

Tabular Database Systems

①

Tabular Data and Database Systems

March 17, 2026
















**Torsten Grust
Universität Tübingen, Germany**

1 | Welcome!

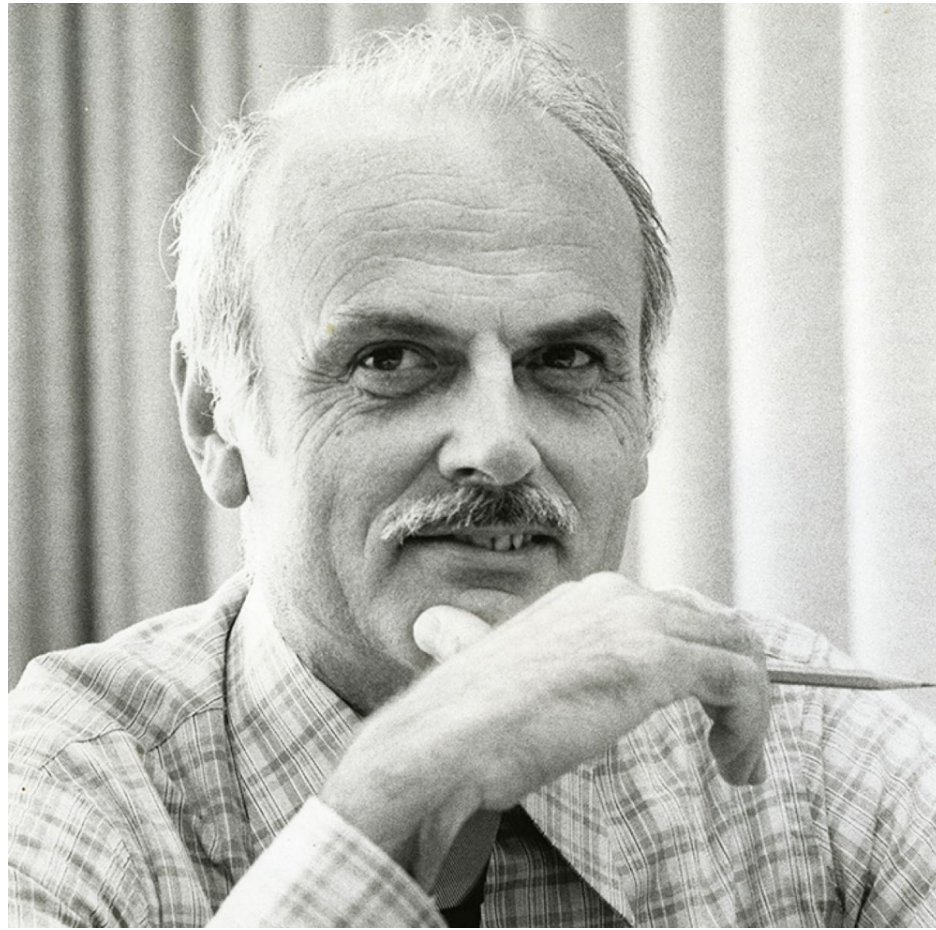
Welcome to the rectangular world  of **tabular database systems**.

Reshaping all kinds of data to fit into rows + columns may appear restrictive and arcane (it certainly did to me).

This course will investigate how **tables**...

- ... are a **versatile and flexible data representation**
(also for non-tabular data               ),
- lead to **compact data storage**
(in main memory as well as secondary memory), and
- admit **super-efficient data processing**.

Tables are Versatile



Edgar F. Codd (© IBM)

“Twelve distinct ways to represent data at the logical level are eleven too many.”

— Edgar F. Codd, ACM Turing Laureate

Tabular Database Systems: The Payoff is Substantial

Q: “But are **1** (encode), **2** (load), and **5** (decode) worth it?”

A: “*You bet!*”

- The DBMS-internal data format is **compact** (compression).
- **SQL is a declarative language** whose primary data type are tables. SQL queries **3**+**4** tend to be concise, often elegant.
- DBMS internals rely on regular table structures: **SQL is very efficient**, typically *way* faster than handcrafted programs.
- DBMS **coordinates concurrent access** to its tables—many users may operate on the same tables using well-defined semantics.
- DBMS **safely persists data** under its control, preventing data loss, e.g., through bugs in apps or system outages.

2 | This Course (Tabular Database Systems, short: TaDa or 🎉)

- We will focus on **table-centric data representation and processing**. There are other kinds of DBMSs (e.g., graph-based DBMS: $\mathbb{G} \rightarrow \mathbb{G}$, $\mathbb{A} \rightarrow \mathbb{A}$), but we will ignore those.
- Whenever we can approach the theory or **pragmatics** of a concept, we typically choose the latter. (Yet, tabular database systems are rooted in a rich and elegant mathematical foundation.)
- We will get our hands dirty using the tabular DBMS **DuckDB** 🎧 and its extensive **SQL** dialect.
- We will explore selected aspects of **DuckDB's internals** (implementation techniques used inside the `{}`).
- We will draw from a variety of data 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/TaDa 

- 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-create-table.sql`).
- **NB.** Code and extra material provide essential content (e.g., details on SQL syntax and semantics).
 -  +  = : Only slides + code provide a complete picture.

Material

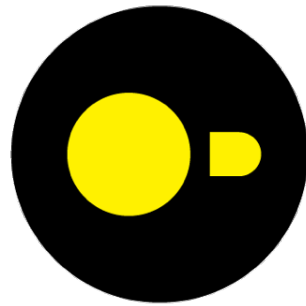
This course is *not* based on a single textbook. Rather, we build on

- 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 (quizzes, puzzles, 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.¹

¹ Run the DuckDB CLI in the web browser: <https://shell.duckdb.org>. Suffices for quick SQL experiments.

3 : The Tabular Data Model

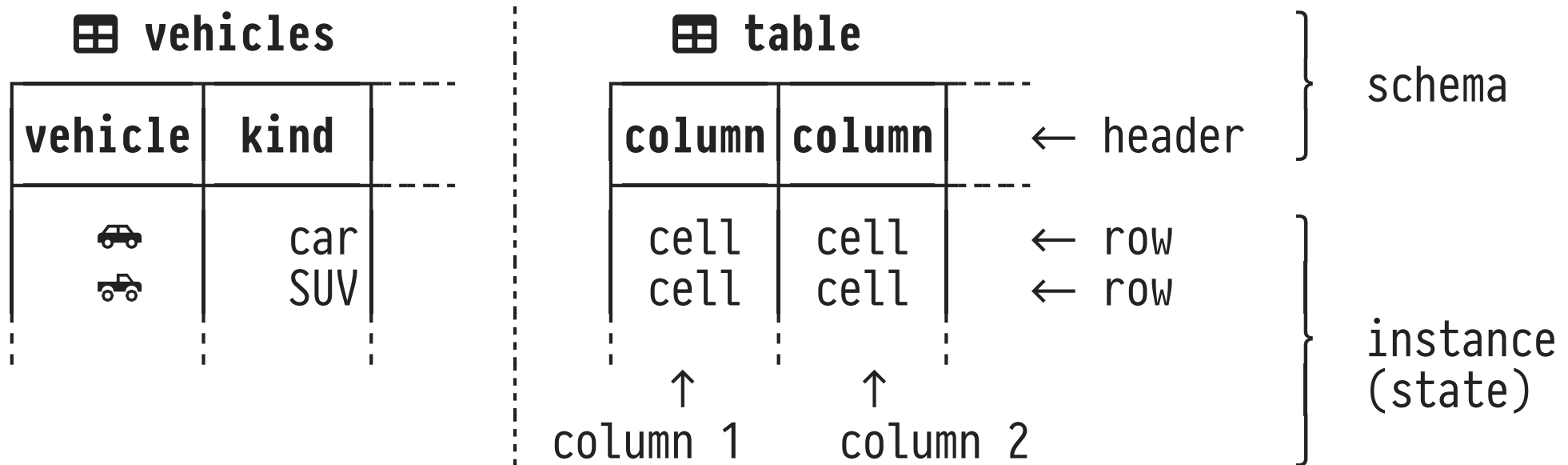
The **tabular data model** arranges all information in a rigorous grid-like fashion. Consider table `vehicles` below:

 `vehicles`

vehicle	kind	seats	wheels?
	car	5	true
	SUV	3	true
	bus	42	true
	bus	7	true
	bike	1	true
	tank	□	false
	cabrio	2	true


- This **table** has four **columns**, seven **rows**.
 - Table and columns are **named**, rows are not.
 - Each column holds **cell** values of a single **type**.
 - Cell value **NULL** (□) signals absence of information.
- Tables are vertical: feature **few columns** (≤ 20 is typical), but may contain **large numbers of rows** ($10-10^6$ rows are typical).

Tabular Data Model: Terminology



- Each **column** $i \in \{1,2,\dots\}$ is assigned a **name** c_i and its **type** τ_i .
- **Types** τ_i are simple (**text**, **int**, **boolean**, **float**, **int[]**, ...): **first normal form** or **1NF** (columns may *not* hold nested tables).
 - *All* cell values in column i are of type τ_i (**NULL** if allowed).
- The **schema** determines the structure of a table (say t): $t(c_1 \tau_1, \dots, c_n \tau_n)$, in short: $t(c_1, \dots, c_n)$. Typically *fixed*.
Ex: **vehicles(vehicle text, kind text, seats int, "wheels?" boolean)**.
- The **instance** is the bag of rows held in the table. *Dynamic*.

SQL: Creating Tables

- The SQL statement `CREATE TABLE $t(c_1 \tau_1, \dots, c_n \tau_n)$` creates
 - a table with given **schema** $t(c_1 \tau_1, \dots, c_n \tau_n)$ and
 - an **empty instance** (no rows yet).
- In the DuckDB  CLI:

optional
table name
└──────────┘
↓

```

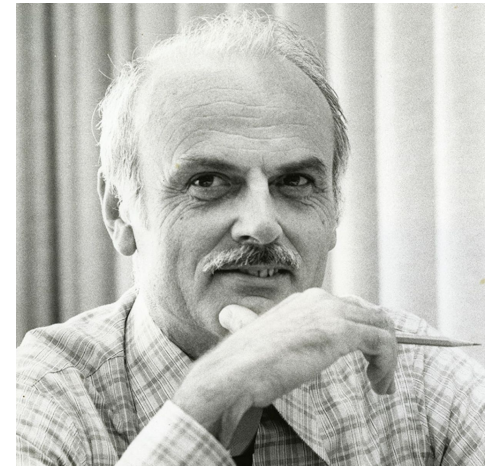
D CREATE OR REPLACE TABLE vehicles (
  vehicle  text  NOT NULL,
  kind     text  NOT NULL,
  seats    int   ,
  "wheels?" boolean NOT NULL
);
  
```

↑
↑
↑
columns
types
constraints [optional]

refers to code/extra material in a file named 001- →*

Tables vs. Relations ①

Origin and foundation of the tabular data model is the **relational model** (Edgar F. Codd, 1970).



Edgar F. Codd (© IBM)

Tables and **relations** (subsets of the Cartesian product of domains) are closely related.

- Relation $<$ (*less than*) on integers:

$$< \equiv \{ (x,y) \mid x \in \mathbb{Z}, y \in \mathbb{Z}, \exists n \in \mathbb{N}, n \neq 0: x+n = y \} \subseteq \mathbb{Z} \times \mathbb{Z}$$

- Notation: $(x,y) \in <$ or $\underline{\leq(x,y)}$ or $x < y$, e.g., $<(0,42)$

- Can define relations **intensionally** (as above) or **extensionally**.
Game *Rock, Paper, Scissors* (finite domain $RPS = \{\text{✊}, \text{✋}, \text{✂}\}$, with extensional “beats” relation: $x < y \Leftrightarrow y \text{ beats } x$):

$$< \equiv \left. \begin{array}{c|c} \mathbf{x} & \mathbf{y} \\ \text{✊} & \text{✊} \\ \text{✋} & \text{✋} \\ \text{✂} & \text{✂} \end{array} \right\} \begin{array}{l} \text{extent of } < \text{ (a subset of } RPS \times RPS\text{):} \\ \text{lists the set of all ordered pairs included in } < \end{array}$$

Tables vs. Relations ②

relation <

$\{$

 $(\text{✊}, \text{✋}),$

 $(\text{✋}, \text{✊}), \leftarrow$ tuple (pair)

 $(\text{✋}, \text{✋}) \}$

$\underbrace{\hspace{10em}}$

set of tuples

 $(\subseteq \text{RPS} \times \text{RPS})$

- $\text{✊} \in \text{RPS}, \text{✋} \in \text{RPS}, \text{✋} \in \text{RPS}$
- tuple components ordered
- no order among tuples
- no duplicate tuples

- Read symbol $::$ as “has type”.

beats

lose	win
✊	✋
✋	✊
✋	✋

\leftarrow row

$\underbrace{\hspace{10em}}$








bag (multiset) of rows

- **lose** $::$ text, **win** $::$ text
- columns ordered
- no order among rows
- possibly duplicate rows

Identifying Rows in Tables: Keys


Tables are unordered: cannot refer to rows by position (~~1st/2nd/...~~ row). Instead, use **cell values that uniquely identify its row.**

 vehicles

vehicle	kind	seats	wheels?
	car	5	true
	SUV	3	true
	bus	42	true
	bus	7	true
	bike	1	true
	tank	0	false
	cabrio	2	true

↑ ↑ ↑ ↑
 unique ~~unique~~ ? ~~unique~~

- A **key** is a column (combination) whose values identify rows:

*“Return the row with a **vehicle** value of .*”
- **vehicle** is a key for **vehicles**. Columns **kind** or **wheels?** are not.
- ?? Column **seats** is not:
 1. column contains **NULL** and
 2. will values remain unique once more vehicles are added?

- A table *t* may contain multiple candidate key columns. Choose one of these candidates to be the **PRIMARY KEY** for *t*.

Identifying Rows in Tables: Good Keys

Choosing good primary keys is vital in **database schema design**:








1. Key values identify rows uniquely in **any table state** (no matter how many rows will be added in the future).
 - Key values should be immutable during row lifetime.
2. Key columns will be used in **WHERE** filter predicates:
 - Aim for **narrow keys** (ideally: single-column).
 - Prefer key columns whose values can be **compared efficiently** (e.g., prefer type **int** over **text**, **date**, or arrays).
3. Consider **introducing additional/artificial key columns** if the application domain does offer natural/narrow/efficient keys.
 - Examples: book ISBNs, product EANs, social security IDs.

