

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)_/~

- Division.** Group numeric c into buckets of **equal width $1/N$** :
 - If c is integral (int): **GROUP BY $c // N$**
 - Otherwise: **GROUP BY $\text{round}(c / N)$**
- 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)
- 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?"	
FROM	t			
GROUP BY	x			
HAVING	$\text{count}(\text{DISTINCT } y)$		> 1	

empty result indicates that $x \rightarrow y$ does hold

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

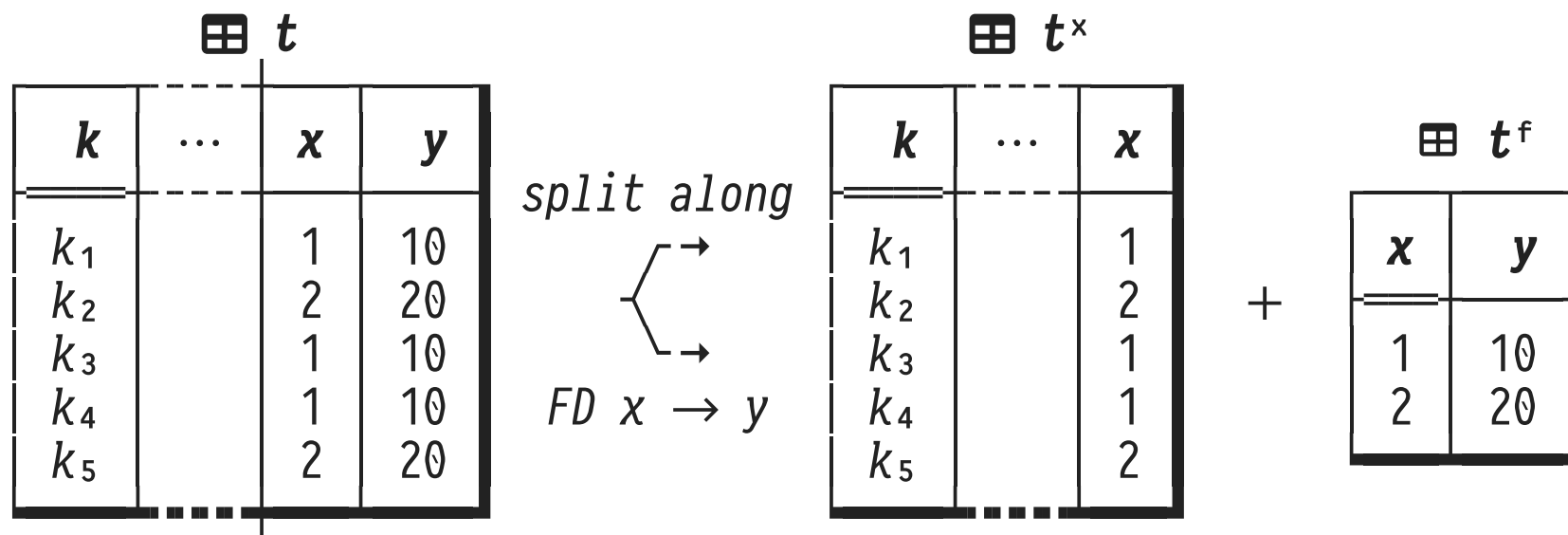
Schedule of multi-hop road trips taken (with drivers and vehicles in use).



- 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)

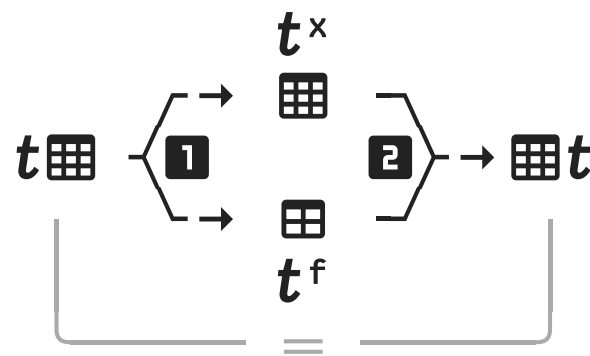
FD splits improve constraint checking and reduce redundancy. But users or apps may require/expect the original non-split table...

- 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



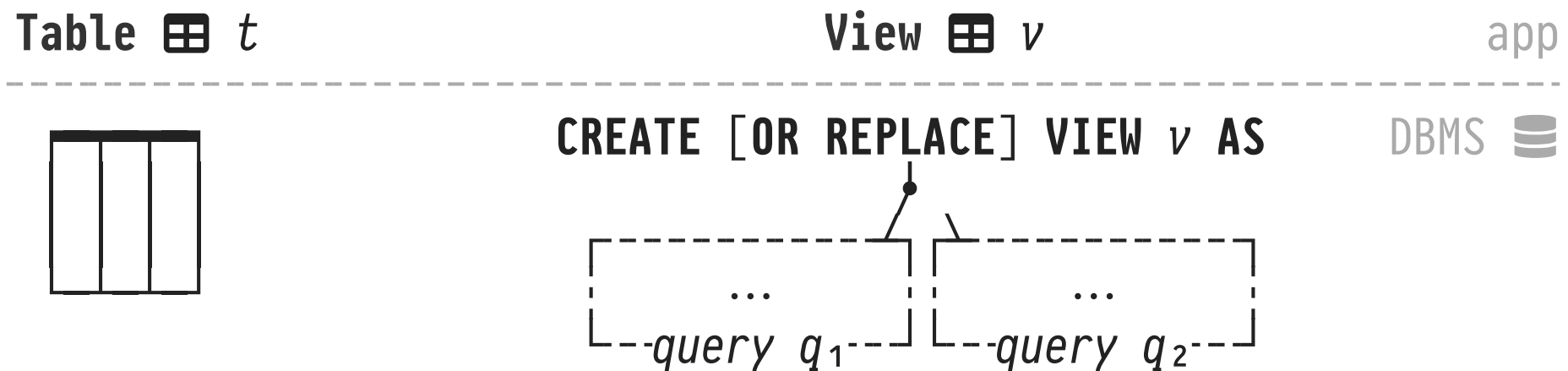
- A **table view** can render this reconstruction of $\boxplus t$ transparent.

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  