdx + %rax × 4] ←32 %ecx
    incq %rax
    cmpq $2048, %rax         # 2048 loop iterations
    jne loop                 # exit if %rax = 2048
    retq
```

- **NB.** One loop exit test per vector element computed.

⁴ Use Matt Godbolt's  [Compiler Explorer](#) to inspect assembly code generated by contemporary compilers.

Explicit Loop Unrolling

- Manually perform **loop unrolling** to
 1. improve the ratio (*useful work*)/(*loop exit test*),
 2. expose independent work that may be executed in `//`:

```
void PROJECT_sub_int_col_int_col(int* col1, int* col2, int* res)
{
    int i;

    for (i = 0; i + 3 < STANDARD_VECTOR_SIZE; i += 4) {
        res[i] = col1[i] - col2[i];
        res[i+1] = col1[i+1] - col2[i+1];
        res[i+2] = col1[i+2] - col2[i+2];
        res[i+3] = col1[i+3] - col2[i+3];
    }
}
```

↓

independent, execute in
any order or even in `//`

#026

- **NB.** This code does not handle the case of
`STANDARD_VECTOR_SIZE mod 4 ≠ 0`.

#027

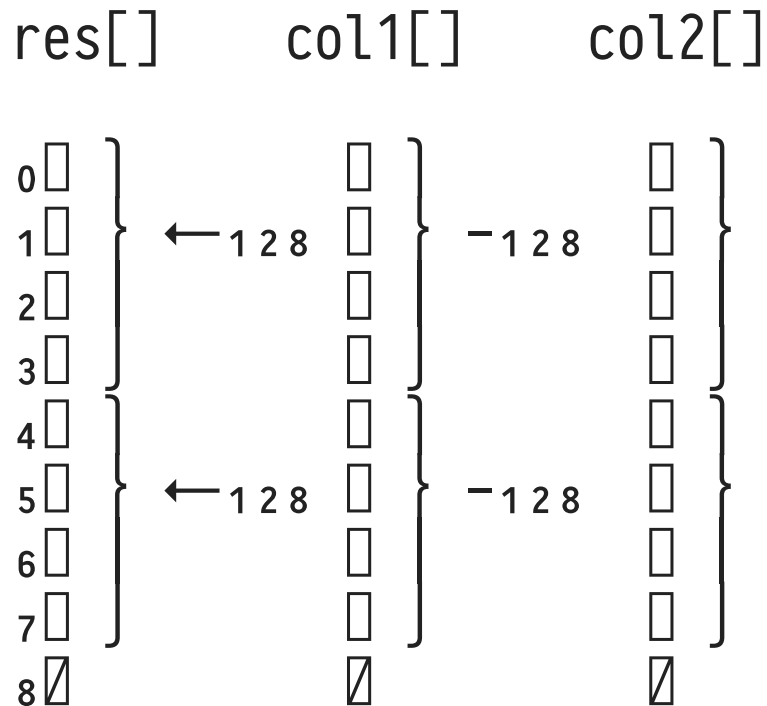
Loop Unrolling by the Compiler

Compiler `clang` (options `-O2 -fno-vectorize -funroll-loops`) unrolls four loop iterations (easy for CPU to //ize):

```
PROJECT_sub_int_col_int_col:
    xorl %eax, %eax                # i ←32 0
loop:
    movl (%rdi,%rax,4), %ecx       # %ecx ←32 col1[i]
    subl (%rsi,%rax,4), %ecx       # %ecx ←32 %ecx -32 col2[i]
    movl %ecx, (%rdx,%rax,4)       # res[i] ←32 %ecx
    movl 4(%rdi,%rax,4), %ecx      # %ecx ←32 col1[i+1]
    subl 4(%rsi,%rax,4), %ecx      # %ecx ←32 %ecx -32 col2[i+1]
    movl %ecx, 4(%rdx,%rax,4)      # res[i+1] ←32 %ecx
    movl 8(%rdi,%rax,4), %ecx      # :
    subl 8(%rsi,%rax,4), %ecx
    movl %ecx, 8(%rdx,%rax,4)
    movl 12(%rdi,%rax,4), %ecx
    subl 12(%rsi,%rax,4), %ecx
    movl %ecx, 12(%rdx,%rax,4)
    addq $4, %rax                 # } i ←32 i+4
    cmpq $2048, %rax              # } 2048 / 4 = 512 loop iterations
    jne loop                       # exit if %rax = 2048
    retq
```

} no control or data **dependencies** between these code blocks

Data-Parallelism Through SIMD⁵ Instructions



- Read/compute/write 4 vector elements (of width 4×32 bits = 128 bits) at a time in **data-parallel** fashion.
- Relies on CPU's SIMD support (e.g., Intel® SSE registers `%xmmi` and instruction `move double quad word`).

- **NB.** Requires care if
 - vectors `res[]` and `col1[]/col2[]` overlap in memory,
 - residual vector elements (see) are to be processed.

⁵ SIMD: Single-Instruction, Multiple Data. Support is vendor-specific, rapidly evolving in contemporary CPUs. Intel's AVX2 registers: 512 bits, ARM NEON registers: 128 bits. Compilers implement CPU architecture flags.

Data-Parallelism Through SIMD Instructions (Vector Overlap)

C compiler `clang` (options `-O2 -fvectorize`) uses SIMD registers and instruction set.

- Extra prelude code checks for **vector overlap** on function entry. If so, jumps to non-vectorized (yet unrolled) version of code.
- Declare function arguments with modifier `__restrict__` to inform C compiler that vectors do not overlap (⚠ omits all checks):

```
void PROJECT_sub_int_col_int_col(  
    int *__restrict__ col1, int *__restrict__ col2, int *__restrict__ res)  
{  
    int i;  
    ⋮  
}
```

Data-Parallelism Through SIMD Instructions (Core Loop)

Process 16 vector elements per iteration (SIMD + 2 loops unrolled; assumes no overlap of vectors `res[]` and `col1[]/col2[]`):

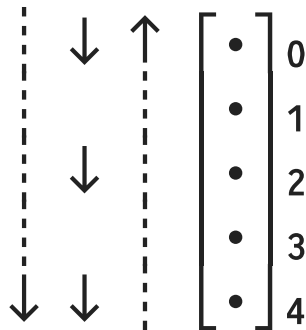
```
PROJECT_sub_int_col_int_col:
  ⋮
  overlap check omitted
  ⋮
  xorl %eax, %eax                4 × 32 bits = 128 bits wide
loop:
  movdqu (%rdi,%rax,4), %xmm0    # %xmm0 ←128 col1[i+0...i+3]
  movdqu 16(%rdi,%rax,4), %xmm1 # %xmm1 ←128 col1[i+4...i+7]
  movdqu (%rsi,%rax,4), %xmm2    # %xmm2 ←128 col2[i+0...i+3]
  psubd %xmm2, %xmm0            # %xmm0 ←128 %xmm0 -128 %xmm2
  movdqu 16(%rsi,%rax,4), %xmm2  # %xmm2 ←128 col2[i+4...i+7]
  psubd %xmm2, %xmm1            # %xmm1 ←128 %xmm1 -128 %xmm2
  movdqu %xmm0, (%rdx,%rax,4)    # res[i+0...i+3] ←128 %xmm0
  movdqu %xmm1, 16(%rdx,%rax,4) # res[i+4...i+7] ←128 %xmm1 .....
  movdqu 32(%rdi,%rax,4), %xmm0 # %xmm0 ←128 col1[i+8 ...i+11]
  movdqu 48(%rdi,%rax,4), %xmm1 # %xmm1 ←128 col1[i+12...i+15]
  movdqu 32(%rsi,%rax,4), %xmm2 # %xmm2 ←128 col2[i+8...i+11]
  psubd %xmm2, %xmm0            # %xmm0 ←128 %xmm0 -128 %xmm2
  movdqu 48(%rsi,%rax,4), %xmm2 # %xmm2 ←128 col2[i+12...i+15]
  psubd %xmm2, %xmm1            # %xmm1 ←128 %xmm1 -128 %xmm2
  movdqu %xmm0, 32(%rdx,%rax,4) # res[i+8 ...i+11] ←128 %xmm0
  movdqu %xmm1, 48(%rdx,%rax,4) # res[i+12...i+15] ←128 %xmm1
  addq $16, %rax                # }
  cmpq $2048, %rax              # } 2048 / 16 = 128 iterations
  jne loop                       # exit if %rax = 2048
  ⋮                               # (non-vectorized code not shown)
```

Loop #n
.....
Loop #n+1

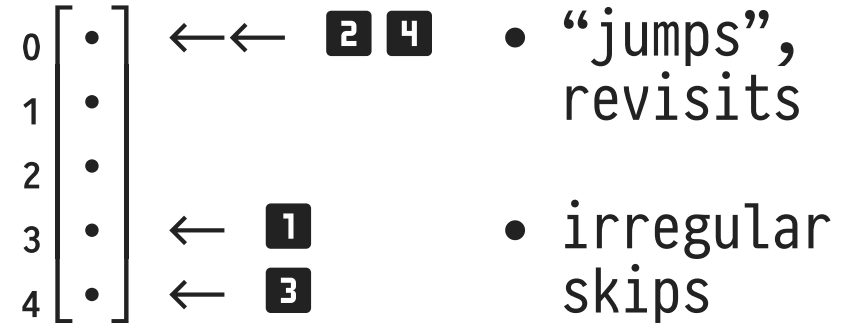
6 : Predictable Memory Access and Prefetching

Predictable Access Patterns

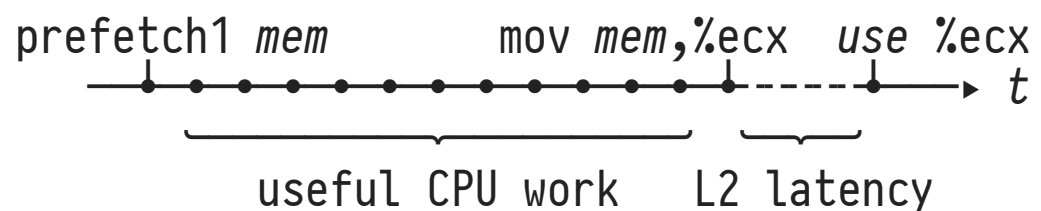
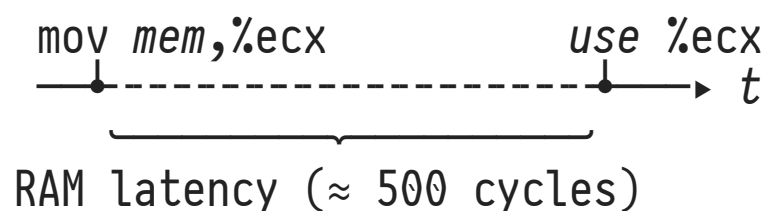
- forward scans (possibly with skips)
- backward scans



Unpredictable Access Patterns



- **Predictable:** CPU automatically issues **asynchronous memory prefetch** operations to preload caches and hide memory latency.
- **Unpredictable:** DBMS code adds explicit **software prefetch⁶ instructions** for memory addresses needed in the future: #028



⁶ No-ops with side effect on the data cache. Example: `prefetcht1` loads into the L2 cache on Intel® Core i7.

7 : Implementing Selection in Tight Loops

The core of the `FILTER` operator is a **conditional statement** that is evaluated for each input vector element (⚠️ simplified again):