Foreign Keys: Identifying (or: Pointing to) Rows in *Other* Tables





Splitting data between tables helps

- to **avoid redundancy** (and thus data integrity issues), and
- to design tables focused on specific **domain concepts**.

vehicles

vid	vehicle	kind	seats	wheels?	driver
v_1		car	5	true	p_4
v_2		SUV	3	true	p_4
v_3		bus	42	true	□
v_4		bus	7	true	□
v_5		bike	1	true	p_2
v_6		tank	□	false	p_3
v_7		cabrio	2	true	p_4

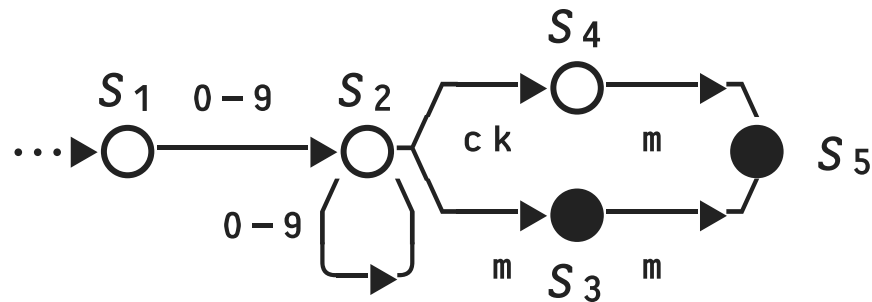
peeps

pid	pic	name	born
p_1		Cleo	2013
p_2		Bert	1968
p_3		Drew	□
p_4		Alex	2002

↑
foreign key (references rows in vehicles, \subseteq peeps(pid))

- Notation: `vehicles(vid, vehicle, ..., driver → peeps)`.

Example: Two Tables Represent a Finite State Machine (FSM)



- FSM accepts metric lengths: '2mm', '135km', '3cm', '42m'.
- Character labels define deterministic transitions.

states

state	start?	final?
s ₁	true	false
s ₂	false	false
s ₃	false	true
s ₄	false	false
s ₅	false	true

transitions

from	to	labels
s ₁	s ₂	[0,...,9]
s ₂	s ₂	[0,...,9]
s ₂	s ₃	[m]
s ₂	s ₄	[c,k]
s ₃	s ₅	[m]
s ₄	s ₅	[m]

Q: How would you define the **keys**?

Q: Identify **foreign keys** (if any).

Q: Provide SQL **CREATE TABLE** statements for both tables.

4 : DuckDB? 🦆?

In case you were wondering:

DuckDB has been named after *Wilbur*, the *Duck*, which has been living as a pet with Hannes Mühleisen²—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)

² Hannes originally is from the Stuttgart area. Back then his car had the license plate **S:QL 1337**.

Tabular Database Systems

②

Tabular Data in CSV Files

March 17, 2026

Torsten Grust
Universität Tübingen, Germany

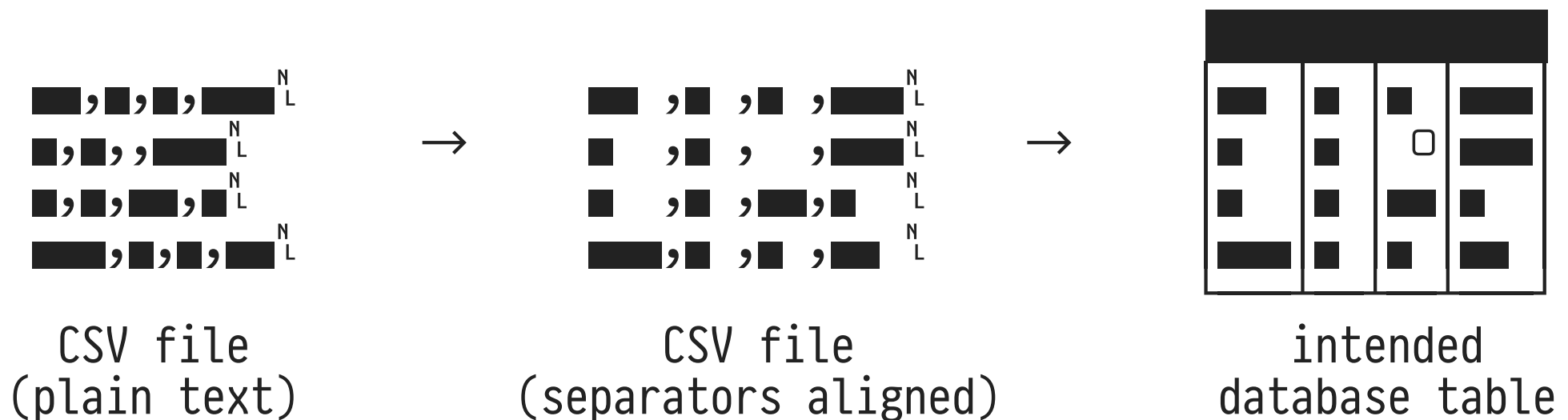
1 | Most Data Lives Outside DBMSs

Reality check: The vast majority of data out there does *not* live in files and formats directly accessible for DBMSs:

- Applications read/write their specific data formats which were probably devised before DBMS *X* became available.
- Users process their data using a variety of systems. Why lock-in to the proprietary database format of DBMS *X* (or *Y* or 🐧)?
 - Rather build on a “lowest common denominator” format that allows **data interchange between systems**.
- Simplicity wins! DBMS-internal data formats are intricate, heavily optimized for access speed and space efficiency.
 - Rather use formats that admit **simple read/write routines**.
 - Let **humans read/write the data format**, ideally using no tools but their favorite **plain text** or spreadsheet editor.

2 : Tabular Data in Comma-Separated Value (CSV) Files ①

CSV files (*.csv) provide one such lowest common denominator DBMS-external data format. CSV encodes tabular data in **line-structured plain text** files:



1. Newline (NL) separates lines \equiv , a **line** encodes one table **row**.
2. Delimiters (often comma $,$) split lines into **columns**.
 - We expect each line to hold the same number of columns, leading to a rectangular grid of cells $\begin{matrix} \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare & \blacksquare \end{matrix}$.

Tabular Data in Comma-Separated Value (CSV) Files ②

CSV file `007-vehicles.csv` to define the column names and contents of table `vehicles` (see Chapter ①):

```
vehicle, kind, seats, wheels?  
🚗, car, 5, true  
🚙, SUV, 3, true  
🚌, bus, 42, true  
🚌, bus, 7, true  
🚲, bike, 1, true  
🚛, tank, , false  
🚗, cabrio, 2, true
```

 #007

- **Notes:**

- Column names held in first CSV line “by convention.”
 - Otherwise, **header and data rows are identical.**
- Columns are **untyped** (`wheels?` “looks like” type `boolean`).
- `NULL` represented by adjacent delimiters(`,` `,`).

Reading CSV Files Like Database Tables

- In DuckDB, read a CSV file much like a database-internal table:

```
┆ FROM '007-vehicles.csv';
```

vehicles	kind	seats	wheels?
----------	------	-------	---------

- This is a shorthand for an explicit call to built-in function `read_csv()`:

```
┆ FROM read_csv('007-vehicles.csv', configuration parameters);
```











vehicles	kind	seats	wheels?
----------	------	-------	---------

- Non-default *configuration parameters* may be required to disambiguate CSV file contents regarding presence of headers, choice of delimiters, column types, ...)

3 | CSV Enables Data Exchange

A variety of systems export data in CSV format. Paired with the DBMS's CSV import, this enables **CSV-based data exchange**.

- Example: Spreadsheets. On Google Sheets, use command **File ▶ ⌵ Download ▶ Comma-Separated Values (.csv)**:

	A	B	C	D
1	vehicle	kind	seats	wheels?
2		car	5	<input checked="" type="checkbox"/>
3		trolley	40	<input checked="" type="checkbox"/>
4		trolley		<input checked="" type="checkbox"/>
5		truck	3	<input checked="" type="checkbox"/>
6		helicopter	6	<input type="checkbox"/>
7		bus	42	<input checked="" type="checkbox"/>
8		tractor	1	<input checked="" type="checkbox"/>
9		boat	4000	<input type="checkbox"/>
10		scooter	2	<input checked="" type="checkbox"/>
11		bike	1	<input checked="" type="checkbox"/>

[Access this spreadsheet on Google Sheets](#) ▶

- **NB.** On export, **data validation rules** (e.g., number ranges, value constraints) , formulæ, or cell links **are lost**.

4 : Importing Data Into Tables, Exporting Query Results Into Files

1. Save repeated imports, use **COPY** to **copy a file into a table**:

```
COPY table  
FROM path_to_file (configuration parameters)
```

- Can now query *table* like a regular database table.

2. **Copy SQL query results into a file** external to the DBMS:

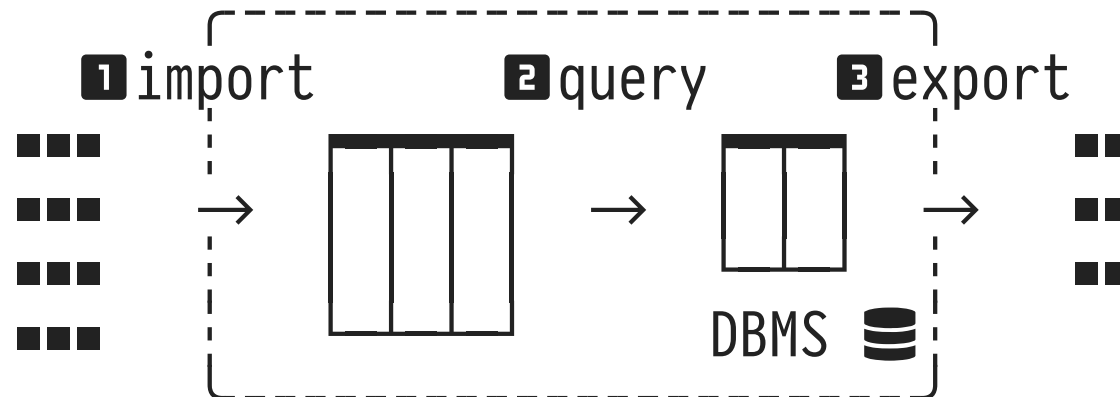
```
COPY (query)  
TO path_to_file (configuration parameters)
```

The *configuration parameters* determine file format (e.g., CSV), delimiters, header output, compression, overwrite/append mode, ...

Using DuckDB as a CSV Processor

DBMS may act as a processing engine for database-external data:

- + SQL is versatile, expressive, and executed efficiently.
- Need to pay for data import and export.



 #010

```

1 D COPY (SELECT ...
      FROM   input_file
      WHERE ...
      ORDER BY ...)
    TO   output_file;
  
```

```

D COPY table
FROM input_file;
2 D COPY (SELECT ...
      FROM   table
      WHERE ...
      ORDER BY ...)
    TO   output_file;
  
```

5 : Is CSV Too Simple? Ubiquitous, But Fragile

The CSV format is documented in [RFC 4180](#) . Still, CSV files in daily practice deviate in a variety of ways¹, including

- choice of **delimiters** (frequent: `,` `;` `|` (pipe) `␣` (tab)),
- **quoting** (enclose columns with significant spaces `_` in `"` or `'`),
- character **escape** sequences (what if `,` `"` `␣` are in columns?),
- absence/presence of **headers or preamble text** above data rows,
- choice of representation of **missing values** (e.g., `...`, `,,`), or
- a **mixture of data types** in a column.

DuckDBs aims to auto-adapt to these CSV dialects using a “multi-hypothesis sniffing” algorithm².

¹ DuckDB's CSV reader comes with 25+ configuration options. To get an impression of the variety of CSV dialects out there, see the paper [Characteristics of Open Data CSV Files](#)  (2016).

² See the blog post [DuckDB's CSV Sniffer: Automatic Detection of Types and Dialects](#)  and the paper [Multi-Hypothesis CSV Parsing](#)  (2017) by Hannes Mühleisen *et al.*

DuckDB's CSV "Sniffer"

1. **Sample 20480 lines** from all over the file (unless we need to read top to bottom, e.g., if input is `stdin` or compressed).
2. Try **CSV dialects** based on 24 selected combinations of the options below. Choose the CSV dialect options that lead to a *consistent and maximum number of columns* per row:

parameter	options
<code>delim</code>	<code>,</code> <code>␣</code> ;
<code>quote</code>	<code>"</code> <code>'</code> (empty)
<code>escape</code>	<code>"</code> <code>'</code> <code>\</code> (empty)

3. **Detect types.** Try to cast values to candidate types *in order*: `[(NULL,) boolean, int, double, time, date, timestamp, text]`.
4. **Header detection:** Do columns in first row match detected types (*is data*, thus generate own header) or not (*is header*)?

DuckDB's CSV "Sniffer": Dialect and Type Detection

Dialect detection:

 #011

flights.csv	# columns/delimiter			
	,		H T	;
FlightDate Carrier Origin Destination ^N _L	1	4	1	1
1988-01-01 AA New York, NY Los Angeles, CA ^N _L	3	4	1	1
1988-01-02 AA New York, NY Los Angeles, CA ^N _L	3	4	1	1
1988-01-03 AA New York, NY Los Angeles, CA ^N _L	3	4 ↑	1	1

Type detection (file reading order →):

band.csv	candidate types for	
	column 0	column 1
Name, Age	(skip)	(skip)
,	[NULL,boolean,...,text]	[NULL,boolean,...,text]
Mike Lindup, 65	[text]	[int,double,...,text]
Mark King, 66.2	[text]	[double,...,text]
	↓	↑

Tabular Database Systems

③




Reading Data at the Speed of ~~Light~~ Memory

March 17, 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  / 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 its 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 (📁 Secondary/📁 Primary)	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.**
 - **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	#012
	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 ./012-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 with open(sys.argv[1], 'r') as file: for line in file: sum = sum + int(line.split(' ')[4]) print(sum)</pre>		<pre>reset sum</pre>	 #013
		<pre>open file (reading) read line by line extract 5th col, cast to int, add output</pre>	

