

# Advanced SQL

---

## 03 — Arrays and User-Defined Functions

**Torsten Grust**  
**Universität Tübingen, Germany**

# 1 | Arrays: Aliens(?) Inside Table Cells

---

SQL tables adhere to the **First Normal Form (1NF)**: values  $v$  inside table cells are *atomic* w.r.t. the tabular data model:

...	A	...
...	$v$	...

Let us now discuss the **array** data type:

- $v$  may hold an ordered array of elements  $\{x_1, \dots, x_n\}$ .<sup>1</sup>
- SQL treats  $v$  as an atomic unit, but ...
- ... array functions and operators also enable SQL to query the  $x_i$  individually (there's some ↯ with 1NF here).

<sup>1</sup> To the PostgreSQL developer who decided to use  $\{\dots\}$  to denote *arrays*: **No dessert for you today!**

## 2 | Array Types

---

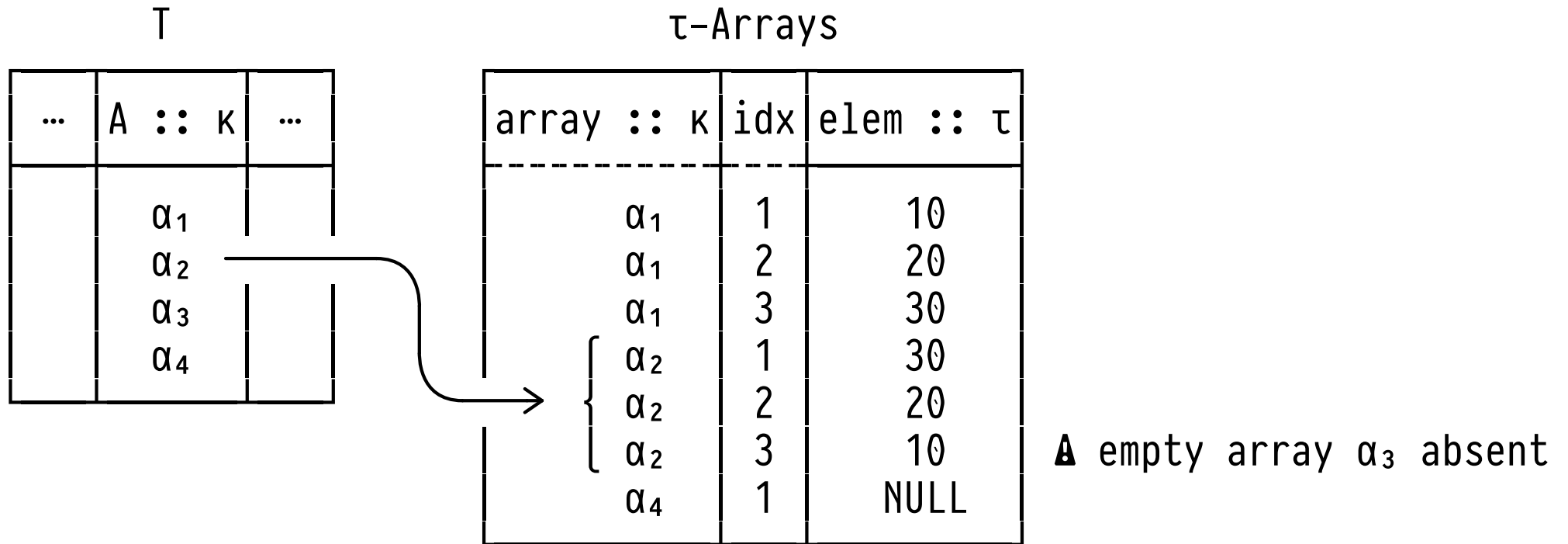
- For type  $\tau$ ,  $\tau[]$  (or  $\tau$  array) is the type of **homogenous arrays of elements of  $\tau$** .
  - $\tau$  may be built-in or user-defined (enums, row types).
  - Array size is unspecified — the array is dynamic.  
(PostgreSQL accepts  $\tau[n]$  but the  $n$  is ignored.)

T

...	A :: int[]	...
...	{10,20,30}	...
...	{30,20,10}	...
...	{}	...
...	{NULL}	...

# “Simulating” Arrays (Tabular Array Semantics)

---



- $\kappa$  denotes a suitable key data type.
- Arrays indexes are of type `int` and 1-based.


### 3 | Array Literals

---

#### One-dimensional array literals of type $\tau$ []:

<code>array[] :: <math>\tau</math> []</code>	empty array of elements of type $\tau$
<code>array[<math>\langle x_1 \rangle, \dots, \langle x_n \rangle</math>]</code>	} all $\langle x_i \rangle$ of type $\tau$
<code>'{<math>\langle x_1 \rangle, \dots, \langle x_n \rangle</math>}' :: <math>\tau</math> []</code>	