```
int FILTER_lt_date_col_date_val(int* res, int* col, int val) 📄 #023
{
    int i, o = 0;

    for (i = 0; i < STANDARD_VECTOR_SIZE; i += 1) {
        if (col[i] < val) { // test filter condition (col < val)
            res[o] = i; // build selection vector
            o += 1;
        }
    }

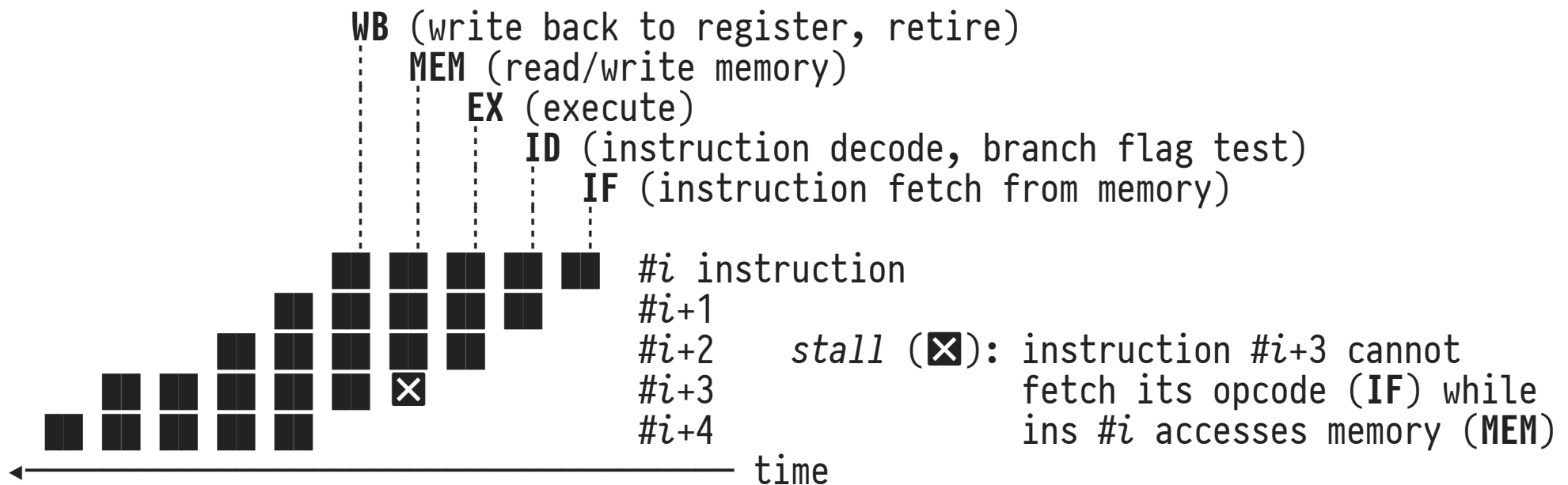
    return o; // output cardinality
}
```

- Conditionals (and loops) lead to **branches** in generated machine code that disrupt the CPU's processing of straight-line instruction sequences. 🗨️

Instruction Pipelining in Modern CPUs

Contemporary CPUs process instructions in a pipeline of execution stages that compete for CPU resources (e.g., the arithmetic logic unit or the memory bus).


- A simple five-stage pipeline:

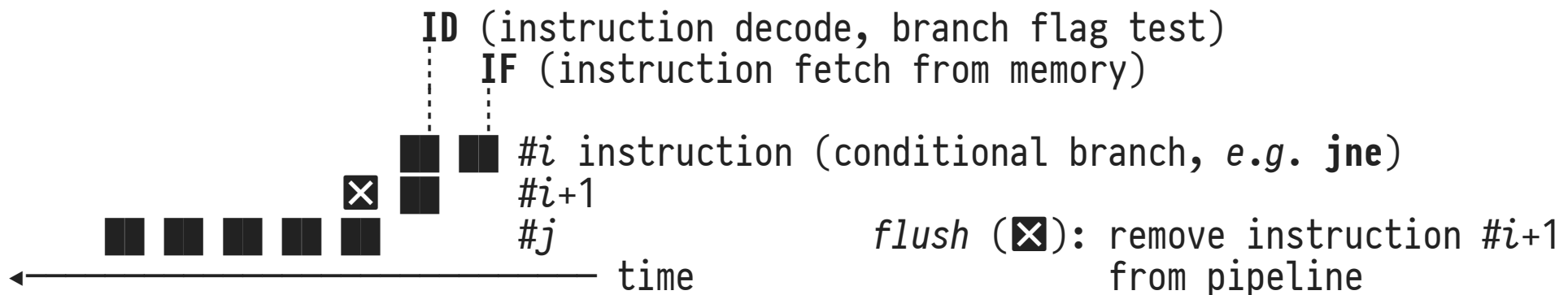


Branch Taken? Yes, Flush Pipeline

Control flow branches (`for`, but particularly `if`) are a challenge for modern pipelining CPUs.

The simple pipeline decides the outcome of branch $\#i$ (at end of **ID**) only *after* instruction $\#i+1$ has already been fetched (**IF**).

- If the branch is taken, **flush** instruction $\#i+1$ from the pipeline , instead fetch instruction $\#j$ at jump target:



Avoiding Branch Mispredictions

A **mispredicted branch** leads to

1. pipeline flushes—effectively a stall **×**—and
2. (possibly) CPU instruction cache misses.

The resulting runtime penalty is significant (≈ 15 cycles).

💡 DBMS code aims to avoid branch mispredictions in tight loops:

- Prefer **branch-less** implementations of operator logic, or at least
- reduce the number of random/hard-to-predict branches.

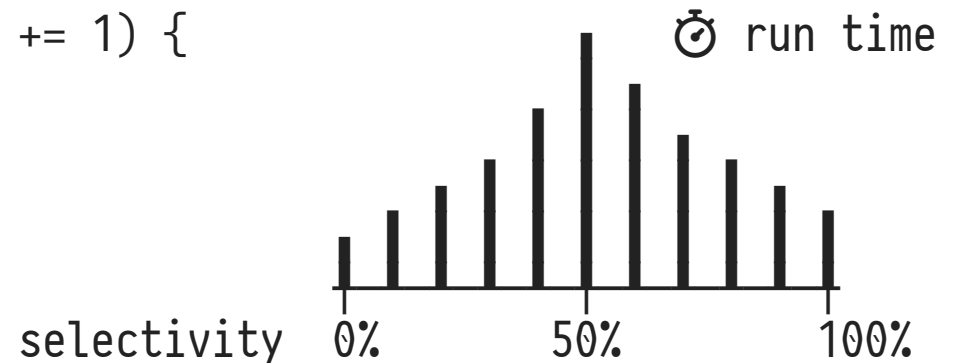
Branch-Less Filtering

Measure wall-clock run time of core **FILTER** loop for different **selectivities** of predicate `col < val`:

```

❶ for (i = 0; i < STANDARD_VECTOR_SIZE; i += 1) {
    if (col[i] < val) {
        res[o] = i;
        o += 1;
    }
}

```

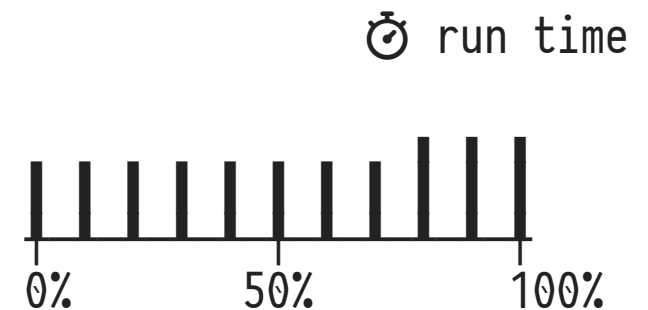


```

❷ for (i = 0; i < STANDARD_VECTOR_SIZE; i += 1) {
    res[o] = i;
    o += (col[i] < v);
}

```

$\equiv 1$ if predicate satisfied, else 0



❷: Only well-predictable loop control flow (**for**) remains. 👍

Mixed-Mode Selection


There is an entire space of possibilities to implement composite predicates (e.g., the conjunction p_1 AND p_2):

- Use branch-less selection via $o += p_1 \ \& \ p_2$ (NB. uses of C's bit-wise *and* operator $\&$).
- Identify the *more selective*⁷ (and thus more predictable) conjunct p_1 , say, then use

```
if (p1) {
    res[o] = i;
    o += (p2);
}
```

 #030

(This approach particularly fits DBMSs that generate code from SQL queries—which DuckDB does *not* do.)

⁷  **This is important.** If p_2 is unpredictable then the mixed-mode selection `if (p2) { ... o += (p1) }` will suffer from branch misprediction.

Design and Implementation of DuckDB Internals

⑧

Query Rewriting and Optimization

April 7, 2026

Torsten Grust
Universität Tübingen, Germany

1 | SQL Needs a Query Optimizer

The **query optimizer** is essential in delivering SQL's promise to function as a *declarative* query language:

- SQL queries describe *what* is to be computed, ...
- ... but do not—in fact: cannot—detail the *how*.

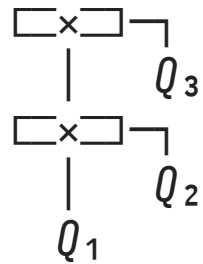
To this end, DuckDB 🐥 implements a SQL processing regime that the DB community has engineered and refined since the late 1970s:

1. Translate the SQL query text into a **canonical plan** 🐼👉 (valid, yet its execution time may well render it unfeasible).
2. Apply equivalence-preserving **plan rewrites** 🐼→⁺🐼, each aiming to reduce execution time and/or resource usage.
3. Map **optimized plan**¹ 🐼👉 into **physical plan**. Execute.

¹ With `PRAGMA explain_output = 'optimized_only'`; DuckDB's `EXPLAIN` displays that final optimized plan. Additionally use `PRAGMA disable_optimizer;` (see below) to inspect the canonical plan.

Canonical Plans

SQL syntax systematically translates into **canonical plans** in which *plan fragments directly correspond with SQL query clauses*:



- A clause **FROM** Q_1, \dots, Q_n on n tables/subqueries translates into a left-deep tree of $n-1$ **CROSS_PRODUCT** (\times) operators.
- The front of this tree exactly mirrors the sequence of the Q_i in **FROM** (no reordering).
- **WHERE** yields **FILTER** operators that sit *above* the $\square \times \square$ tree.
 - **FILTER** processes inputs of potentially enormous cardinality.
- **SELECT** translates into a **PROJECTION** close to the plan root.
 - Unneeded columns are only removed late during execution.

2 : Optimization Passes

DuckDB's query optimizer proceeds in **passes**, each of which consume and produce a valid plan $\text{⋈} \rightarrow \text{⋈}$.

- Each pass is a “specialist” in one particular plan aspect (e.g., Boolean expressions, join ordering, CTE inlining, ...).
- Passes run *once in a pre-determined order*.
 - Predictable optimization effort 👍.
 - Avoid spending dozen of ms to optimize a sub-ms query.
 - Typical pass duration: 10–50 μs , passes overall: 1 ms.
 - Optimization potential may be missed 👎.
 - Passes are *not* iterated (e.g., until fixpoint reached).
 - Pass ordering problem (e.g., statistics lookup may yield `EMPTY_RESULT` too late for propagation into upstream plan).

Controlling the DuckDB Query Optimizer