- 🕒 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) {  
    delim = line;  
  
    for (column = 1; column < 5; column++) {  
        delim = strchr(delim, '|');  
        delim++;  
    }  
  
    sum = sum + atoi(delim);  
}
```

 #014

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%

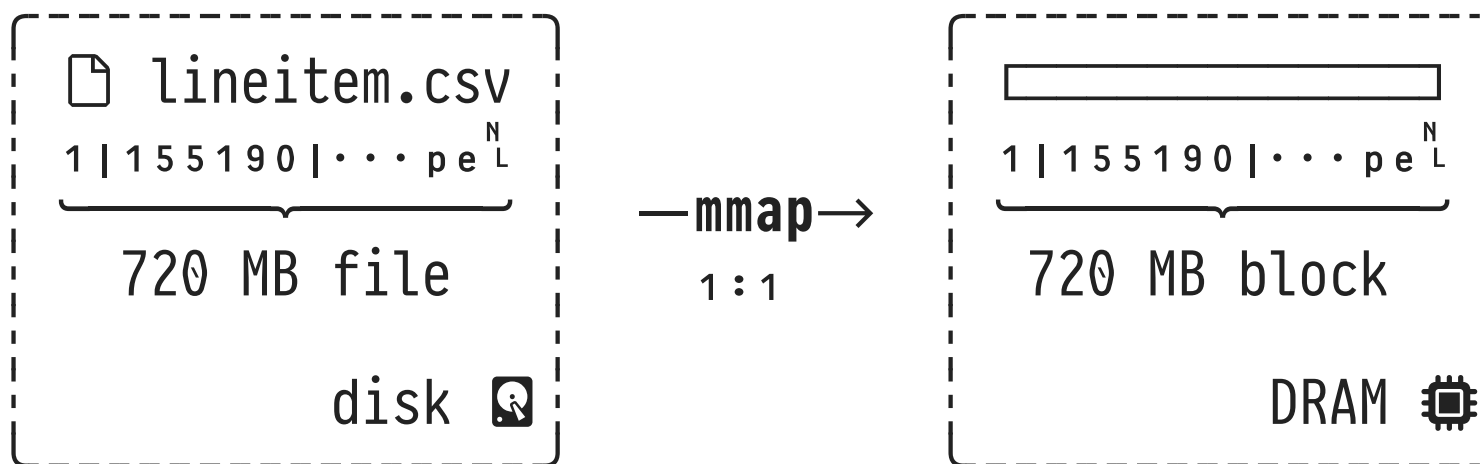
- Reading the CSV file line by line is too slow.
 - 💡 Read entire file at once, impose line structure ourselves.

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

 #015

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
(all other C code)	(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

 #016

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 `nl` 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 & 0x00008000000000000000ULL) { data += 7; continue; }  
    data += 8;  
}
```

 #017

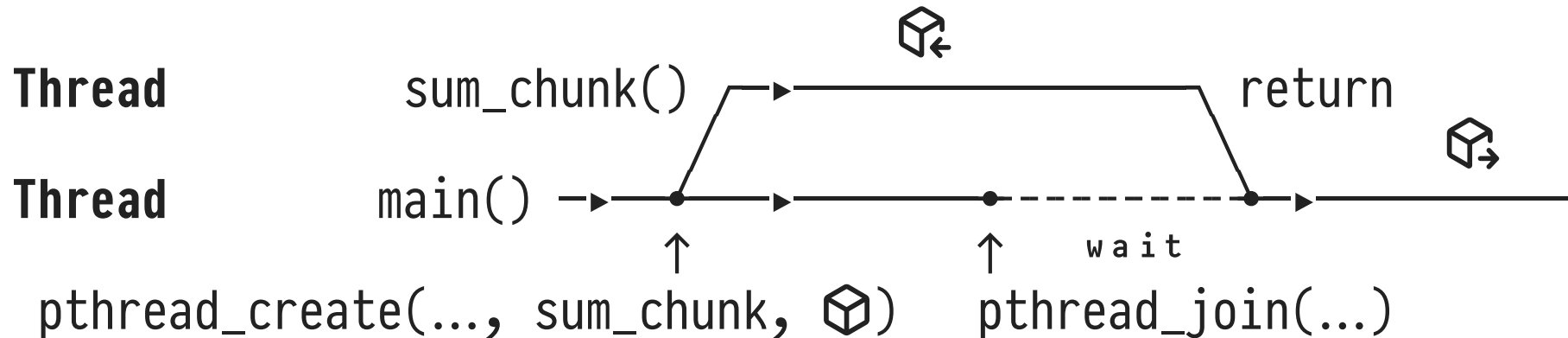
Sum of Quantities ⑤ — C with `mmap(2)` + Block-Wise `ℓ` Search

- 🕒 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.

- ? 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 ⑥ 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):

 #018

```
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 : Interlude: SQL Aggregate Functions

SQL aggregates⁵ condense a bag of rows into a *single* value:

 vehicles

vehicle	kind	seats	wheels?
	car	5	true
	SUV	3	true
	bus	42	true
	bus	7	true
	bike	1	true
	tank	□	false
	cabrio	2	true

- Most aggregates ignore (= skip) NULL (□) values.
 - An optional **FILTER** clause can control value inclusion.
- Order-aware aggregates use clause **ORDER BY** to be deterministic:

`list(vehicle ORDER BY seats)`

`bool_and("wheels?")` ≡ false
`max(seats)` ≡ 42
`arg_max(kind, seats)` ≡ 'bus'
`count(vehicle)` ≡ 7

⁵ DuckDB's SQL dialect features an extensive list of [built-in aggregate functions](#) ↗.

Interlude: SQL Aggregate Functions

- Let us use annotations to explain SQL constructs (here: query clause **row cardinality**, **1 2 3**: clause processing order⁶):

<pre>SELECT <i>expr</i>,...,<i>expr</i> FROM <i>t</i> WHERE <i>p</i></pre>	<pre> 3 <i>m</i> 1 <i>n</i> (= <i>t</i> , cardinality of <i>t</i>) 2 <i>m</i> (≤ <i>n</i>) └── # rows returned by query clause</pre>
---	--

- SQL aggregates reduce row cardinality to 1:

<pre>SELECT <i>agg</i>,...,<i>agg</i> FROM <i>t</i> WHERE <i>p</i></pre>	<pre> 3 1 1 <i>n</i> (= <i>t</i>) 2 <i>m</i> (≤ <i>n</i>) └── # rows</pre>
---	---

- ⇒ **SELECT** clause *cannot mix* scalar *expr* and aggregates *agg*.

 #019

⁶ SQL clause processing order ≠ SQL syntactic order. DuckDB implements a “friendly SQL dialect”  to partly rectify this (e.g., allowing FROM-SELECT queries). We'll use such DuckDB-specific features sparingly.

10 : Sum of Quantities ⑦ — SQL

Aggregate function `sum()` is what we need to formulate a SQL variant of the benchmark query over the CSV file:⁷

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

#020

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**.

⁷ 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**⁸ that exposes details of the system's **query evaluation plans**:

- Supports query and performance debugging.
- **Shows order of query operations** explicitly, making effects of query optimization visible (e.g., projection pushdown).
- **Measures 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;
```

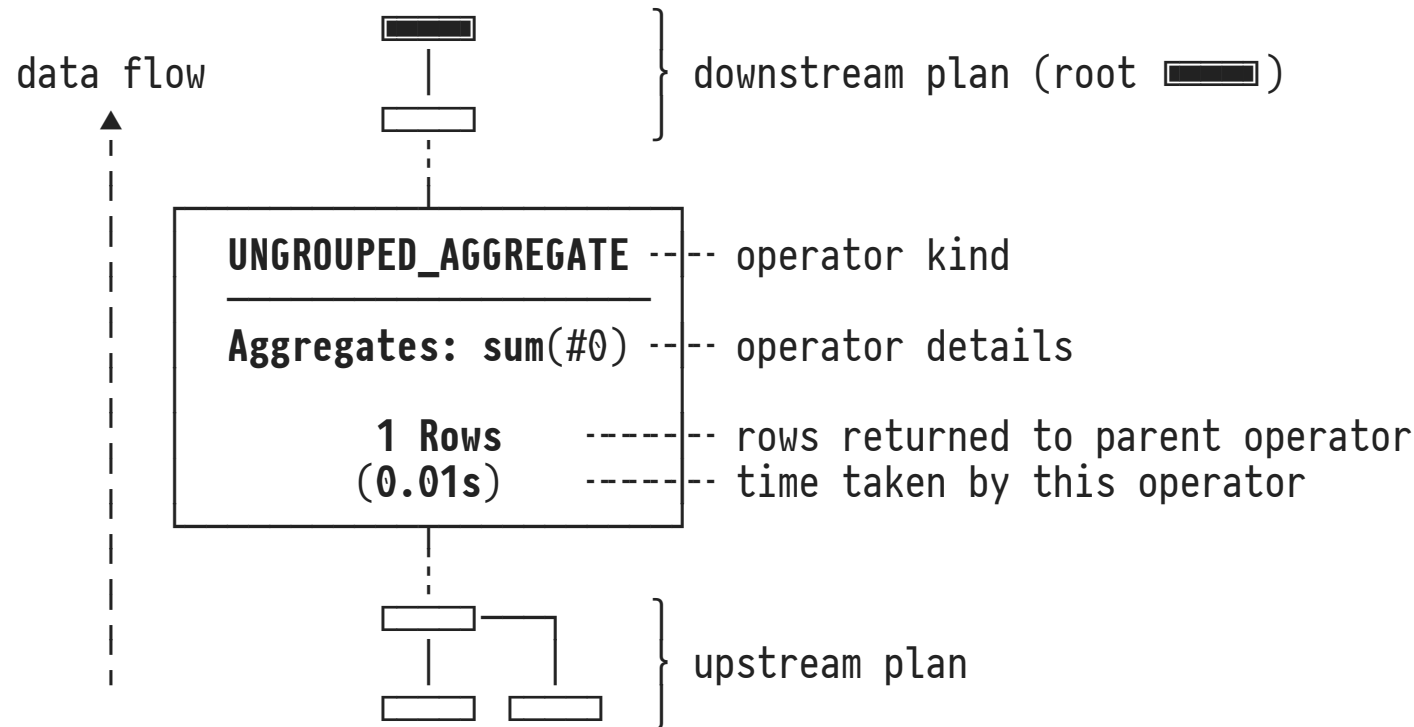
Actually run query. Measure times/rows, annotate plan.

⁸ 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**:⁹


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



⁹ 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

```
D SELECT SUM(l_quantity)
   FROM   lineitem;
```

 #020

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
- the data for table `lineitem` does reside in DRAM (not in secondary storage).

Tabular Database Systems

④

Columnar Table Storage

March 17, 2026

Torsten Grust
Universität Tübingen, Germany

1 | Transient vs. Persistent Databases

- **Transient in-memory databases**

By default, tables are created and stored in volatile **DRAM** 🧠. When the DuckDB process ends, all table data is lost.

```
$ duckdb
  |
Connected to a transient in-memory database.
D |
```

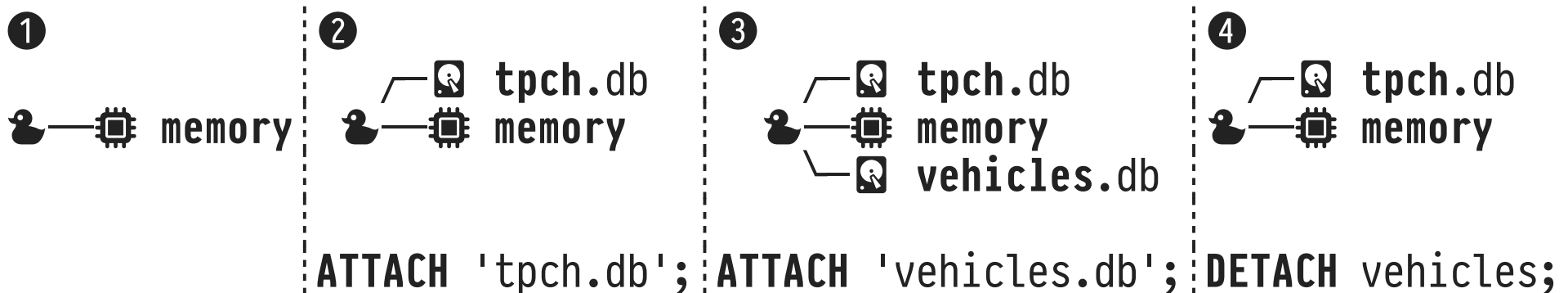
- **Persistent on-disk databases**

Table data is held in a **file on secondary storage** 📁¹. When the DuckDB process ends and later restarts, the current state of the database **persists** (= is preserved between DBMS sessions).

¹ DuckDB stores the instances of all tables and associated data in a self-contained *single regular file*, conventionally named `<database>.db`, ready to be posted/shipped. Other DBMSs (like PostgreSQL) store databases in subdirectory structures or directly in disk blocks (“raw” storage, no filesystem).



One DBMS, Multiple Databases

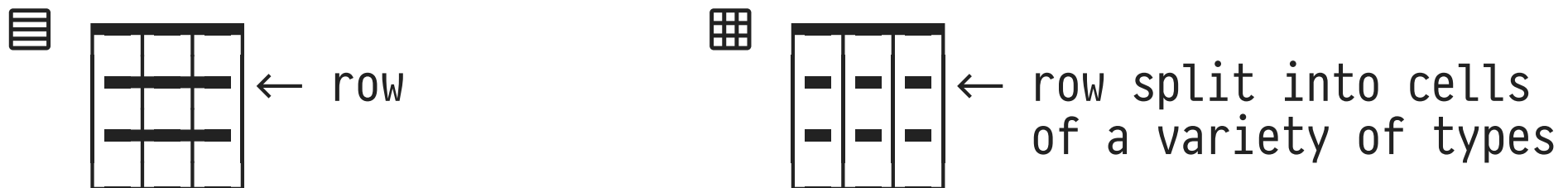
- A DuckDB process can *attach* to **multiple databases** at one time.
 - A **database** comprises
 - schema and state (i.e., bag of rows) of all tables,
 - other user-defined objects (types, views, macros), and
 - installed DBMS extensions.





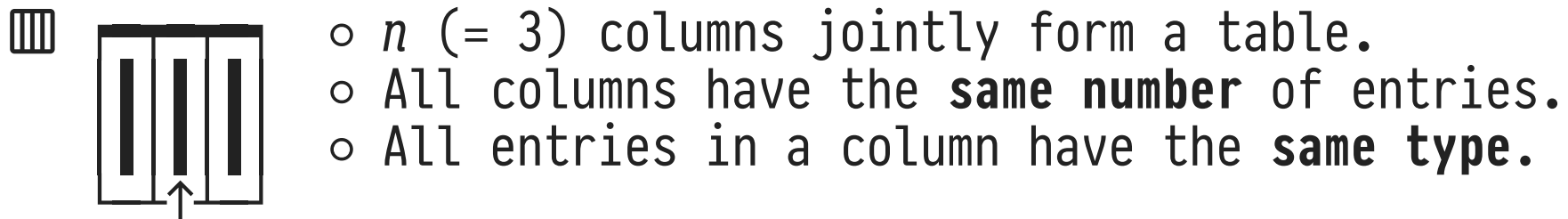
- Objects (e.g., tables) live in database-owned **namespaces**:
 - Name `t` refers to table `t` in the default database (see [USE](#)).
 - Name `tpch.t` refers to table `t` in database `tpch`.

2 : Table: A Collection of Rows or Columns?

- The tabular data model appears to have its focus on **rows**:
 - : Instances of tables are bags of **rows** **— — —**.
 - : Each **row** is split into individual cell values **- - -**.
 - SQL clause **FROM** *t* iterates over the **rows** of *t*, evaluates **WHERE/SELECT** clause per **row**, **ORDER BY/LIMIT** sort/count **rows**.










- Yet, DBMSs like  organize the **storage** of tables by **column**:
 - : Values in a **column** appear consecutive in memory/on disk.



Queries That Focus on Few Columns (But Read All Rows)

- Observation: a large class of SQL queries touch **few columns** but (almost) **all rows** of a table. These are known as **OLAP² queries**.
 - The benchmark query over table `lineitem` from Chapter 03.
 - Aggregation queries scan all entries of their source column:

vehicles

vehicle	kind	seats	wheels?
 ○	car ○	5 ✓	true ○
 ○	SUV ○	3 ✓	true ○
 ○	bus ○	42 ✓	true ○
 ○	bus ○	7 ✓	true ○
 ○	bike ○	1 ✓	true ○
 ○	tank ○	0 ✓	false ○
 ○	cabrio ○	2 ✓	true ○

Q: **SELECT** max(seats)
FROM vehicles

- ✓ relevant cells (25%)
- no contribution to result

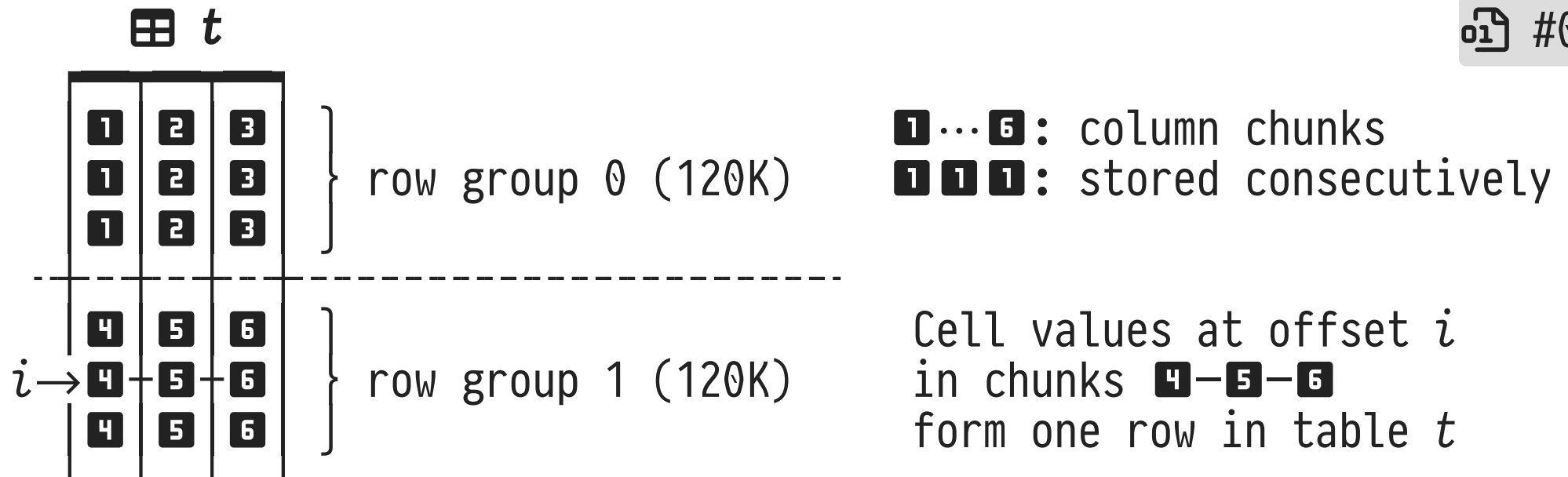
- Only read column **seats** in `TABLE_SCAN` in plan for query Q.
- Row-wise table storage would read all cells (but ignore 75%).

² Online Analytical Processing, as opposed to OLTP (Online Transactional Processing) queries which touch few (even single) rows but read all those rows' columns.

3 : DuckDB : Columnar Table Storage



DuckDB storage format for table t with r rows and c columns:

1. **Row groups:** Partition table t horizontally into groups of 120K ($120 \times 1024 = 122880$) rows $\Rightarrow t$ has $\lceil r / 120K \rceil$ row groups.
2. **Column chunks:** Inside each row group, store a chunk (of 120K cell values) for each of the c columns.



4 : Column Compression

Tables may exhibit repeating cell values or redundancies across rows that make instances worthwhile targets for **compression**.

- Compression (and decompression) **adds CPU effort** but reduces the storage footprint and may thus **save disk I/O bandwidth**.
Overall, performance may improve.
 - DuckDB  applies compression to on-disk databases only.
- **Columnar table storage works well with compression:**
 - Within a column (chunk) **1 1 1**, values are homogenously typed and potentially similar. This aids compression algorithms.
 - (In row-wise storage, values of different types are interleaved **1 2 3**, leading to lower compression rates.)
- For some tables, the effect of compression can be drastic. 

Compression: General vs. Lightweight

1. **General-purpose compression** algorithms (e.g., *gzip*, *zstd*):
 - Detect and exploit **patterns in arbitrary bit sequences**:
 - `11111111`: predictable, compressible as `8×1` (run length).
 - `10110010`: noisy, random, not compressible.
 - High compression rates, but (de)compression is **costly**.
 - Work best on **large chunks of data** (>> 256kB)—decompressing sizable chunks renders accessing individual rows expensive.
2. DuckDB builds on a family of **lightweight compression** schemes:
 - Detect **patterns in typed data** in a column chunk:
 - `10 12 9 10 8 11`: compressible as `8 ⊕ 2 4 1 2 0 3` (values close to reference `8`).
 - (De)compression is **cheap** and incurs light CPU load only.
 - Effective on small data: DuckDB scans column chunks to select best compression algorithm **per row group** (120K rows).

DuckDB: Lightweight Compression Schemes ①

- **Constant Encoding:**

- Applies if **every value** in column chunk is the **same**.
- Example: ranges of **NULLs** or rarely changing values (e.g., **year** in the timestamp of log entries).

Uncompressed	Constant Encoding
2025 2025 2025 2025 ...	2025

- **Run-Length Encoding (RLE):**

- Compress **groups of repeated values** using *(count,value)* pairs.
- Example: sorted or partitioned data.

Uncompressed	Run-Length Encoding
a, a, a, a, b, b, c, c, c	(4,a), (2,b), (3,c)

DuckDB: Lightweight Compression Schemes ②

- **Bit Packing:**

- Exploit that values do **not span full domain of their type**.
- Max value determines bit width, then elide leading 0 bits.

Uncompressed (type smallint, int2: 16 bits)			Bit Packing (6 bits)
00000000000101010	0000000000010011	0000000000000100	101010, 010011, 000100
42	19	4	

- **Frame of Reference (FOR):**

- Store **Δ s from a reference (minimum)**, not absolute values.
- Example: dates close to one point in time. Absolute values are days since 1970-01-01, Δ s to min date will be smaller.

Uncompressed	FOR Encoding
1968-08-26, 1968-08-24, 1968-08-27	$\underbrace{1968-08-24}_{\text{reference}} \oplus \underbrace{2 \ 0 \ 3}_{\Delta s}$

DuckDB: Lightweight Compression Schemes ③

- **Dictionary Encoding:**

- Place **frequent values in dictionary**, only store # of dictionary entry. Effective if values are wide (strings).

Uncompressed	Dictionary Encoding
Zelda, Mario, Mario, Zelda, Zelda	[₀ Zelda, ₁ Mario] (dictionary) 0, 1, 1, 0, 0 (entry #s)

- **Fast Static Symbol Table Encoding (FSST):**

- Extends dictionary encoding to capture **frequent substrings**.
- Example: URLs or e-mail address strings.

Uncompressed	FSST Encoding
www.archive.org, www.duckdb.org	[₀ www, ₁ .org, ₂ archive, ₃ duckdb] (symbol table) (0 2 1), (0 3 1) (entry #s)

- **Compression of IEEE 754 Floating-Point Values (ALP³).**

 #023

³ XORing similar `doubles` typically leads to a high number of leading/trailing 0 bits. See the paper [ALP: Adaptive Lossless Floating-Point Compression](#) (2023). ALP is implemented in DuckDB.

Tabular Database Systems

⑤

Database-External Data in Parquet Files 🗄️


March 17, 2026

Torsten Grust
Universität Tübingen, Germany

1 | Columnar Compressed Data Storage Outside the DBMS: Parquet

Fragility, space requirements, and parsing effort render CSV files problematic for applications that read/write GBs or TBs of data.

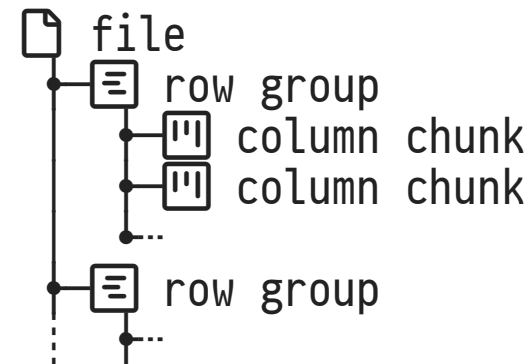
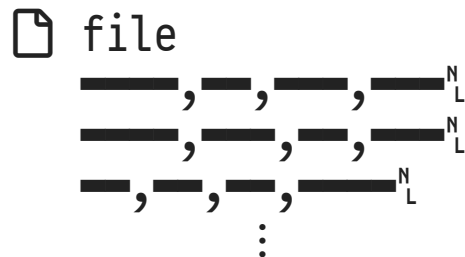
Alternative **database-external file formats** have been developed to address these issues. Among these, **Parquet**¹  is widely used.

- Developed in the open since 2013 (initiated by Twitter *et al*).
- Stores **columns of typed data**.
- Built-in **compression** (on file and column levels).
- Incorporates **rich metadata** that supports projection and selection pushdown.
- Supported by libraries for a wide range of programming languages. Directly readable/writable by DuckDB .

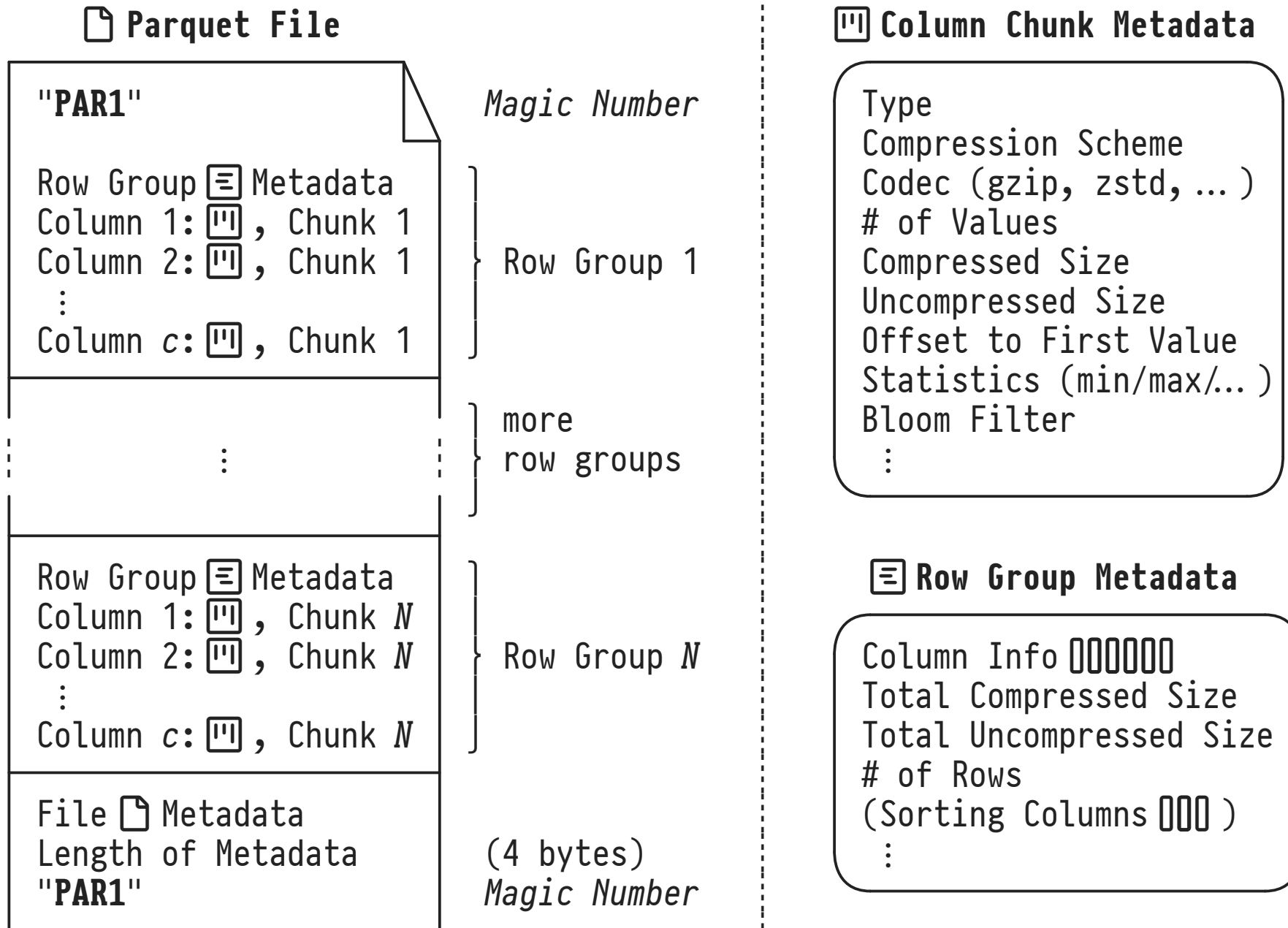
¹ [Apache Parquet](#)  (sponsored by the Apache Software Foundation): open source, column-oriented data file format designed for efficient data storage and retrieval.