```
D FROM duckdb_optimizers();
```

name
expression_rewriter
filter_pullup
⋮
join_elimination
window_self_join

```
-- DuckDB v1.5 implements 30+ passes
-- (listed here in arbitrary order)
```

- Disable and re-enable the entire query optimizer:

```
D PRAGMA disable_optimizer; -- ⚠ execute canonical plans
D PRAGMA enable_optimizer;
```


- Disable and re-enable all selected optimization passes:

```
D SET disabled_optimizers = 'filter_pullup, join_order';
D SET disabled_optimizers = '';
```

Order of Optimization Passes in DuckDB (version 1.5)

#	Pass	Specializes in...
1	expression_rewriter	simplifying/evaluating scalar expressions
2	cte_inlining	choosing between inlining or materializing CTEs
3	sum_rewriter	rewriting aggregates of the form sum (<i>e</i> + <i>const</i>)
4+5	filter_pullup+pushdown	moving FILTER operators downstream/upstream
6	cte_filter_pusher	moving FILTERs into materialized CTEs
7	regex_range	turning reg.exp. matches into range predicates
8	in_clause	turning IN (...) into regular filter or join
9	delimiter	removing redundant DELIM_JOINS/DELIM_GETs
11	empty_result_pullup	simplifying plans in presence of empty results
12	window_self_join	turning window computations into self-joins
13	join_order	turning × into joins, decide join order
14	join_elimination	detecting and removing redundant (semi-)joins
16	unused_columns	identifying unused columns to make scans narrower
17	duplicate_groups	removing redundant grouping criteria
18	common_subexpressions	evaluating common scalar subexpressions only once
19	column_lifetime	introducing PROJECTIONS to remove unused columns
20	build_side_probe_side	placing join inputs on lhs (build)/rhs (probe)
21	common_subplan	converting common subplans into materialized CTEs
22	limit_pushdown	moving LIMIT below PROJECTIONS (limit early)
23	row_group_pruner	identifying row groups that need not be scanned
25	top_n	merging ORDER_BY and LIMIT into TOP_N
26	late_materialization	narrowing inputs, joining in columns late in plans
27	statistics_propagation	using column statistics to simplify plans
30	reorder_filter	evaluating cheap clauses first in con/disjunctions
31	join_filter_pushdown	filtering probe side based on build side values

💡 Some **scalar expressions** e (e.g., in **SELECT/WHERE** clauses) can be **statically evaluated**—maybe partially only—*without* base data access.

- Savings are minimal for any individual evaluation of e , but e may be evaluated millions of times during plan execution.
- Simplification may reveal opportunity for further plan rewrites.
- Implemented in DuckDB:
 - **Arithmetic/Boolean simplification:**  #033
 - $const - const$, $e + 0$, $e // 1$, $e * NULL$.
 - $const \wedge const$, $e \vee false$, $NOT (e_1 < e_2)$, distributivity.
 - **CASE simplification:** constant guards, branch pruning.
 - **LIKE/regular expression match simplification:**
 - Rewrite **LIKE** into range or **suffix/contains** predicates.
 - Trade **regexp_matches** for the computationally lighter **LIKE**.

4 | Reordering Clauses in AND/OR Predicates

reorder_filter

DuckDB evaluates conjunctions² e_1 AND e_2 AND ... AND e_n left to right using Boolean shortcut.

💡 Reorder clauses e_i from cheap to costly.

- Heuristic expression cost (read $<$ as “is cheaper to evaluate”):
 - Ops: $+$ $<$ $*$ $<$ $/$ $<$ `round` $<<$ `~~`. Types: `int` $<$ `double` $<$ `text`.
 - $cost(e_1 \otimes e_2) = cost(e_1) + cost(e_2) + cost(\otimes)$.

$$cost(e :: \tau) = \begin{cases} cost(e) & \text{if } e \text{ has type } \tau \\ cost(e) + 200 & \text{if } e :: \text{text or if } \tau = \text{text} \\ cost(e) + 5 & \text{otherwise} \end{cases}$$

- ⚠ Do not reorder if any e_i can possibly fail (e.g., cast to `int`): moving e_i rightwards may hide an observable side effect.

² This also applies to disjunctions (OR).

5 : Using Statistics to Simplify Plans : statistics_propagation

DuckDB maintains basic **statistics for base table columns**.

💡 Propagate statistics through plan operators to derive min/max value ranges and constraints *without* access to the base data.

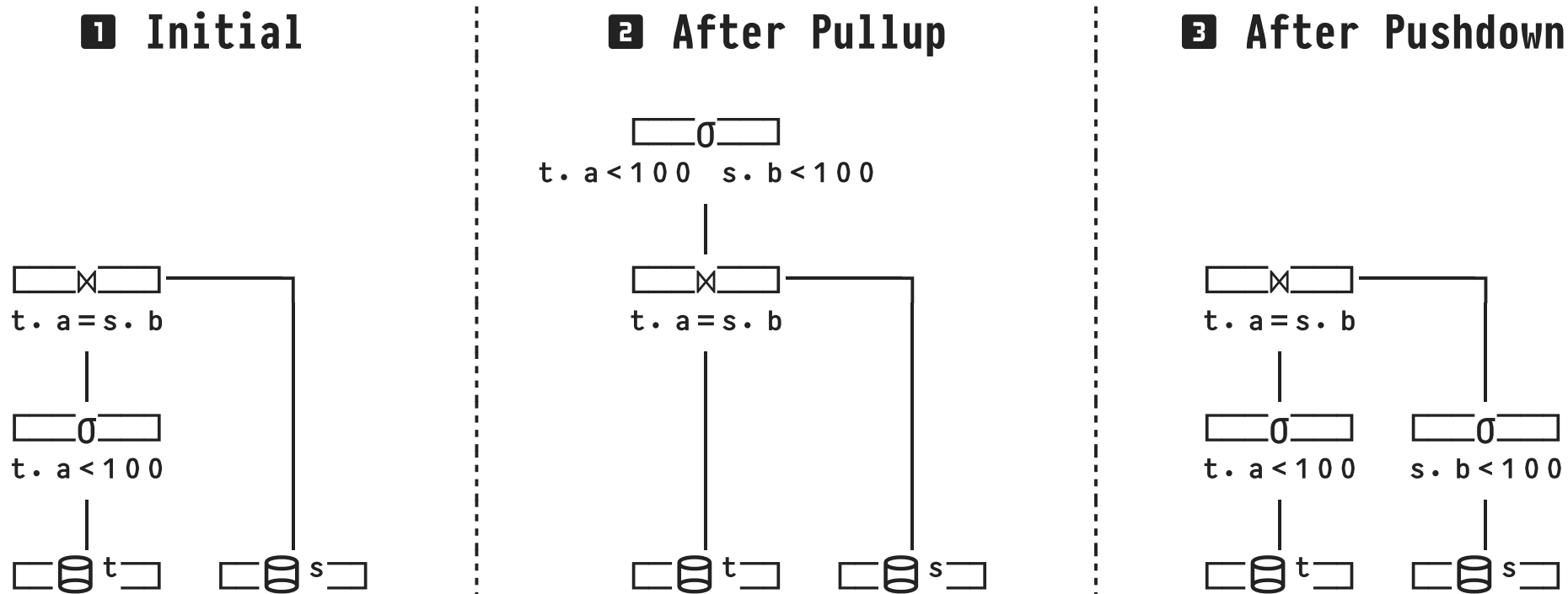
	Propagated statistics (excerpt)
intermediate tables	<ul style="list-style-type: none"> • estimated cardinality • maximum cardinality (helps allocation)
numeric columns	<ul style="list-style-type: none"> • min/max value (may be unknown or uninteresting if entire domain of column type is covered) • is constant? • contains NULL?
text columns	<ul style="list-style-type: none"> • min/max value (first 8 characters) • maximum string length • contains Unicode (non-ASCII) characters?

- Propagated statistics help to identify **FILTER** predicates that will be *always true* (remove) or *always false* (**EMPTY_RESULT**).

6 : Moving Filters Downstream + Upstream : filter_pullup+pushdown

💡 Move **FILTERs** upstream (“**filter pushdown**” towards/into base table scans) to reduce intermediate table cardinalities early.

- DuckDB additionally implements **filter pullup** (propagate restrictions through equality predicates):³



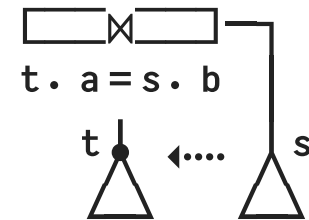
³ Here and in the sequel let us adopt common operator abbreviations: $\square \text{⊗} \square$: SCAN, $\square \sigma \square$: FILTER, $\square \times \square$: JOIN.

Dynamic Join Filter Pushdown

! join_filter_pushdown

Join execution proceeds in two (alternating) phases:

1. Read rows from build/outer input s .
 ⚡ Record min/max values m^b/M^b in column $s.b$.
2. While reading rows from the probe/inner input t ,
 ⚡ use m^b/M^b on column $t.a$ to pre-filter rows.

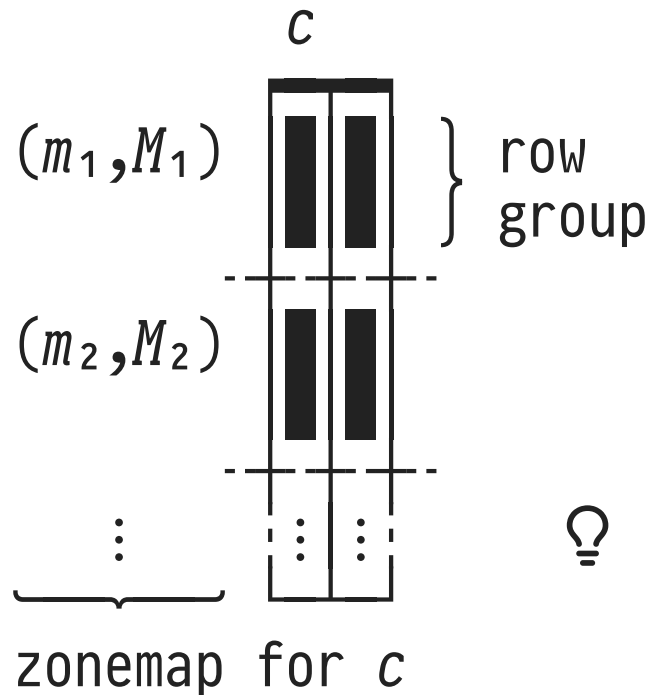


- This **join filter pushdown** (“*sideways information passing*” ←…) is a *dynamic* optimization performed *during* plan execution: only then are m^b/M^b known.⁴
 - During query optimization, pass `join_filter_pushdown` only identifies the scan • which will apply the dynamic predicate.