CSV vs. Parquet

CSV	Parquet 📄
monolithic row-oriented (lines separated by \backslash) plain text, uncompressed	split into horizontal row groups column-based (chunk-by-chunk) supports compression: <ul style="list-style-type: none"> • file-level (gzip, zstd, ...) • column-level (dictionary, RLE, ...)
untyped (requires parsing)	typed: <ul style="list-style-type: none"> • scalar (INT32/64, FLOAT, BYTE_ARRAY, ...) [• nested records]
no metadata (optional header row)	metadata for file/row groups/columns: <ul style="list-style-type: none"> • file format version • min/max/count statistics • cardinality, (un)compressed byte sizes



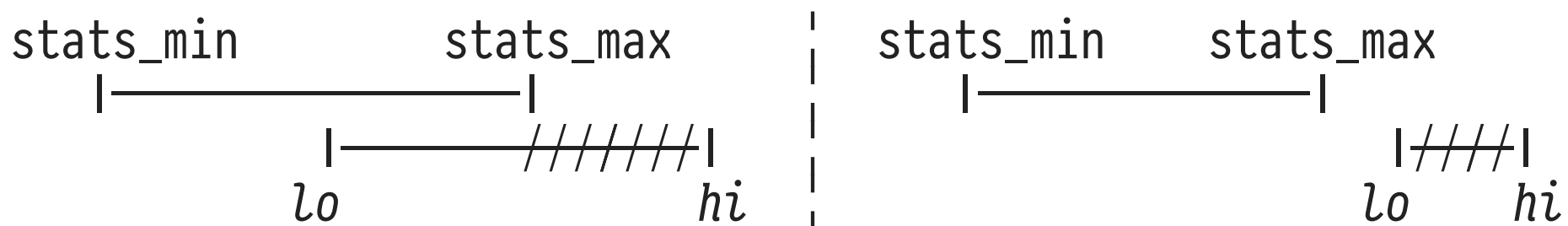
A Sketch of Parquet's Storage Format



2 : Pushing Projection and Filtering Down into Parquet Reading

Parquet's file structure + metadata enables readers (like DuckDB's `PARQUET_SCAN` operator) to **only access relevant data subsets**.

- **Projection pushdown:** Exploit **column-based layout**. In each row group, use `data_page_offset` to navigate the file to only read required column chunk(s).
- **Filter pushdown:** Exploit **statistics metadata**. Entirely skip row group if `stats_min/stats_max`² for column `c` indicate that filter predicates like `lo ≤ c ≤ hi` will always fail (///).

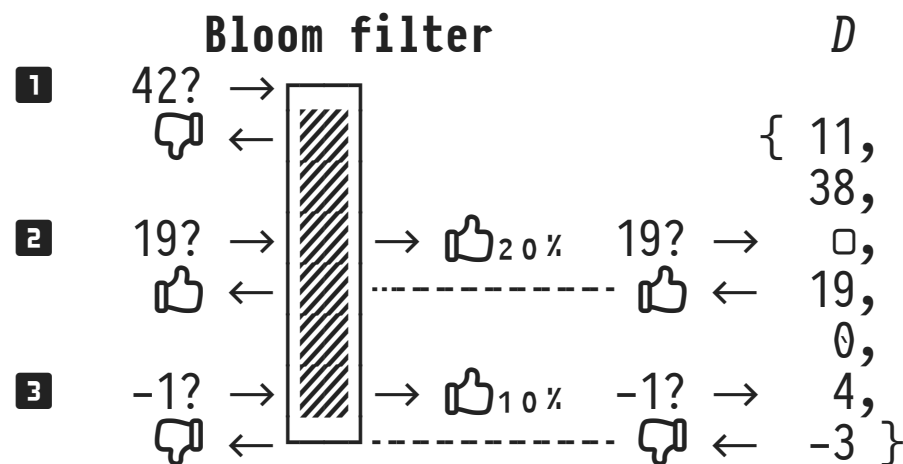


² In the DB research literature, the `(stats_min, stats_max)` pairs in all row groups are sometimes collectively referred to as **zone map** (for column `c`).

3 : Optional in Parquet Files: Bloom Filters

If column values are randomly shuffled, all values may occur in all row groups and the effectiveness of zone maps is largely lost.

- **Bloom filters** are *compact* data structures that *over-approximate* the set D of distinct values³ in a row group.
- Given the question $x \in D?$, Bloom filters either respond with
 - *definitely no* (👎) or
 - *probably yes* (👍 _{$p\%$} , where p is a false positive rate).

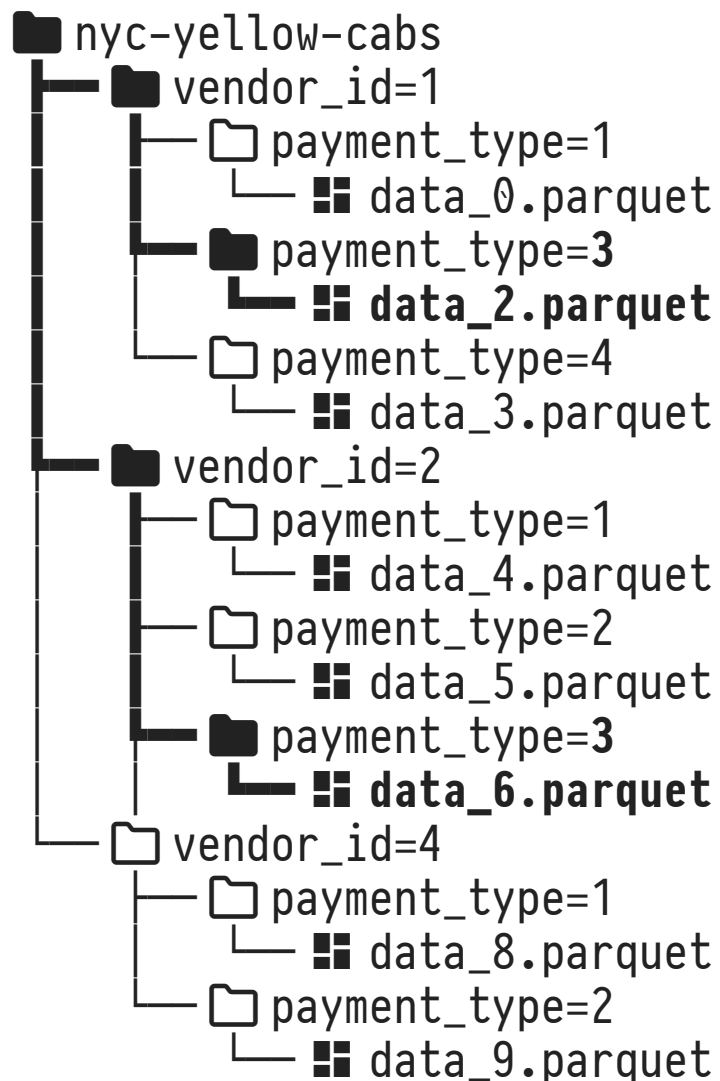


 #026

³ Set D is the **active domain** of the row group. Example: if the values in a column have type `int32` with domain $[-2^{31}, 2^{31})$, we typically have $D \subseteq [-2^{31}, 2^{31})$ (read \subseteq as “significantly smaller”).

4 : Partitioned Files and Filter Pushdown (Hive Partitioning)

If data volumes get truly large, it may make sense to **partition rows** into a family of files. Which partition will a row belong to?



- Choose $n \geq 1$ columns c_1, \dots, c_n as partition criteria (left: $n=2$, `vendor_id` + `payment_type`).
- Arrange partitions in **tree hierarchy**:
 - Depth of tree: $n+1$.
 - Width of tree \leq product of the size of active domains of columns c_i .
- Inner nodes: OS directories, leaf nodes: data (Parquet or CSV).
- Filter query: traverse relevant subtrees/read relevant data files only (left: \equiv `payment_type = 3`).
- Particularly useful when I/O is slow.

Tabular Database Systems

⑥

The Structured Query Language (SQL)

March 17, 2026

Torsten Grust
Universität Tübingen, Germany

1 | The Origins of SQL



Don Chamberlin



Ray Boyce (†1974)

Don Chamberlin and **Ray Boyce**, co-inventors of SQL, back in 1972/73 both members of IBM's research project *System R*.

The Origins of SQL

- Development of the language started in 1972, first as **SQUARE**, from 1973 on as **SEQUEL** (*Structured English Query Language*). In 1977, SEQUEL became **SQL** because of a trademark dispute. (Thus, both “S-Q-L” /,ɛskjuːˈɛl/ and “*sequel*” /ˈsiːkwəl/ are okay pronunciations.)
- First commercial implementations in the late 1970s/early 1980s. By 1986, the ANSI/ISO standardization process begins.
- Since then, SQL has been in active development and remains the **“Intergalactic Dataspeak”**¹ to this day.

Current SQL standard (as of April 2025): SQL:2023.

¹ Due to Mike Stonebraker, inventor of Ingres (1972, precursor of Postgres, PostgreSQL)

2 | SQL: Row Variables


Like most regular programming languages, SQL features a construct to **bind variables to values**.

- `FROM` is the *only* SQL clause that can introduce variables:

```

⋮
FROM t AS v -- bind variable v to the rows in table t
⋮

```


- If $|t| = m$, `v` iterates over all m rows of `t` in some order. `v` is thus known as **row variable** (sometimes: table alias .
- If table `t` has schema `t(c1 τ1, ..., cn τn)`, then `v` has type `row(c1 τ1, ..., cn τn)`.²
- Can access column `ci` of `v` using **dot notation**: `v.ci` (`:: τi`).

² In DuckDB, `row(c1 τ1, ..., cn τn)` is a synonym for type `struct(c1 τ1, ..., cn τn)`: think of a table row like a record/struct with n fields.

FROM generates Cross Products

- If a SQL query needs to read data from **multiple tables** t_1, \dots, t_m , list all tables in the **FROM** clause:

```
⋮  
FROM  $t_1$  AS  $v_1$ , ...,  $t_m$  AS  $v_m$  -- row variable names  $v_i$  unique  
⋮
```

- **Cross product** semantics: This generates **all** $|t_1| \times \dots \times |t_m|$ **combinations** of bindings for row variables v_i in some order.
- Row variable names v_i are unique, tables t_j may repeat (**FROM** t **AS** v_1 , t **AS** v_2 allows to combine each row of t with its peers).
-  Yes, **FROM** clauses may generate *MANY* bindings. Typical queries will use predicates on the v_i to only focus on the combinations that make sense.

3 | SQL: Joining Tables

Cross products between tables are (too?) general: it is common to draw rows from tables based on *conditions* that identify (un)wanted row combinations. SQL supports several such **table joins**.

- **Inner Join:**³

```












⋮
FROM  $t_1$  AS  $v_1$  [INNER] JOIN  $t_2$  AS  $v_2$  ON ( $p$ )
⋮

```



- Draw all combinations of rows from tables t_1 , t_2 .
- Only keep those bindings of v_1 , v_2 that satisfy **predicate p** (v_1 , v_2 typically do occur free in p).

³ Inner join is often referred to as simply *join* or *θ -join* (initially, DB literature used the fancy greek theta θ instead of p to denote the join predicate).

Inner Join Follows Foreign Keys

vehicles						peeps			
vid	vehicle	kind	seats	wheels?	pid	pid	pic	name	born
v ₁		car	5	true	p ₄	p ₁		Cleo	2013
v ₂		SUV	3	true	p ₄	p ₂		Bert	1968
v ₃		bus	42	true	□	p ₃		Drew	□
v ₄		bus	7	true	□	p ₄		Alex	2002
v ₅		bike	1	true	p ₂				
v ₆		tank	□	false	p ₃				
v ₇		cabrio	2	true	p ₄				

```
SELECT v.vehicle, p.name AS driver, p.pic
FROM vehicles AS v INNER JOIN peeps AS p ON (v.pid = p.pid)
                        "equi-join" ↗
```

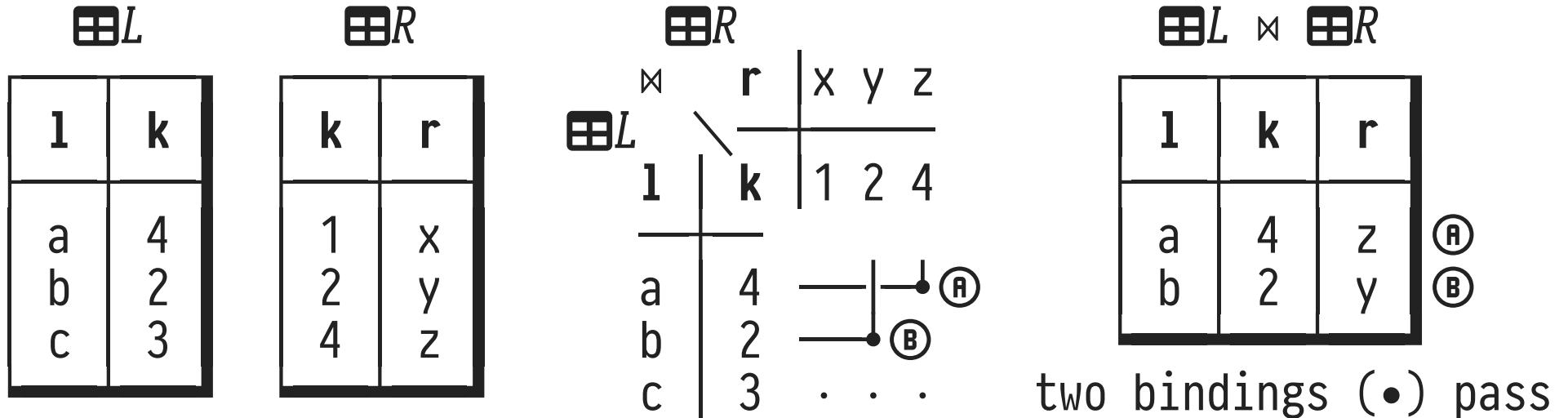
- Some rows in table `peeps` find multiple *join partners* in `vehicles` (like p₄ ) , some find none (p₁ ).
- General: A join between `t1` and `t2` may yield $0 \dots |t_1| \times |t_2|$ rows.

More Join Operations ①

Inner joins (in particular: equi-joins) are frequent. Daily SQL practice calls for a whole **family of join variants**, however.

To introduce the join family, let us use a sketch of their semantics (see below). We assume join predicate $p \equiv L.k = R.k$.

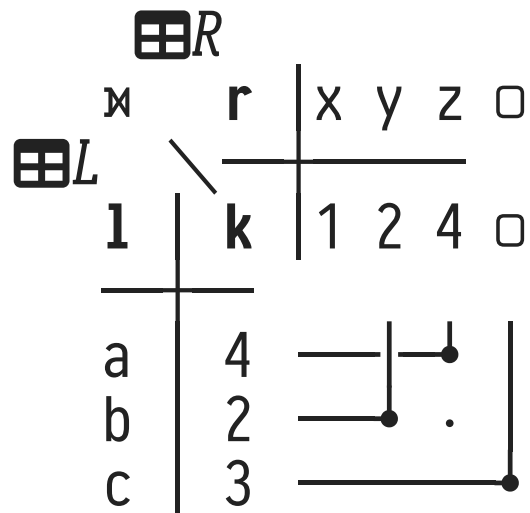
- **INNER JOIN** (\bowtie , `FROM L INNER JOIN R ON (L.k = R.k)`):



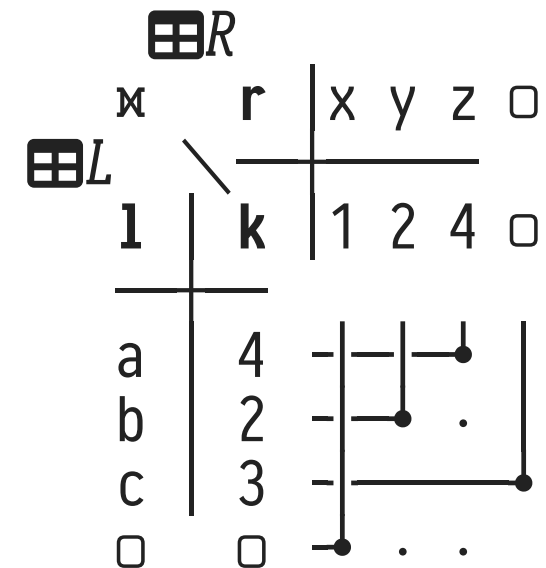
More Join Operations ②: Outer Joins

- LEFT/RIGHT/FULL OUTER JOIN (\bowtie , \bowtie , \bowtie)

LEFT OUTER JOIN



FULL OUTER JOIN



- Left outer join: *All rows* of the left table $\bowtie L$ are kept. Rows from $\bowtie L$ that find no join partner (like $(c,3)$) are paired with an artificial all-NULL (\square) row.