 - ⚠ Only **EXPLAIN ANALYZE** will reveal this optimization.

⁴ If the distinct values in $s.b$ are $\{x_1, \dots, x_n\}$ (with $n \leq \text{dynamic_or_filter_threshold}$, default 50), the dynamic predicate will be `IN (x1, ..., xn)` to filter t precisely.

7 : Scanning Row Groups in Order, Exit Early

: row_group_pruner



- DuckDB can **scan row groups in an order** defined by the zonemap entries (m_i, M_i) :
 - Scan row groups by ascending m_i (min).
 - Scan row groups by descending M_i (max).

💡 Use these scan orders to limit the row groups read to process TOP_N (**ORDER BY** c + **LIMIT** N).

FROM t **ORDER BY** c **DESC** **LIMIT** N -- rows with N largest c



📄 #037

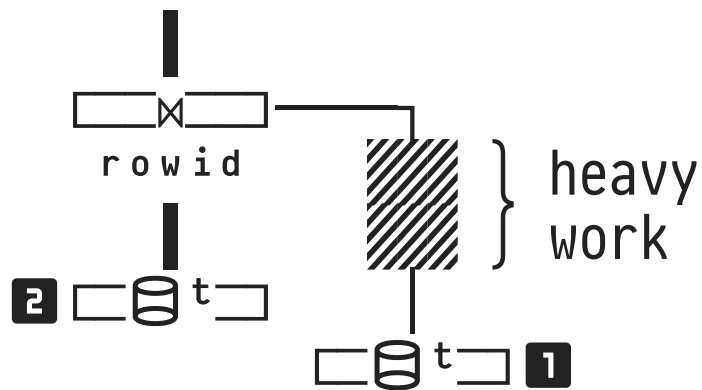
💡 Scan row groups in descending M_i (max) order. **Stop criterion:**



⚠️ $(m_i, M_i) \cap (m_j, M_j) = \emptyset?$ (no overlap of c values between row groups?) $\left\{ \begin{array}{l} \text{yes, stop scanning row groups once} \\ \quad R \geq N \text{ rows have been read} \\ \text{no, stop after at most } N \text{ row groups} \end{array} \right.$

8 : Attaching Columns Late During Execution : late_materialization

Execution benefits if rows are kept *narrow* as long as possible:

1. **Push down projections** to scan relevant columns (say, for filtering, grouping, sorting, see  below) only.
2.  *Late in the plan* use **row IDs to join in payload columns** to build (a potentially wide) query result.



- TABLE_SCAN **1** performs projections, leaving only columns relevant in .
- TABLE_SCAN **2** receives the list of relevant row IDs (“rowid pushdown”) to access only few yet wide rows.
-  is a left semi-join on row IDs.

rows: | narrow, | wide

9 | Rewriting `sum()` Into `sum() + count()`

| `sum_rewriter`

💡 Use associativity and commutativity of `+` to save `sum` aggregate evaluation from the repeated addition of constant values:

$$\text{sum}(e + \text{const}) \rightarrow \text{sum}(e) + \text{const} \times \text{count}(e)$$

• Notes:

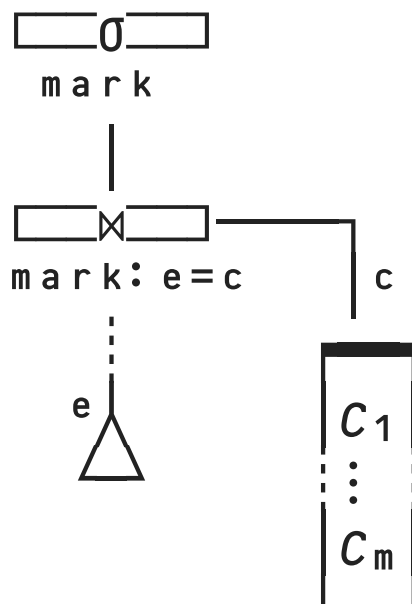
- In IEEE754 floating point arithmetic, `+` is not associative. DuckDB applies pass `sum_rewriter` only if `e` and `const` are integer-typed.
- Both sides are equivalent also if `e` may be `NULL`. (In SQL, aggregates `sum` and `count` ignore `NULL` values.)
- Expression `e` is only evaluated once (by a dedicated `PROJECTION` operator) and then re-used in both aggregates.

10 : Rewriting IN (...) Clauses

: in_clause

Naively evaluating predicates like $e \text{ IN } (c_1, \dots, c_m)$ may be expensive: $O(n \times m)$ if the predicate is tested for n rows.

- ! If m is large, rewrite predicate into an $O(n + m)$ **HASH_JOIN**:
 - Build side: single-column table of m rows carrying the c_i .
 - Probe side: provides n bindings for e .



COLUMN_DATA_SCAN

#040

- Optimizer creates a COLUMN_DATA_SCAN pipeline source that delivers the c_i .
- Since e and the c_i are irrelevant once the predicate is evaluated, $\square \times \square$ is a **MARK** join:
 - Generate a Boolean *mark* that holds the result of the equality comparison.
 - Process *mark* as needed (here: $\square \sigma \square$).

11 : Turning Window Computations Into Self-Joins : window_self_join

SQL **windows** establish frames of related rows⁵ in a table *t* that a function *f* can then process jointly. ? For simple frame specs and *f*, the same operation can be expressed by a self-join of *t*:

<pre> 1 SELECT agg(e) OVER (PARTITION BY c) FROM t </pre>	<pre> 2 WITH frame(part,agg) AS (SELECT c AS part, agg(e) FROM t GROUP BY c) SELECT agg FROM t SEMI JOIN frame ON (c IS NOT DISTINCT FROM part) </pre>
---	--

- Pass `window_self_join` rewrites **1** into **2** on the plan level for
 1. window functions *f* that are SQL aggregates and
 2. frames that are exclusively specified by `PARTITION BY`.

⁵ Rows relate based on position in an ordering (as in `ORDER BY ... ROWS BETWEEN ... AND ...`) and/or inside partitions (`PARTITION BY ...`).

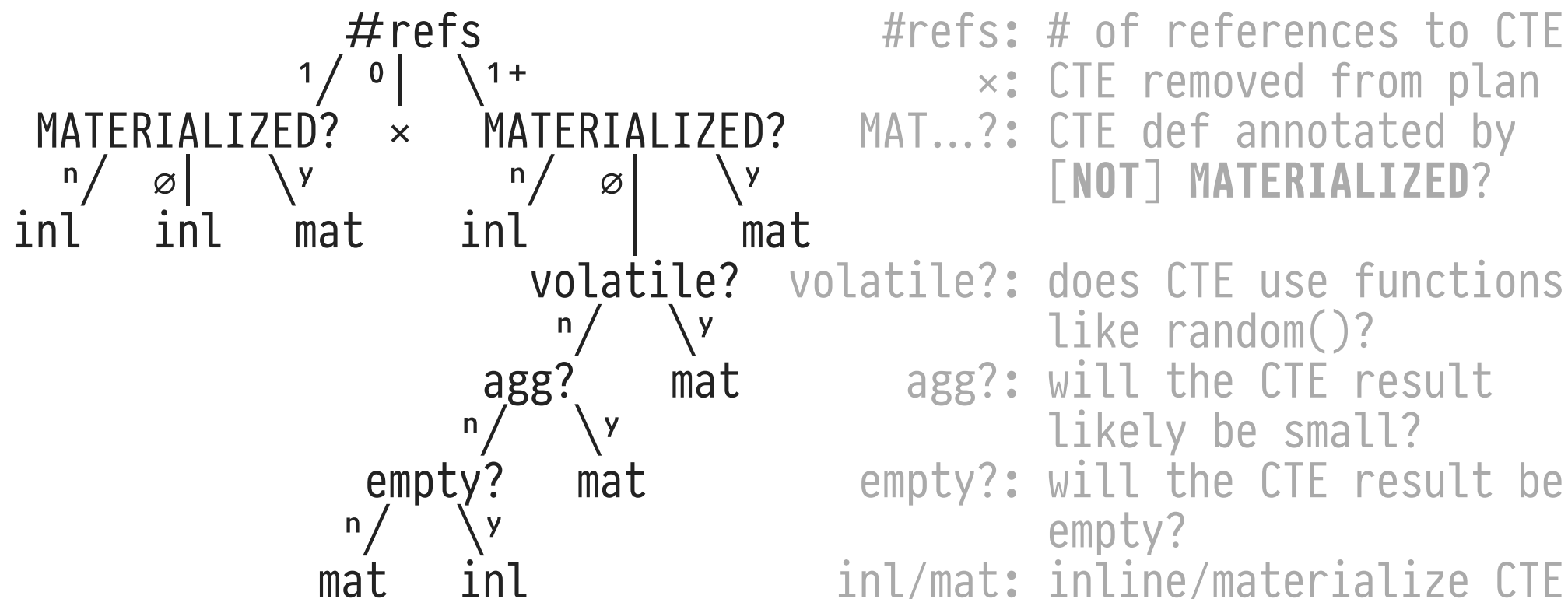
12 : Inlining or Materializing CTEs?

: cte_inlining


Common table expressions (CTEs) may be referred multiple times. Shall the CTE be **materialized once or inlined** at every reference?

💡 DuckDB uses a decision tree to answer that question:

 #042



SQL queries may exhibit query fragments that are syntactically similar (albeit not identical). Can evaluation be shared?

1.  Identify **identical plan—not query—fragments** (*subplans*),
 2. evaluate subplan once and materialize its result,
 3. replace each fragment by a scan of the materialized result.
- Pass `common_subplan` builds on DuckDB's CTE infrastructure:
 - Use plan operator `CTE` to evaluate and materialize the subplan (under fabricated CTE name `__common_subplan_<n>`).
 - Use plan operator `CTE_SCAN` (multiple times) to read the materialized result.
 - Do not share subplans that call on volatile functions.

14 : Introducing and Reordering Joins

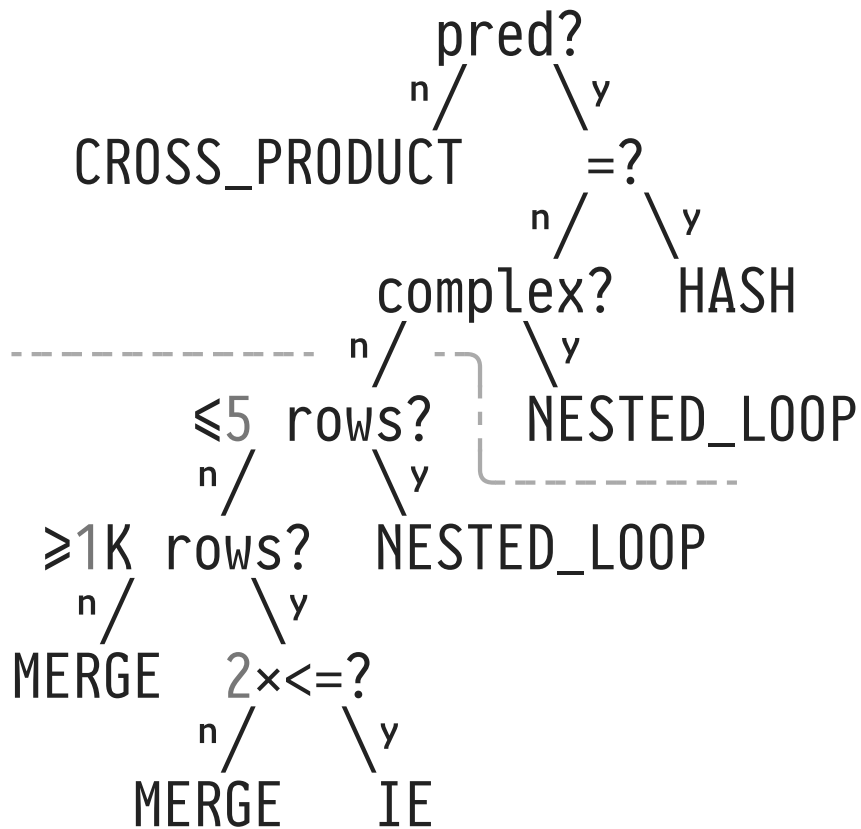
Join (operator \bowtie , algebraic: $S \bowtie_p T$) is fundamental in SQL and query optimization. Remains a hot topic in database research.

DuckDB's join optimizer considers a whole range of aspects:

1. Join predicates p come in many forms. When building the physical plan, **pick an implementation** for \bowtie that fits the given predicate(s).
2. Join can reduce or potentially *multiply* the **cardinalities** of its input tables: $0 \leq |S \bowtie T| \leq |S| \times |T|$.
3. Join operates on two tables. In an n -table query, it is essential to select an efficient **order of the $n-1$ joins**.
4. $S \bowtie T$ is commutative. In most join *implementations*, however, S and T assume **assymmetric** roles. Choose $S \bowtie T$ or $T \bowtie S$ wisely.

Join Predicates Determine Implementations of \bowtie in DuckDB

Inspect join predicate p to select a fitting join implementation:



```

pred?: any predicate at all?
  =?: any equality? (equi-join)
complex?: complex predicate? (~~)
≤5 rows?: nested_loop_join_threshold
≥1K rows?: merge_join_threshold
2x=? : two or more inequalities
  
```

- Subtree below ---- exclusively concerned with inequalities.
- Equi-joins are predominant in typical SQL workloads \Rightarrow HASH_JOIN is DuckDB's workhorse.

Join Implementations for $S \bowtie_p T$ in DuckDB

NESTED_LOOP_JOIN

```

 $T_m \leftarrow \text{materialize}(T)$ 
for  $s \in S$             $S$ : outer
  for  $t \in T_m$         $T$ : inner
    if  $p(s,t)$ 
      emit  $s\#t$  #: output col.s
  
```

- **BLOCKWISE_NL_JOIN**: process chunks of S/T at a time.

PIECEWISE_MERGE_JOIN

```

 $T_s \leftarrow \text{sort}\langle p \rangle(T)$ 
for chunk  $C \in S$ 
   $C_s \leftarrow \text{sort}\langle p \rangle(C)$ 
  for  $(s,t) \in \text{merge}(T_s, C_s)$ 
    if  $p(s,t)$ 
      emit  $s\#t$  scan  $T_s, C_s$  once top-down
  
```

HASH_JOIN

 #044

```

 $H \leftarrow \text{table}(\text{hash}\langle p \rangle, T)$     $T$ : build
for  $s \in S$             $S$ : probe
   $h \leftarrow \text{hash}\langle p \rangle(s)$ 
  for  $t \in H(h)$        scan bucket in  $H$ 
    if  $p(s,t)$ 
      emit  $s\#t$ 
  
```

IE_JOIN ($p \equiv S_1.aS_2.a \wedge S_1.b > S_2.b$)

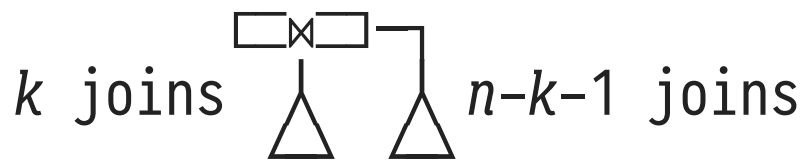
```