More Join Operations ③: Semi/Anti/Cross Joins

• SEMI/ANTI/CROSS JOIN (\ltimes , $\bar{\ltimes}$, \times)

SEMI JOIN


		\ltimes	$\bar{\ltimes}$	\times
		\ltimes	$\bar{\ltimes}$	\times
\ltimes	$\bar{\ltimes}$	\times	$\bar{\ltimes}$	\times
\ltimes	$\bar{\ltimes}$	\times	$\bar{\ltimes}$	\times
1	k	1	2	2
a	1	•	.	.
b	2	—	—	•
c	3	.	.	.

ANTI JOIN

		\ltimes	$\bar{\ltimes}$	\times
		\ltimes	$\bar{\ltimes}$	\times
\ltimes	$\bar{\ltimes}$	\times	$\bar{\ltimes}$	\times
\ltimes	$\bar{\ltimes}$	\times	$\bar{\ltimes}$	\times
1	k	1	2	4 /
a	4	.	.	o .
b	2	.	o	. .
c	3	—	—	•

CROSS JOIN

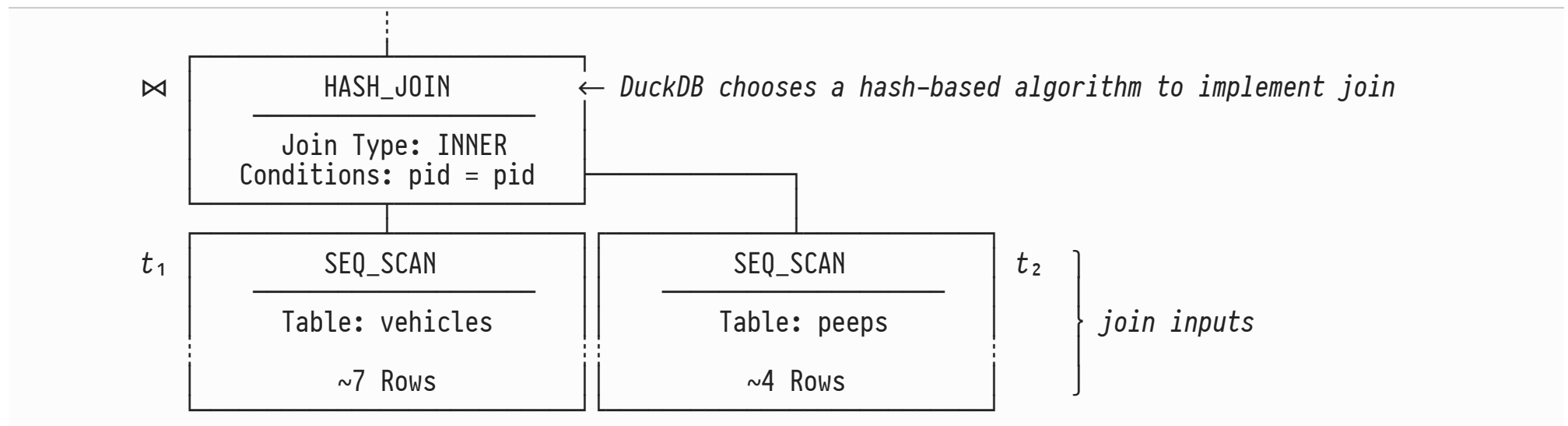
		\ltimes	$\bar{\ltimes}$	\times
		\ltimes	$\bar{\ltimes}$	\times
\ltimes	$\bar{\ltimes}$	\times	$\bar{\ltimes}$	\times
\ltimes	$\bar{\ltimes}$	\times	$\bar{\ltimes}$	\times
1	k	1	2	4
a	4	•	•	•
b	2	•	•	•
c	3	•	•	•

- Semi join: Keep rows from $\ltimes L$ with *at least one join partner*.
- Anti join: Keep rows from $\ltimes L$ that have *no join partner*.
- Semi/anti join return a subset of rows from $\ltimes L$:  #032
join predicate p aside, row variables bound to $\ltimes R$ are inaccessible in the SQL query (thus above: \rightarrow instead of \rightarrow).

Joins in Query Plans

Joins—in textbooks often: \bowtie (or $\bowtie \bowtie \bowtie \bowtie \bar{\bowtie} \bowtie \bowtie \dots$)—operate on two tables. A join between n tables requires $n-1$ binary \bowtie operations).

- **FROM vehicles NATURAL JOIN peeps** in **EXPLAIN** output (excerpt):



- Inner join is *commutative* ($t_1 \bowtie t_2 \equiv t_2 \bowtie t_1$) and *associative* ($(t_1 \bowtie t_2) \bowtie t_3 \equiv t_1 \bowtie (t_2 \bowtie t_3)$).

Exploited by the DBMS's **query optimizer** to rearrange joins and inputs to find an efficient query plan alternative.

4 | SQL: Bag Algebra

A table instance is a **bag** of rows (no order, duplicate rows may occur). SQL features the binary operations of the bag algebra:

[Likewise: **INTERSECT ALL** \cap_+ , **EXCEPT ALL** \setminus_+]

bag of rows

$\underbrace{query_1 \text{ UNION ALL } query_2}_{\text{two bags of rows}} \quad \text{-- } query_1 \uplus query_2$

two bags of rows

- Rows returned by *query*₁, *query*₂ (and the result) conform:⁴
 - Number of columns are equal.
 - Types in the same column position match (or are castable).
 - *query*₁ determines the column names of the result.

 #034

⁴ Alternatively, DuckDB can match columns in a bag union operation by name (not position): **UNION ALL BY NAME**. This is non-standard. See the [DuckDB documentation on bag \(and set\) operations](#) .

SQL's Bag Algebra Respects Row Multiplicities

Operations of the bag algebra respect *multiplicities*. So does SQL:

$$\begin{array}{l}
 \{\text{A A A B C}\} \cup \{\text{A B B D}\} = \{\text{A A A A B B B C D}\} \quad (\text{add}) \\
 \{\text{A A A B C}\} \cap_+ \{\text{A B B D}\} = \{\text{A B}\} \quad (\text{minimum}) \\
 \{\text{A A A B C}\} \setminus_+ \{\text{A B B D}\} = \{\text{A A C}\} \quad (\text{subtract, } \geq 0)
 \end{array}$$

A B C D denote rows, $\{\dots\}$ denotes a bag (table) of rows.

- To ignore multiplicities, omit the **ALL** modifier. Bag operations then discard row duplicates ($\delta(\{\text{A A A B C}\}) \stackrel{\text{def}}{=} \{\text{A B C}\}$):

$$\text{query}_1 \text{ OP query}_2 \equiv \delta(\delta(\text{query}_1) \text{ OP ALL } \delta(\text{query}_2))$$

- Duplicate elimination** (δ) can be generally useful in queries. SQL thus supports the **DISTINCT** modifier in the **SELECT** clause:

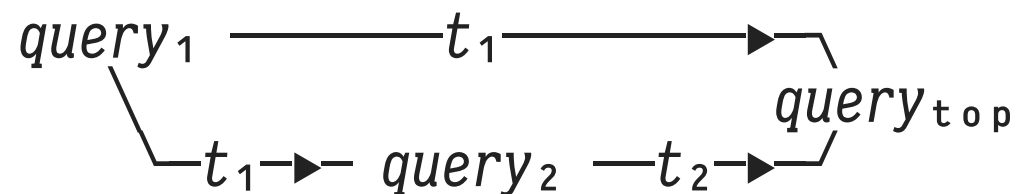
$$\begin{array}{l}
 \text{SELECT DISTINCT } \dots \\
 \text{FROM } \dots \\
 \vdots
 \end{array}
 \equiv
 \delta \left(\begin{array}{l}
 \text{SELECT } \dots \\
 \text{FROM } \dots \\
 \vdots
 \end{array} \right)$$

5 : SQL: Assembling Queries From Pieces (Common Table Expressions)

If a complex SQL query is best understood in terms of *intermediate result tables*, name these intermediates and refer to them later:


WITH	Common Table Expression (CTE)
t_1 AS ($query_1$),	• locally define name t_i as result of $query_i$
$t_2(c_1, \dots, c_n)$ AS ($query_2$)	• may name result columns of intermediates
$query_{top}$	• top-level query (= result of CTE)


- In such a CTE, the dependencies between queries and visibility (or: scopes) of the t_i form a DAG. No cycles:



- Decision: *materialize* t_i or *inline* $query_i$ at its use site(s)?

6 : A *Star Wars* Database

kaggle  is an extensive online repository of open data sources, many of those in CSV and/or Parquet formats.

- The “*Star Wars Dataverse*” collects bits of information on Disney's *Star Wars* franchise in 15 Parquet files.
- We created a native  database file '*035-starwars.db*' from these Parquet sources and cleaned the data:
 - Normalized in-universe dates⁵, map to type `int` (`'3 BBY'` → `-3`, `'2 ABY'` → `2`).
 - Turned comma-separated `text` values into `text[]` arrays (`'Luke, Leia'` → `['Luke', 'Leia']`).
 - Normalized film titles across tables. (`'Episode IV: A New Hope'` → `'A New Hope'`)

Load the *Star Wars* database to exercise your SQL skills.

 #035

⁵ BBY/ABY: Years *Before/After* the *Battle of Yavin* in which the Rebel Alliance  destroyed the Death Star .

DuckDB UI: Browsing a Database

Load the *Star Wars* database and use the browser-based **DuckDB UI**⁶ to get familiar with the schemata/instances of the 15 tables:

1. Install and load the DuckDB UI extension (first use only):

```
D INSTALL motherduck;  
D LOAD motherduck;
```

2. Load *Star Wars* database, invoke the UI (starts browser):

```
D ATTACH '035-starwars.db' AS starwars;  
D CALL start_ui();
```

```
UI started at http://localhost:4213/
```

⁶ Once the UI extension is installed, may start DuckDB UI from the shell via command `duckdb -ui <database>.db`.

SQL Training (Star Wars Database)

Use SQL to formulate these queries against the *Star Wars* database. If needed, use CTEs (**WITH ...**) to cope with query complexity.

1. Which characters have got “*a bad feeling about this*” in what film(s)? [★]
2. Which films (title, year) make up the prequels (released after *Return of the Jedi*, before *The Force Awakens*)? [★★]
3. Which films (title) were created by the congenial duo of director George Lucas and composer John Williams? [★★]
4. Which droids (name) appear in *any* film directed by George Lucas?
5. Which droids (name) appear in *all* films directed by George Lucas? [★★]
[★★★★]

Tabular Database Systems

⑦

More SQL (Subqueries + Embedded SQL)

March 17, 2026

Torsten Grust
Universität Tübingen, Germany

1 | Compositionality (in Programming Languages)

“The meaning of a complex program is determined by the meanings of its constituent programs.”

—Principle of **Compositionality**

The following two (Python) programs are equivalent:

1. Uses literal **21**:

```
print(2 * 21)
```

2. Computes value **21**:¹

```
def twenty_one():  
    return (12 + 12*12 + 20) // (2*4) - 1  
print(2 * twenty_one())
```

¹ Based on a variation of a popular limerick: “A dozen, a gross, and a score / divided by two times the four / decreased just by one / it gives twenty-one / (which is three times seven, no more).”

2 : SQL: Subqueries (Queries Inside Queries)

SQL queries may contain nested **subqueries** enclosed in (...).

1. **Scalar subquery:** Where a query accepts a scalar x , x may be replaced by a subquery (q) that returns a single-cell table:

q evaluates to

x


 $\equiv x$

 #036

NB. A scalar subquery q needs to return...

- ... a **single-column** table (column name irrelevant). Otherwise: error at query *compile time*. 🗑️
- ... **at most one row**.
0 rows \equiv **NULL**, ≥ 2 rows: error at query *run time*. 🗑️ 🗑️

Subqueries Can Relate to Outer Queries: Correlation

- Row variables bound inside a subquery do not “escape”: their scope is local to the subquery (see var `v1` in  #036).
- **But:** Subqueries may relate to row variables bound in the enclosing/outer query:

```

      outer query
SELECT v.*, ( [subquery] )
              [ ... v ... ]
FROM   vehicles AS v;

```

The subquery can access the current row bound to `v` ✓

- The subquery in `[]` *cannot* be evaluated in isolation: depends on outer query `□` to provide a binding for row variable `v`.
- DB jargon: “The subquery is **correlated** (since it uses `v`).”²

² PL jargon: “Variable `v` occurs free in the subquery (but bound in the outer query).”

Correlated Subqueries \equiv Nested Loops? 🤖

NB. Due to correlation, the subquery q in `[]` below acts like a function of type `int \rightarrow text` ($q(pid)$ maps vehicle to driver `name`):

```
-- Query Q: Pair vehicles with their driver (if any)
```

```
SELECT v.*, (
  
    SELECT p.name
    FROM   peeps AS p
    WHERE  p.pid = v.pid
  
) AS driver
FROM   vehicles AS v;
```

- A “*nested loops*” evaluation strategy for Q :

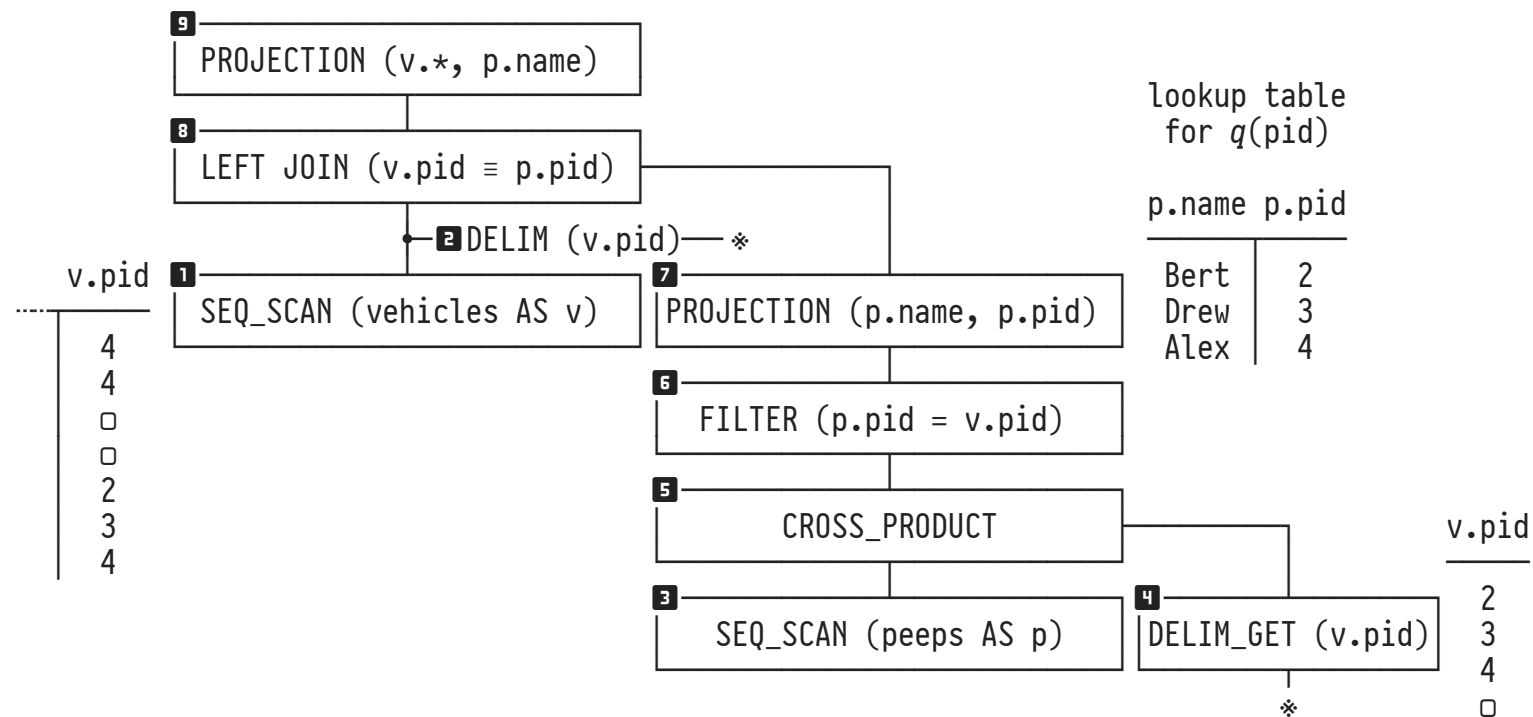
```
for each v  $\in$  vehicles
{
  driver  $\leftarrow$   $q(v.pid)$ 
  result  $\leftarrow$  result  $\cup$  (v.*, driver)
}
evaluated
|vehicles| times
```

- If `vehicles.pid` is not unique, this will evaluate subquery $q(pid)$ repeatedly with identical arguments. 🗨️

DuckDB's Query Optimizer Removes Correlation

Query optimization **decorrelates** the subquery (plan for Q):³

1. **2**, **4**: Compute the set of *distinct arguments* for subquery q .
2. **3**–**7**: Build a lookup table for “function” $q(pid)$.
3. **8**: For each vehicle v , perform a lookup to evaluate $q(v.pid)$.



#038

³ Plan has been simplified. Use `PRAGMA explain_output = 'all';` to see the above Unoptimized Logical Plan. In these plans, the `LEFT (OUTER) JOIN/DELIM` pair is represented as `DELIM_JOIN`.

Scalar vs. Table-Valued Subqueries

SQL interprets subqueries (q) based on their *usage context*:

1. **Scalar**: q returns single-cell table \square that holds one scalar.
2. **Table-valued** (e.g., in **FROM** clause): q returns any table \boxplus :

outer query

```

⋮
FROM t1 AS v1, ( [subquery q]
                   [⋮ v1 ⋮]
                   [-----] ) AS v2, ...
⋮
  
```

 #039

Subquery q in
table-valued
context

- **Correlation**: row variable v_1 may occur free⁴ in subquery q . In this case, \boxplus acts like a table-valued function $q(v_1)$. (DBMS will decorrelate to avoid $|t_1|$ evaluations of q).
- If q is uncorrelated: \boxplus acts like a (computed) table.

⁴ Some SQL implementations require the keyword **LATERAL** (“sideways”) to allow q to refer to v_1 (and thus depend on the evaluation of t_1): **FROM** t_1 **AS** v_1 , **LATERAL** (q) **AS** v_2 , ... (DuckDB infers whether **LATERAL** is required.)

3 | SQL: Existential and Universal Quantification

SQL uses table-valued subqueries (q) to compactly formulate **existentially or universally quantified comparisons**:⁵

EXISTS (q) does q return ≥ 1 rows (is q non-empty)?
NOT EXISTS (q) does q return no row at all (is q empty)?

$expr =$ **ANY** (q) does $expr$ equal any value in q ?
 $expr =$ **ALL** (q) does $expr$ equal all values in q ?

also: $<>$ $<$ $<=$ $>=$ $>$

q evaluates to a
single-column table:



⁵ The SQL keywords **ANY** and **SOME** are synonyms. Syntactic sugar: $expr =$ **ANY** (q) is equivalent to $expr$ **IN** (q) (think of \in or “is element **in**”).

4 | Embedding SQL in Python Programs





- The DuckDB CLI `>-` enables interactive experimentation and the execution of ad-hoc/one-short querying. Definitely valuable.
- **Database-supported applications** embed SQL statements directly in the program source instead:
 - Programs can connect to/disconnect from selected databases.
 - Program flow controls which/how often SQL queries execute.
 - Queries may be constructed/parameterized on the fly.
 - Query results may be consumed by the program:
 - Map SQL data types to programming language's type system.
 - Receive *all* rows at once? *Iterate* over result row-by-row?

Q: Which parts of the app logic are performed by the DBMS?
Which parts are implemented by program code? `~_(\ツ)_/`

Here  ↔ : Use DuckDB's API to **embed SQL queries into Python.**

Embedding SQL in Python: General Setup

Application code mixes SQL query strings  with program code :

```
import duckdb # requires: pip install duckdb
:
# 1 connect to database (in  or on )
with duckdb.connect(database) as con:
    # 2 construct SQL query, submit to DuckDB
    rel = con.sql("""
        [-----
         SQL query
        -----]
    """)
    # 3 DuckDB executes query, retrieve all result rows
    result = rel.fetchall()
    # 4 iterate over list of rows
    for row in result:
        [-----
         code
        -----]
        # Python code operates on a result row
        # represented as a tuple (...,...,...)
        # DuckDB not involved
```

 #042

DuckDB's Python DB API (Overview⁶)

DuckDB Python API Call	Python Result
1 <code>con = duckdb.connect(":memory:")</code> <code>con = duckdb.connect(<i>database file</i>)</code>	DuckDB connection object
2 <code>rel = con.sql(<i>sql</i>)</code> <code>rel = con.sql(<i>sql</i>, params = [...])</code>	effect on DB or DuckDB relation object ⁷ <code>sql</code> may contain parameters <code>\$1</code> , <code>\$2</code> , ...
3 <code>rel.fetchall()</code> <code>rel.fetchmany(<i>n</i>)</code> <code>rel.fetchone()</code> <code>rel.columns</code> <code>rel.types</code> <code>rel.show()</code>	list of all result tuples in table <code>rel</code> list of next <code>n</code> result tuples next result tuple or <code>None</code> list of column names for table <code>rel</code> list of column types <code>None</code> + printed table output (📄)

⁶ [DuckDB's Python DB API \(documentation\)](#) 🖱️

⁷ If `sql` is a **SELECT** query, returns a DuckDB relation object `rel`. Otherwise, applies the effect of the DDL or DML statement to the database represented by DuckDB connection object `con`.

How SQL Data Maps to Python Values/Objects

The type systems of SQL and Python—indeed most PLs—differ.⁸ SQL data (tables, rows, cell values) needs to be **mapped** to Python values and objects:

	SQL		Python	
	table	≡	list	[
	row	≡	tuple	(A ₁ , B ₁ , C ₁),
	cell a _i	≡	val/obj A ₁	(A ₂ , B ₂ , C ₂),
				(A ₃ , B ₃ , C ₃)
]

Data and Type Mapping in DuckDB's Python DB API

- DuckDB implements a best-effort mapping from cell values a_i (of types `int`, `text`, ...) to Python's value and objects A_i . 📄 #043

⁸ This tension—not only regarding types—between SQL as a query language and programming languages (PLs) is known as the **impedance mismatch**.

Constructing SQL Queries at Program Run Time

Embedded SQL queries are regular strings: **programs can construct queries at run time** by interpolation 👉 or concatenation 🗑️.

1. **Interpolation** (Python values replace parameterized SQL values \$**1**, \$**2**, ... in a template query):

```
"SELECT $1 || p.name FROM peeps AS p AS WHERE p.born < $2"
```

🐍 [**1**:"Driver: ", **2**:2000]

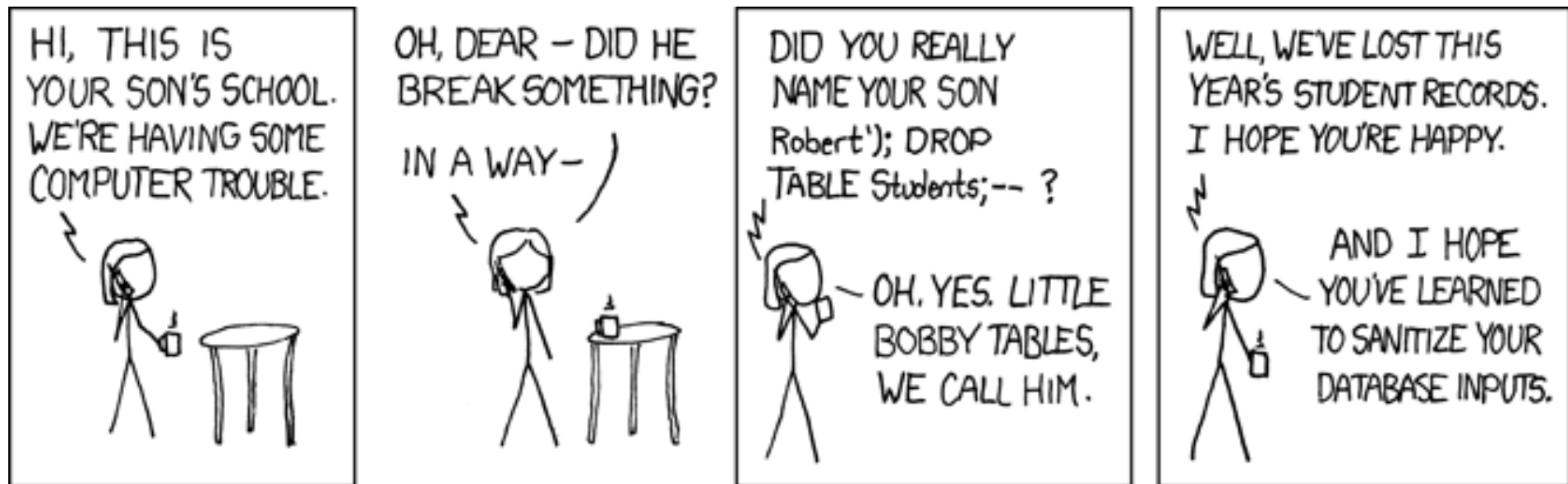
2. String concatenation (⚠️ Risk of SQL injection):

```
"FROM peeps AS p WHERE p.name = ' ' + N + ' ' AND p.born < 2000"
```

🐍 string concatenation

If an attacker controls the value of Python string variable *N*, the DBMS may be tricked into executing arbitrary queries.

Parameterized Queries Protect Against Little Bobby's Mom



"Exploits of a Mom" ↗, © xkcd

 #044 + #045

Move Your Computation Close to the Data!

A rule of thumb 👍: *If you can, express data-related computation using SQL.⁹ Do not demote the DBMS to a dumb table storage.*

- Filter/aggregate tables to **reduce result set sizes** before you transport data from the DBMS to the program.
- Common anti-pattern (the “***n+1 query problem***”):

```
outer = con.sql("SELECT ...")      # yields n rows
for row in outer.fetchall():
    inner = con.sql("SELECT ... $1 ...", params = [... row ...])
    :
```

 #046

- The above issues *n+1 separate queries* all of which need to be interpolated, parsed, optimized, executed, and fetched. 👎
- Reformulate using a join or a correlated subquery (which the query optimizer will decorrelate). Issues a *single query*. 👍

⁹ Indeed, with SQL:1999 and the introduction of *recursive common table expressions*, the query language has become Turing-complete. We explore the consequences of this jump in expressiveness in the course *Advanced SQL*.

Tabular Database Systems

⑧

SQL: Grouping + Aggregation and Functional Dependencies

April 7, 2026

Torsten Grust
Universität Tübingen, Germany

1 | SQL: Grouping (Clause **GROUP BY**)

- Up to now, SQL queries operate *row-by-row*:
 - FROM** t binds *individual rows* of t to row variables,
 - WHERE** p drops row var bindings that do not satisfy p ,
 - SELECT** e, \dots, e evaluates e, \dots, e for each row var binding.