 $L_1 \leftarrow \text{sort}\langle a \rangle(S)$    <: ascending
 $L_2 \leftarrow \text{sort}\langle b \rangle(S)$    >: ascending
for  $j \in 1 \dots n$ 
   $s \leftarrow L_2[j]$         $i/j$ : pos of row  $s$ 
   $i \leftarrow s @ L_1$        in  $L_1/L_2$ 
  for  $t \in L_1[1 \dots i-1] \cap L_2[j+1 \dots n]$ 
    emit  $s\#t$ 
  
```

Join Tree Shapes

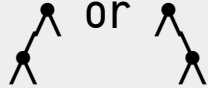
⋮ join_order

A SQL join query with $n+1$ tables requires n binary joins. $\square \bowtie \square$ is *associative*: the number of possible **join trees** grows quickly.



of join tree shapes:

$$C_n = \sum_{k=0}^{n-1} C_k \times C_{n-k-1}$$

n (# of $\square \bowtie \square$)	C_n (# of join tree shapes) ⁶	
0	1	single-table query
1	1	two tables, single join \wedge
2	$C_0 \times C_1 + C_1 \times C_0 = 1 + 1 = 2$	 or \wedge
3	5	
4	14	
5	42	
10	16,796	

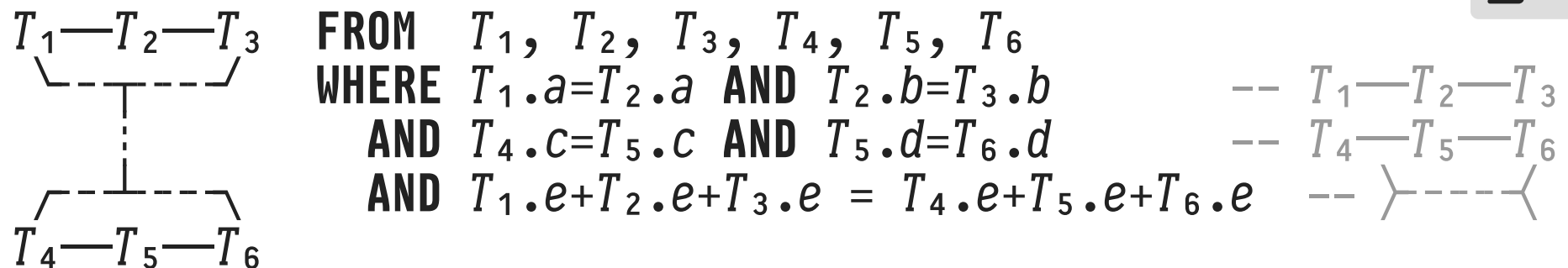
It is not reasonable to fully explore this join tree search space.

⁶ C_n is the n th Catalan number. Closed form: $C_n = (2n)! / ((n+1)! \times n!)$.

Enumerating Join Trees

: join_order

- From the SQL join query, build a **join hypergraph** in which n tables are connected by join predicates:

 #046


- Use dynamic programming⁷ to iteratively build join trees of minimal $cost()$ with $2, 3, \dots, n$ tables. Build join trees of size k from memoized trees of size i and $k-i$.
 - Disregard non-connected tables to avoid cross products.
 - $cost(S \bowtie T) \stackrel{\text{def}}{=} |S| + |T| + |S \bowtie T|$. (all estimated)

⁷ DuckDB's [join_order](#) optimizer is modelled after the hypergraph-aware dynamic programming algorithm *DPhyp* described in G. Moerkotte, T. Neumann, “*Dynamic Programming Strikes Back*”, SIGMOD, Vancouver, Canada, 2008.

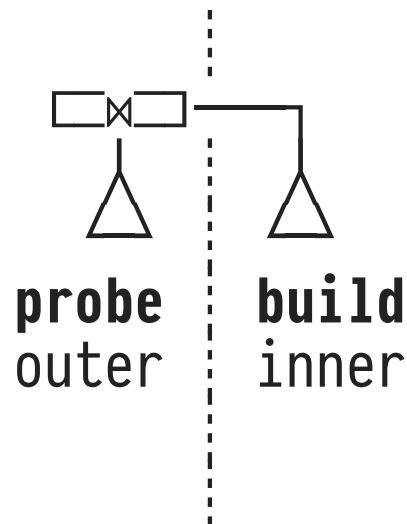
Choose $S \bowtie T$ vs $T \bowtie S$

build_side_probe_side

$S \bowtie T$ is commutative but the chosen operators $\square \bowtie \square$ do *not* treat the inputs symmetrically.

lhs

- Larger inputs
- Simple subplans (scans)
- Inputs that pass many columns/row IDs downstream



rhs

- Smaller inputs **s**
- Complex subplans (other joins)

- Heuristic **s** (smaller inputs on the *rhs* of $\square \bowtie \square$):
 - **HASH_JOIN**: Build and materialize small(er) hash table H .
 - **NESTED_LOOP_JOIN**: Use less space to materialize inner T_m .
 - **PIECEWISE_MERGE_JOIN**: Use less space to hold sorted T_s .

15 : Nested Subqueries and Correlation

: (*mandatory*)

Wherever a SQL query expects a scalar or tabular value, a **nested subquery** (enclosed in (...)) may be used to compute that value:

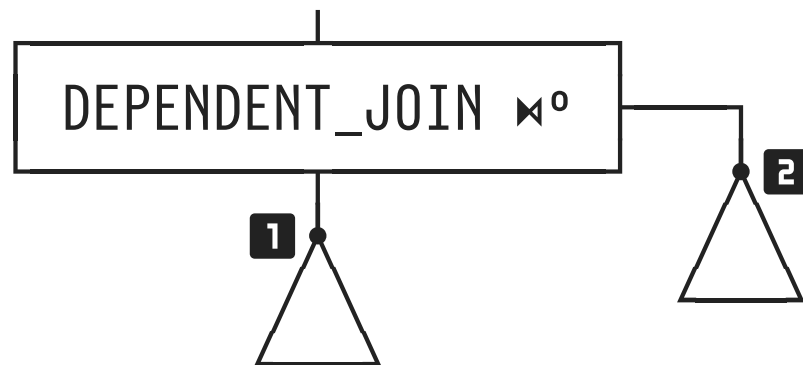
```
SELECT o.o_orderkey, o.o_orderdate, l.l_linenumbr
FROM   orders AS o, lineitem AS l
WHERE  o.o_orderkey = l.l_orderkey
AND    l.l_extendedprice = (
    SELECT min(i.l_extendedprice)
    FROM   lineitem AS i
    WHERE  i.l_orderkey = o.o_orderkey);
```

Q₁: All TPC-H orders along with (the linenumbr of) their cheapest lineitem

- Subquery [...] is **correlated** (row variable `o` occurs free at \blacktriangledown).
- SQL semantics: Re-evaluate query in [...] for each `o,l` pair.
- For sizeable input tables, a naive nested-loop evaluation strategy will lead to unacceptable query run times. 🗨️

Correlation Leads to Dependent Joins

In initial/non-optimized query plans, correlation is represented by the `DEPENDENT_JOIN` operator (symbol: \bowtie^o):




01 #048 1

plan for nested query \square
(contains references to o \blacktriangleright)

plan for outer query
(provides bindings for o)

- **Nested loop semantics:** Evaluate **2** for each row produced by **1**.
- **Query decorrelation:** Remove *all* occurrences of `DEPENDENT_JOIN`,⁸ rewrite plans **1** and **2** such that both are evaluated only once.

⁸ In fact,  does not implement `DEPENDENT_JOIN` as a *physical* plan operator. Query decorrelation thus is an absolute must in DuckDB.

16 : Optimization: Query Decorrelation

A correlated subquery acts like a (table-valued) **function** whose parameters are the free row variable references—in Q_1 : $[\](o)$. Do not compute that function for the same parameters more than once.

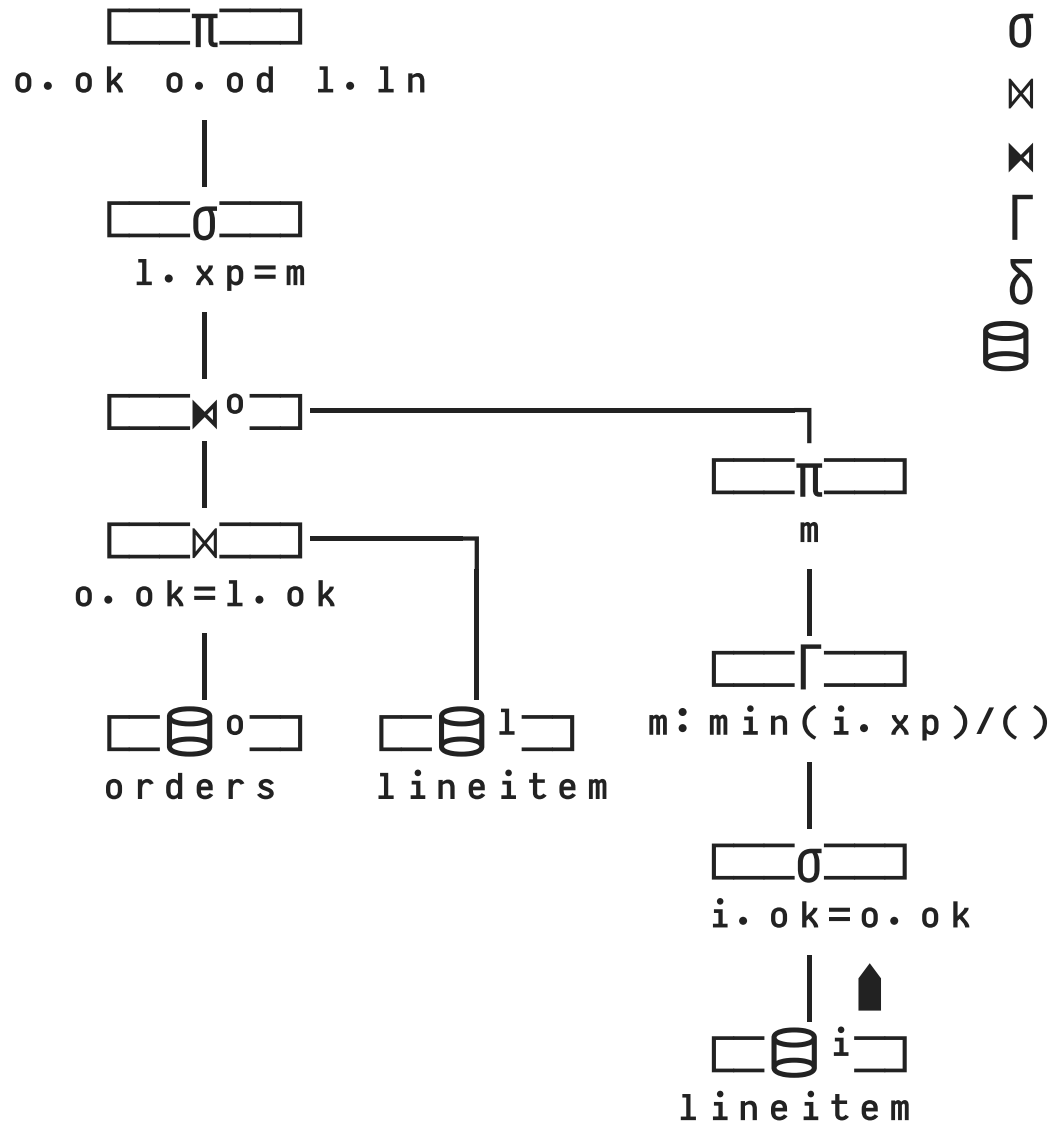
💡 Query decorrelation:


1. Compute the duplicate-free set D of distinct parameters.
2. Evaluate the subquery $[\]$ for *all* parameters $d \in D$ in one go, yielding a lookup table T of $(d, [\](d))$ pairs.
3. Evaluate the outer query by joining with lookup table T .

Thomas Neumann and Alfons Kemper of TU Munich⁹ implement this idea via systematic rewrites that push **DEPENDENT_JOIN** \bowtie down into a query plan until it can be trivially replaced by regular **JOIN** \bowtie .