Clause **GROUP BY** e_1, \dots, e_n identifies **subtables (or: groups)** of rows that yield the same value for expressions (criteria) e_1, \dots, e_n :

🚗 vehicles					groups
vehicle	kind	seats	wheels?	driver	
🚲	bike	1	true	p_2	true
🚗	cabrio	2	true	p_4	
🚙	SUV	3	true	p_4	
🚗	car	5	true	p_4	false
🚌	bus	42	true	□	
🚌	bus	7	true	□	
🚛	tank	□	false	p_3	□

↓

```

SELECT ...
FROM   vehicles AS v
GROUP BY v.seats < 5
          grouping criterion
  
```

} \forall rows v in this group:
 $v.seats < 5 \equiv \text{false}$

SQL: Grouping Requires Aggregation

💡 After grouping, use **aggregation** to represent group contents:

❶ 🚗 vehicles (grouped)

vehicle	...	seats	...	driver
🏍️		1		p_2
🚗	...	2	...	p_4
🚗		3		p_4
<hr/>				
🚗		5		p_4
🚌	...	42	...	□
🚗		7		□
<hr/>				
🚗	...	□	...	p_3

❷ 🚗 vehicles (grouped + aggregated)

seats<5	count(*)	max(seats)	list(driver)
true	3	3	$[p_2, p_4, p_4]$
false	3	42	$[p_4, \square, \square]$
□	1	□	$[p_3]$

- Tables are in 1NF. Cells cannot contain groups/subtables. Thus:
 - SQL never provides access to the grouped table **❶** itself.
 - **Aggregation after grouping **❷** is mandatory.**

GROUP BY Reduces Granularity (Row-by-Row to Group-by-Group)

- The **FROM** **1** and **WHERE** **2** clauses operate row-by-row.
- After **GROUP BY** **3**, the **SELECT** **4** clause operates group-by-group and yields **one row per group**:

```

SELECT  expr1, ..., exprn,
        agg, ..., agg
FROM    t
WHERE   p
GROUP BY expr1, ..., exprn

```

```

3 } g
4 }
1 n (= |t|)
2 m (≤ n)
3 g (≤ m)
# rows

```

 #048

Non-aggregate *expr*_{*i*} in **SELECT** is OK only if (the DBMS can statically deduce that) *expr*_{*i*} is *constant within each group*.

- **Grouping Quiz:** What can you infer if you observe the following?
 - Ⓐ $g = 1$. Ⓑ $g = m$. Ⓒ $g = 0$.
 - Ⓓ Adding a new *expr*_{*j*} to the **GROUP BY** clause does not change *g*.

Controlling Granularity: Derived Grouping Criteria

```
SELECT t.c,agg,...,agg
FROM   t
GROUP BY t.c -- criterion c leads to too many groups...
```

 #049

~_(\ツ)_/~

1. **Division.** Group numeric c into buckets of **equal width** $1/N$:
 - If c is integral (int): **GROUP BY** $c // N$
 - Otherwise: **GROUP BY** $\text{round}(c / N)$
2. **Modulus.** Put integral c into one of N buckets (**order lost**):
 - **GROUP BY** $c \% N$
 - If $N = 2^n$: **GROUP BY** $c \& (1 \ll n) - 1$ (based on n low bits)
3. **Grading.**¹ Place ordered c into buckets with **given borders** b_i :






¹ *Grading* is derived from APL's \uparrow (“grade up”), `list_grade_up()` in DuckDB: `list_grade_up([30,10,20]) = [2,3,1]` are the list indices that would rearrange `[30,10,20]` in ascending order.

A Grouping Playground: The New York City Taxi Database



The *Yellow Cab* database² published by the NYC Taxi & Limousine Commission (TLC) provides a fun playground for **GROUP BY** queries.

A DuckDB database file holds the dataset for year 2024:

 #051

 rides	main table, one row per taxi trip
 zones	division of New York City into 265 zones (“boroughs”), rides start/stop here
 central_park_weather	temperature/wind/precipitation data for Central Park, one row per day

Find plain + embedded SQL queries over this data in

 #050  #052











² Published monthly on the [TLC Trip Record Data page](#) , (formerly in CSV, since 2022 in Parquet format).

2 : Group Preservation

```
-- How many peeps can each driver give a lift?
SELECT p.pid AS driver, max(v.seats) - 1 AS "can lift",
       max(v.seats) - 1 AS "can lift"
FROM   vehicles AS v NATURAL JOIN peeps AS p
GROUP BY p.pid, max(v.seats) - 1 AS "can lift";
```

 #053

vehicles ⋈ peeps

vehicle	kind	seats	wheels?	pid	pic	name	born
	bike	1	true	2		Bert	1968
	tank	0	false	3		Drew	0
	car	5	true	4		Alex	2002
	SUV	3	true	4		Alex	2002
	cabrio	2	true	4		Alex	2002


- Adding `p.name` (`pic`, `born`) to `GROUP BY` ~~max(v.seats) - 1 AS "can lift"~~ preserves groups: rows that agree on `pid`, also agree on `name`: `pid` → `name`.

Functional Dependencies (FDs)

Given table $t(c_1, \dots, x, \dots, y, \dots, c_n)$, the **functional dependency (FD)** $x \rightarrow y$ (“column x determines column y ”) holds in t iff

$$\forall \text{ rows } r_1, r_2 \in t: r_1.x = r_2.x \Rightarrow r_1.y = r_2.y.$$

• Notes on FDs:

- FD $x \rightarrow y$ indicates that t embeds a lookup table $t^f(x, y)$ defining a function $f(x) = y$.³ Embedded table t^f has key x .
- Generalization: $x_1 \ x_2 \ \dots \ x_n \rightarrow y$ defines an n -ary function.
- If $x \rightarrow y$ and $y \rightarrow z$, then $x \rightarrow z$ (transitivity, like $g \circ f$).
- Shorthand notation: $x \rightarrow y$ and $x \rightarrow z \equiv x \rightarrow y \ z$.
- A key  k for table $t(k, c_1, \dots, c_n)$ is like a special, powerful FD: $k \rightarrow c_1 \ \dots \ c_n$ (key k determines all columns).
- Trivial FD (holds for any column x in any table t): $x \rightarrow x$.
- $x \rightarrow y$ does *not* generally imply $y \rightarrow x$.

 #053

³ Example: For the FD $pid \rightarrow pic$ in `vehicles` \bowtie `peeps`, we have $f(2) = \text{☺}$, $f(3) = \text{☹}$, $f(4) = \text{☹}$ (3x).

SQL: HAVING (and a Test for Functional Dependencies)

- SQL clause **HAVING** is evaluated *after* groups have been formed. Predicate q thus may use aggregates to **filter entire groups**:

SELECT	$expr, \dots, expr,$	
	agg, \dots, agg	
FROM	t	
WHERE	p	
GROUP BY	$expr, \dots, expr$	
HAVING	q	

- A **HAVING**-based test: Does FD $x \rightarrow y$ hold in table t ? 📄 #054

SELECT DISTINCT	$'x \not\rightarrow y'$	AS	"FD violated?"	} empty result indicates that $x \rightarrow y$ does hold
FROM	t			
GROUP BY	x			
HAVING	$\text{count}(\text{DISTINCT } y) > 1$			

3 : Apps Define FDs ...

🗪 trips



trip	weekday	driver	vehicle	from	to	dist	via	hop#
t_1	Mon	4	🚚	Alton	Corby	150	Alton	1
t_1	Mon	4	🚚	Alton	Corby	150	Luton	2
t_1	Mon	4	🚚	Alton	Corby	150	Corby	3
t_2	Tue	2	🚲	Derby	Eaton	17	Derby	1
t_2	Tue	2	🚲	Derby	Eaton	17	Eaton	2
t_3	Wed	2	🚲	Derby	Crich	23	Derby	1
t_3	Wed	2	🚲	Derby	Crich	23	Crich	2
t_4	Thu	4	🚚	Alton	Corby	150	Alton	1
t_4	Thu	4	🚚	Alton	Corby	150	Luton	2
t_4	Thu	4	🚚	Alton	Corby	150	Corby	3

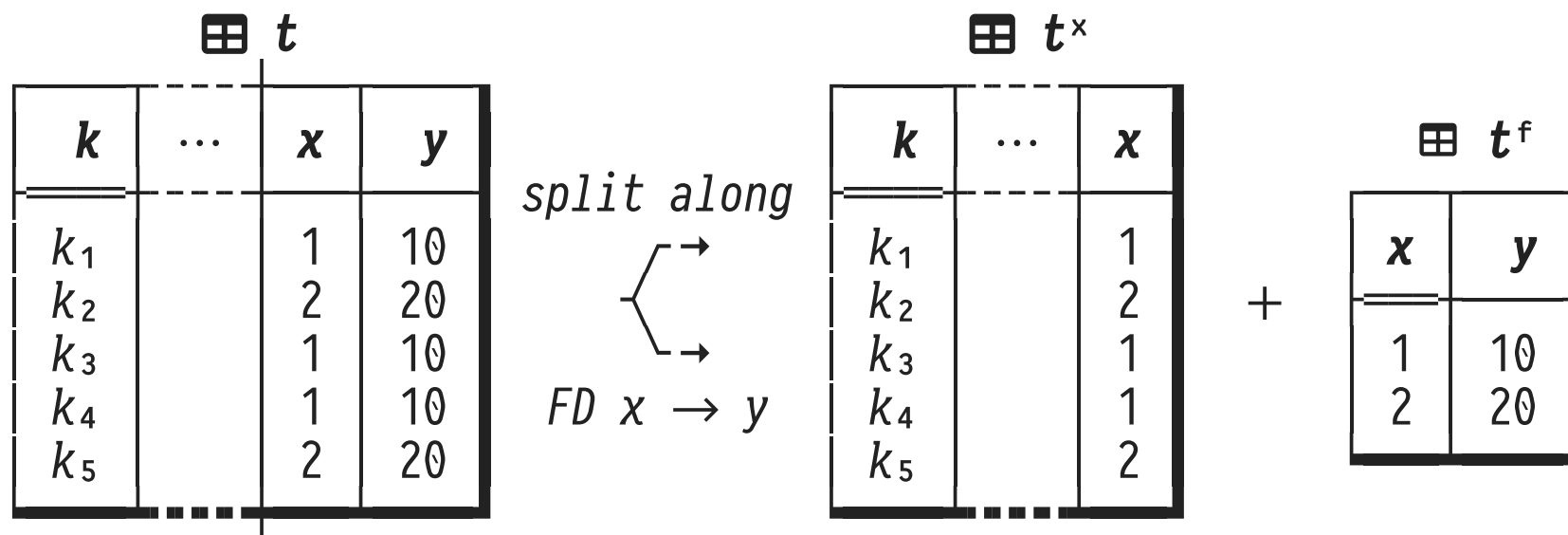
Schedule of multi-hop road trips taken (with drivers and vehicles in use).



- Here are some FDs that hold in table `trips`. Verbalize what they mean in a DB application that manages road trips:


- Ⓐ `vehicle` \rightarrow `driver` Ⓑ `trip` \rightarrow `weekday` `from` `to` Ⓒ `weekday` \rightarrow `driver`
- Ⓓ `from` `to` \rightarrow `dist` `driver` Ⓔ `from` `to` `hop#` \rightarrow `via` `vehicle`

... FDs Suggest Good Table Designs

1. The only FDs actually enforced by DBMSs are primary keys : recall the **PRIMARY KEY** constraint.
2. Non- FDs (embedded lookup tables) indicate data redundancy.



- After the **split** $\mathbb{t} \leftarrow \mathbb{t}^x + \mathbb{t}^f$ along $FD\ x \rightarrow y$:⁴
 - x is key in t^f : the DBMS can now enforce the $FD\ x \rightarrow y$. 
 - **No redundancy** in t^f . Function $f(x) = y$ is now explicit. 

⁴ Advise: When splitting along $FD\ x \rightarrow y$, move *all* columns functionally determined by x into t^f (including transitively determined columns). If t^f contains non- FDs, split t^f again.

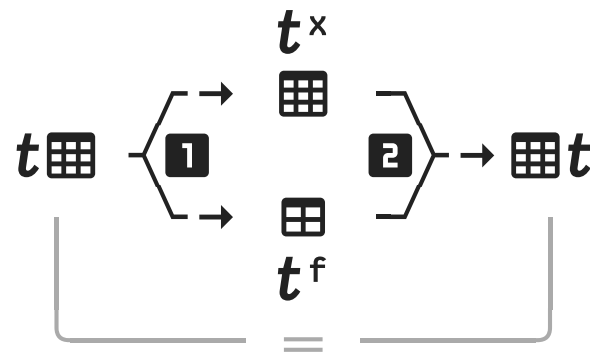
Splits Can Be Undone (Virtually)

FD splits improve constraint checking and reduce redundancy. But users or apps may require/expect the original non-split table...

- An inner join \bowtie can **faithfully reconstruct** the original $\boxplus t$.
Table state is preserved (rows are neither lost nor added):⁵

FD split **1** **2** Inner join \bowtie

 #056



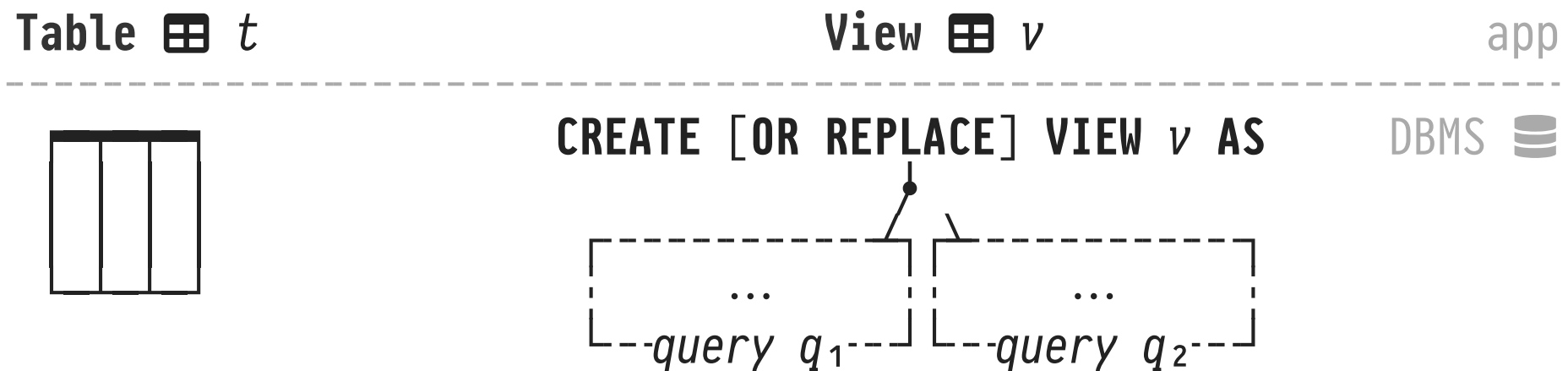
- A **table view** can render this reconstruction of $\boxplus t$ transparent.

⁵ FD splits are **lossless splits**.  But: arbitrary vertical table splits may lose information ($\boxplus \neq \boxplus + \boxplus + \boxplus$).

4 : Views (Logical Data Independence)

A SQL **view** v represents a *virtual* table whose schema and state is **defined by a query** q (instead of extensionally by a bag of rows).

- Views provide **logical data independence**. The actual data sources read by q remain hidden (q may change ↯ ↔ ↰):



- Whenever v is referenced by a query, the DBMS re-evaluates q .
- Views are read-only and cannot be updated.⁶

⁶ In general, it is ambiguous how an “update” of v should affect the underlying source tables read by q .

The End.

Since you've got this far... A companion course on the design and implementation of selected DuckDB internals is available at

<https://github.com/DBatUTuebingen/DiDi>.⁷

⁷ DiDi: Design and Implementation of DuckDB Internals  