⁹ T. Neumann, A. Kemper, “Unnesting Arbitrary Queries”, Proc. of the BTW Conference, Hamburg, Germany, 2015.

Query Decorrelation (Initial Plan)



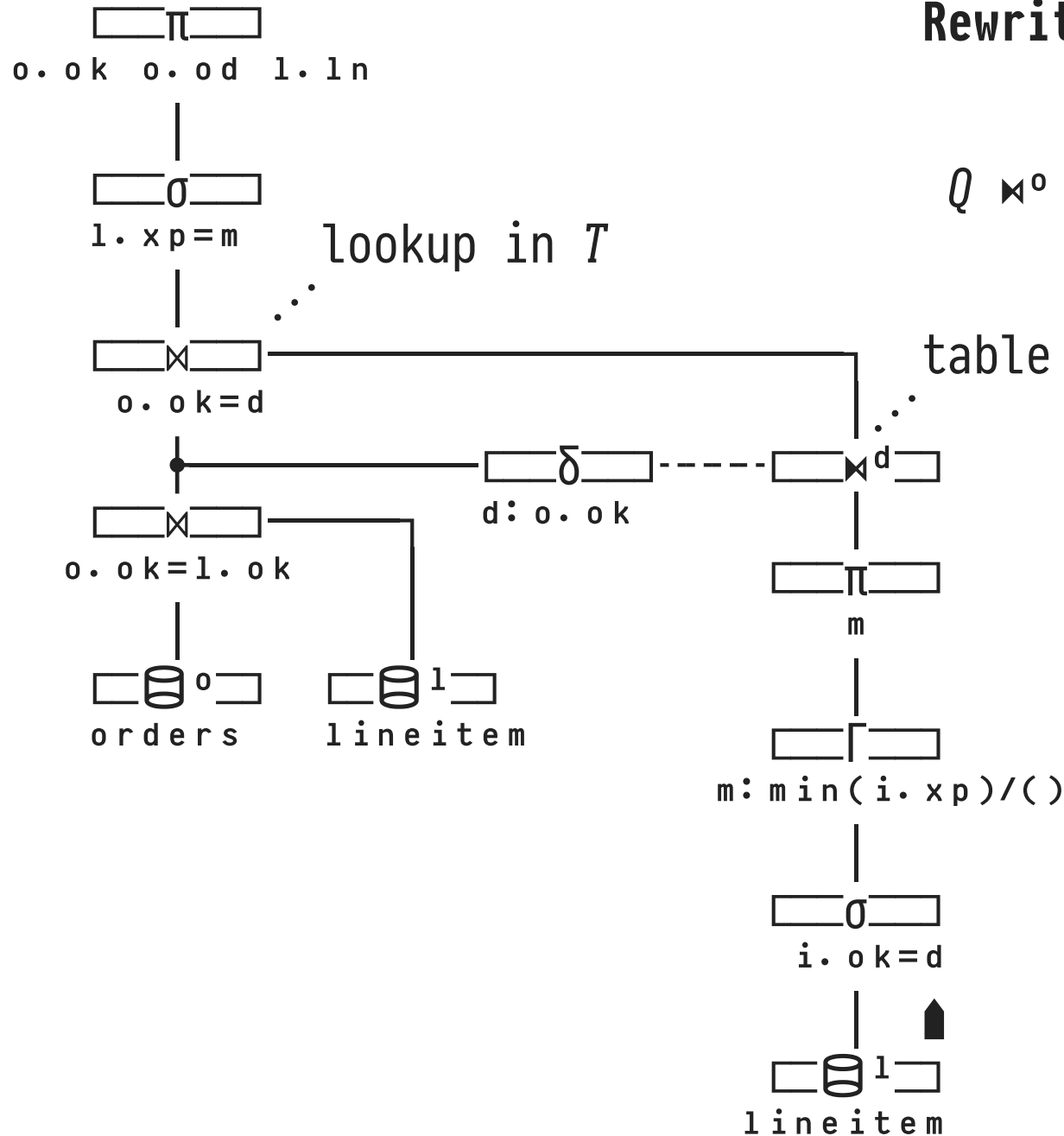
π PROJECTION
 σ FILTER
 \bowtie JOIN
 \bowtie^o DEPENDENT_JOIN
 Γ AGGREGATE/(grouping)
 δ PROJECTION (dup. elimination)
 SCAN

 orders/lineitem columns:

`ok` `l/o_orderkey`
`od` `o_orderdate`
`ln` `l_linenum`
`xp` `l_extendedprice`

 dependency (free row var)

Query Decorrelation 1 (Compute Parameter Set D)



Rewrite rule:

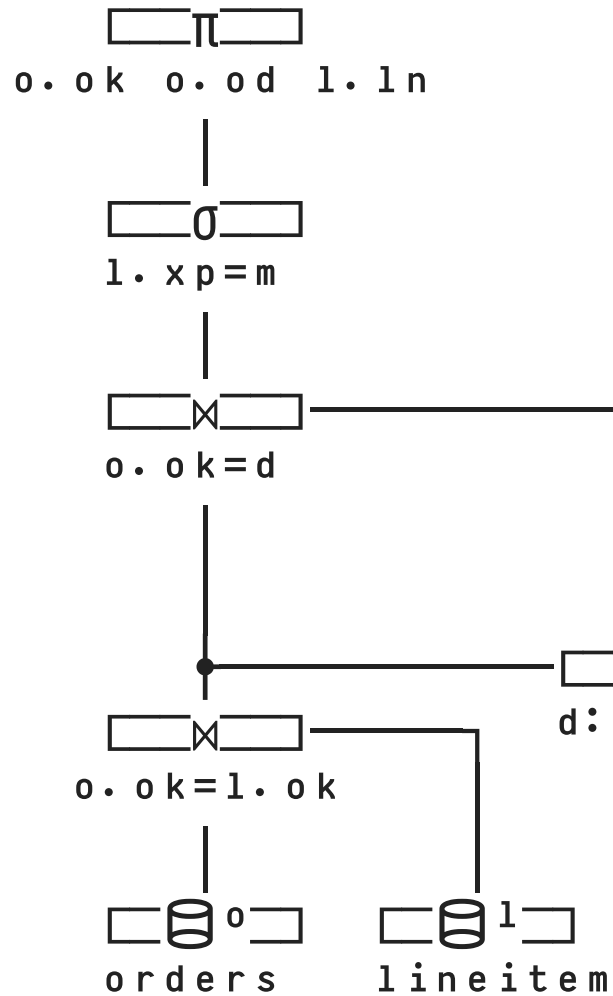
$$D := \delta^{d: o}(Q)$$

$$Q \bowtie^o [\] \equiv Q \bowtie^{o=d} (D \bowtie^d [\])$$

table $T(d,m)$

----- no duplicates

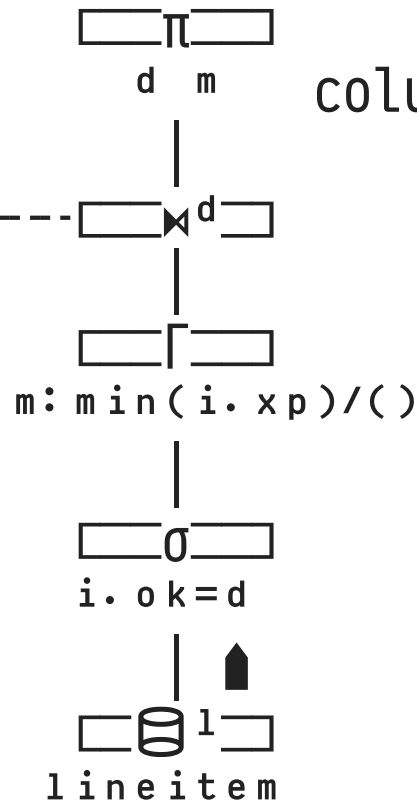
Query Decorrelation (Push \bowtie Down Through π)



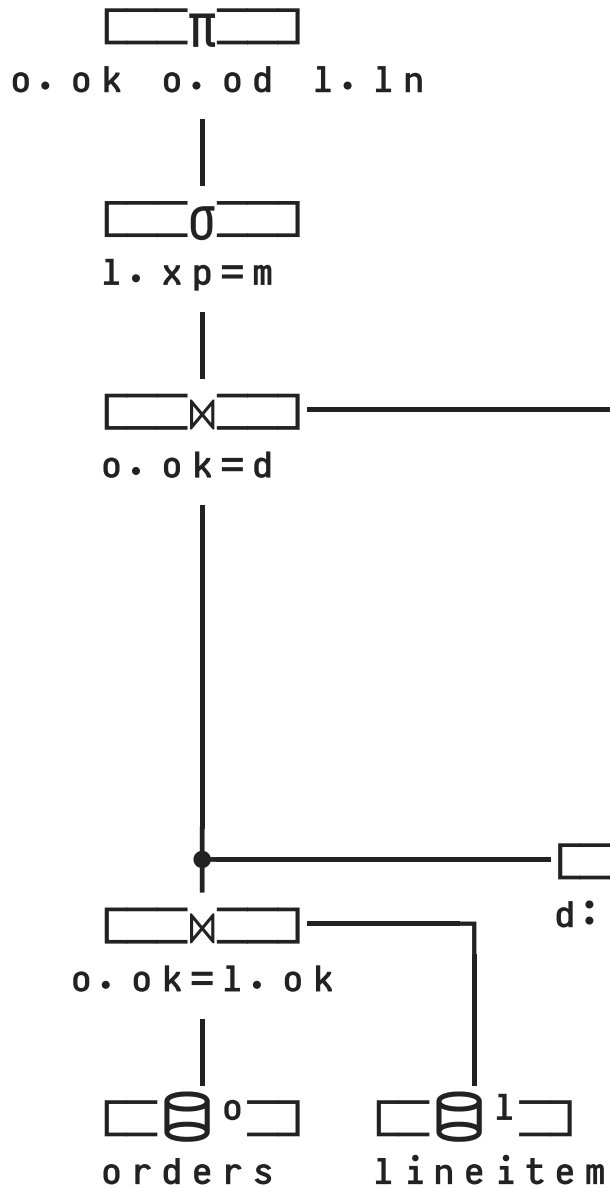
Rewrite rule:

$$D \bowtie^d \pi^m(Q) \equiv \pi^{d,m}(D \bowtie^d Q)$$

column list extended by d



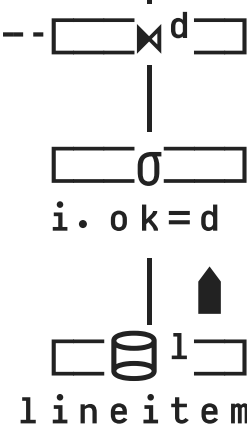
Query Decorrelation \boxtimes (Push \bowtie Down Through Γ)



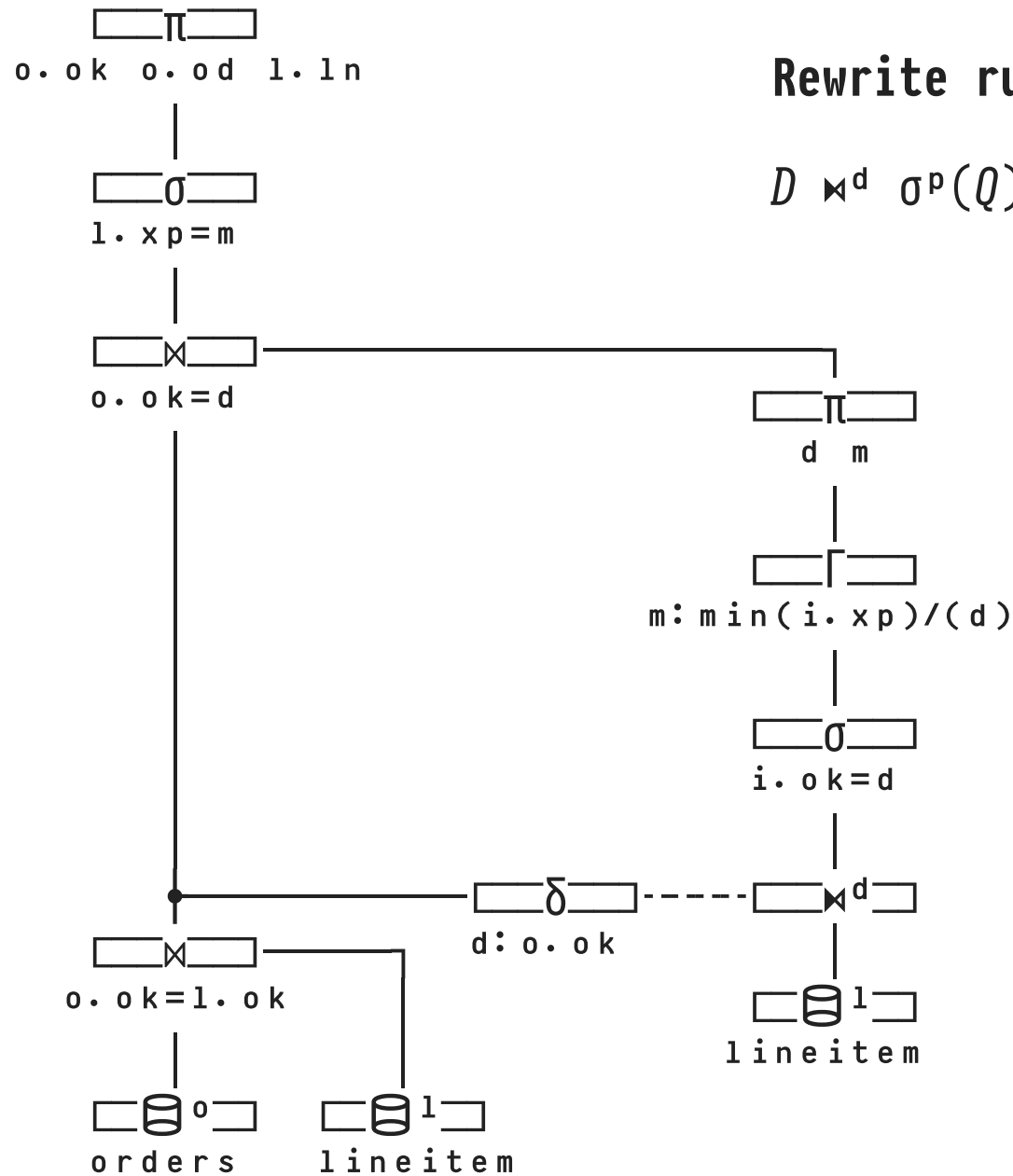
Rewrite rule:

$$D \bowtie^d \Gamma^{m: a/(g)}(Q) \equiv \Gamma^{m: a/(dg)}(D \bowtie^d Q)$$

aggregation now
grouped by d



Query Decorrelation (Push \bowtie Down Through σ)





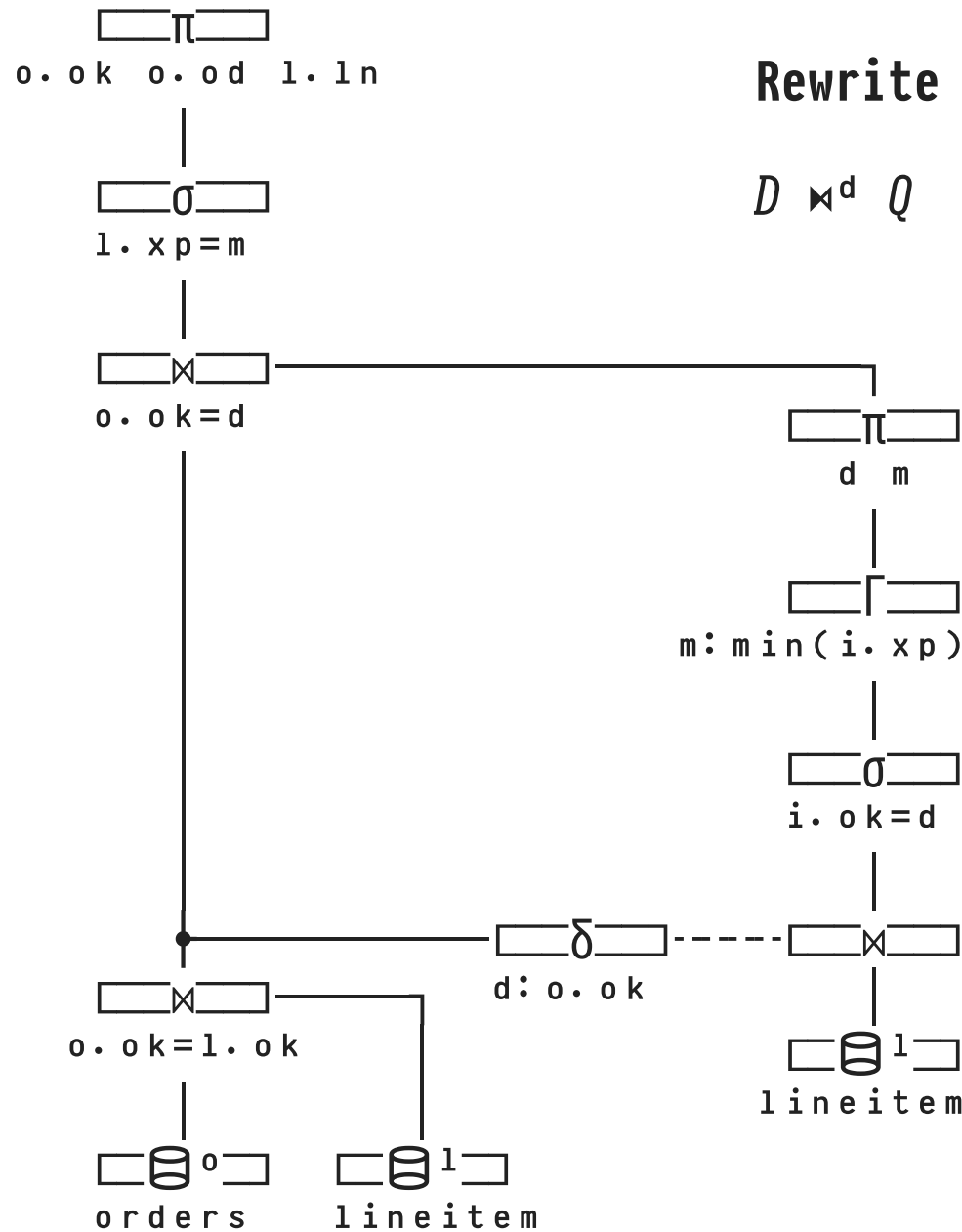
Rewrite rule:

$$D \bowtie^d \sigma^p(Q) \equiv \sigma^p(D \bowtie^d Q)$$

 no dependency on d
below \bowtie^d

Query Decorrelation (Eliminate)

 #048 

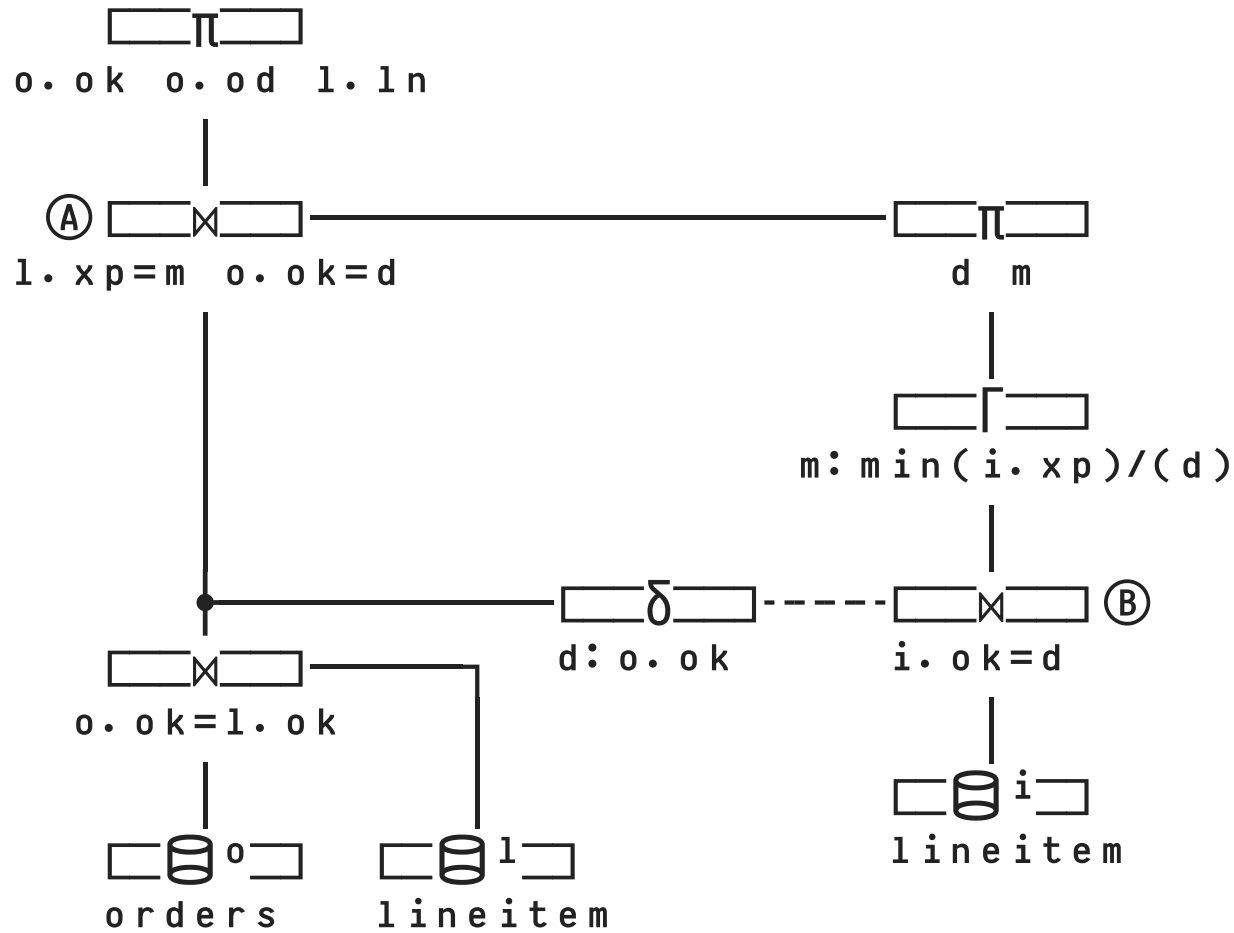


Rewrite rule:

$$D \bowtie^d Q \equiv D \bowtie Q \text{ if } d \text{ not free in } Q$$

\bowtie replaced by a regular \bowtie

17 : Post-Optimizations 6 (Push σ Down Into \bowtie)

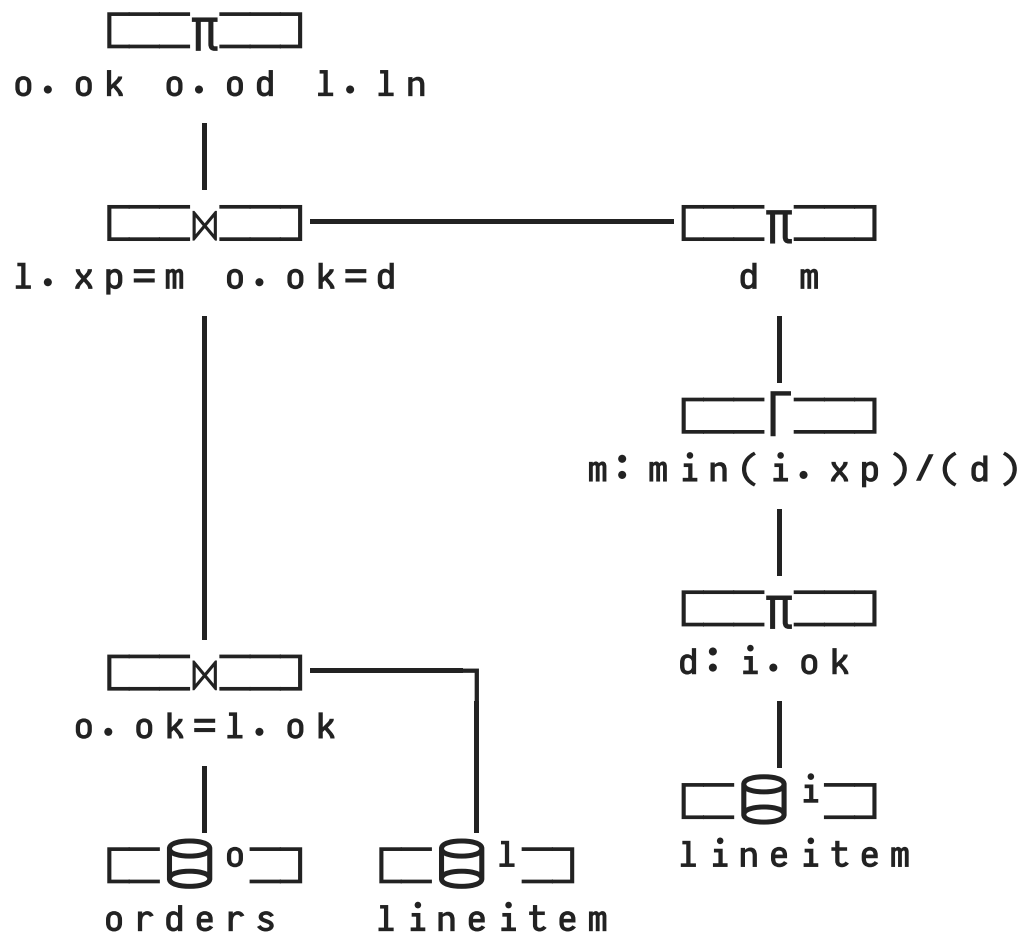


A Pulled $\sigma^{l.xp=m}$ down into \bowtie .

B Pulled $\sigma^{i.ok=d}$ down into \bowtie .

18 : Post-optimization: decoupling

01 #048 3



Optimization (performed by 🐧):

Save $\square\delta\square$ --- $\square\bowtie\square$. Instead, derive d from $i.ok$ via $\square\pi\square$ ($i.ok=d$ holds after $\square\bowtie\square$).

👍 Decouples subquery entirely from the rest of the plan.

⚠️ Need not-NULL-check ($\square\sigma\square$) should $i.ok$ be nullable.

The End.

Since you've got this far... A DuckDB-based companion course on the fundamentals of the relational data model and SQL is available at

<https://github.com/DBatUTuebingen/TaDa>.¹⁰

¹⁰ TaDa: Tabular Database Systems  