#### Multi-dimensional rectangular array literals of type $\tau$ [][]:

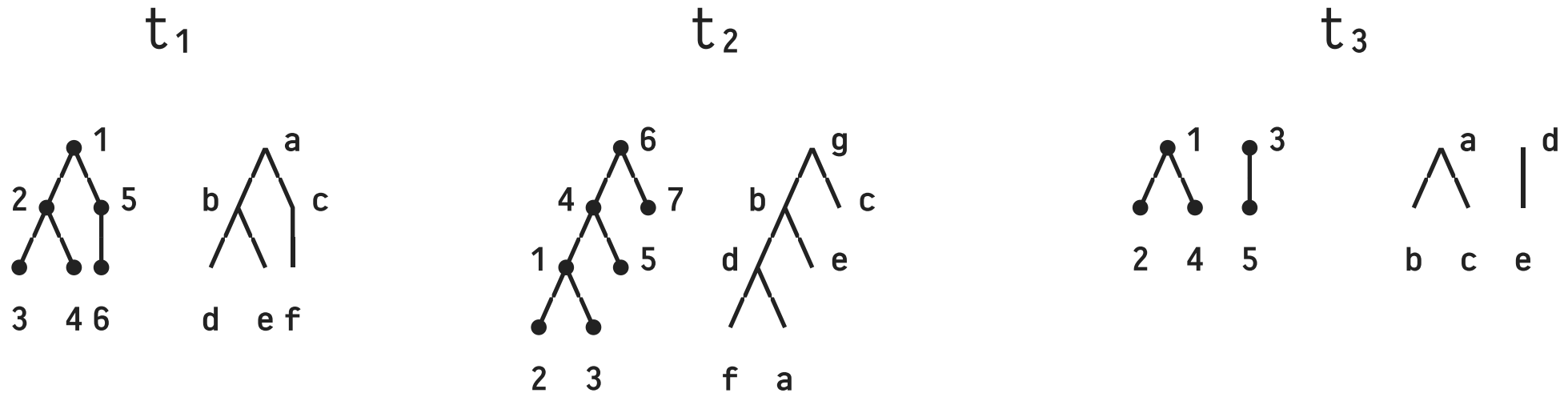
 all sub-arrays need to agree in size

<code>array[<u>array[<math>\langle x_{11} \rangle, \dots, \langle x_{1n} \rangle</math>]</u>, ..., <u>array[<math>\langle x_{k1} \rangle, \dots, \langle x_{kn} \rangle</math>]]</u></code>	1 ■■■■
	⋮ ■■■■
	k ■■■■
	1.....n

`'{{ $\langle x_{11} \rangle, \dots, \langle x_{1n} \rangle$ }, ..., { $\langle x_{k1} \rangle, \dots, \langle x_{kn} \rangle$ }}' ::  $\tau$  [][]`

# Example: Tree Encoding (`parents[i] ≡ parent of node i`)

---



Tree shape and node labels held in separate in-sync arrays:

## Trees

<u>tree</u>	<u>parents</u>	<u>labels</u>
<code>t<sub>1</sub></code>	<code>{NULL,1,2,2,1,5}</code>	<code>{'a','b','d','e','c','f'}</code>
<code>t<sub>2</sub></code>	<code>{4,1,1,6,5,NULL,6}</code>	<code>{'d','f','a','b','e','g','c'}</code>
<code>t<sub>3</sub></code>	<code>{NULL,1,NULL,1,3}</code>	<code>{'a','b','d','c','e'}</code>
	<code>1 2 3 4 5</code>	<code>1 2 3 4 5</code> ← @idx

## Constructing Arrays

---

- **Append/prepend** element ★ to array or
- **concatenate** arrays:

```
array_append (array[x1, ..., xn], ★) ≡ array[x1, ..., xn, ★]  
array_prepend(array[x1, ..., xn], ★) ≡ array[★, x1, ..., xn]
```

```
array_cat(array[x1, ..., xn],  
          array[y1, ..., ym]) ≡ array[x1, ..., xn, y1, ..., ym]
```


- Overloaded operator `||` embraces all of the above:

```
xs ||| ★ ≡ array_append(xs, ★)  
★ ||| xs ≡ array_prepend(xs, ★)  
xs ||| ys ≡ array_cat(xs, ys)
```

## Accessing Arrays: Indexing / Slicing

---

- Array **indexes**  $i$  are 1-based (let  $xs \equiv \text{array}[x_1, \dots, x_n]$ ):

$xs[i]$	$\equiv x_i$	$i \notin \{1, \dots, n\}: \text{NULL}$
$(\text{NULL})[i]$	$\equiv \text{NULL}$	
$xs[\text{NULL}]$	$\equiv \text{NULL}$	
$xs[i:j]$	$\equiv \text{array}[x_i, \dots, x_j]$	$i > j: \text{array}[]$
$xs[i: ]$	$\equiv \text{array}[x_i, \dots, x_n]$	}  requires PostgreSQL 9.6
$xs[ :j]$	$\equiv \text{array}[x_1, \dots, x_j]$	

- Access **last element**  $x_n$ :

$xs[\text{array\_length}(xs, 1)]$	# of elements in dimension 1: $n$
$xs[\text{cardinality}(xs)]$	

↑  
 $\Sigma$  (# of elements) in all dimensions



## Searching for Elements in Arrays

---

Indexing accesses array by position. Instead, **searching** accesses arrays by **contents**.

- Let  $xs \equiv \text{array}[x_1, \dots, x_{i-1}, \star, x_{i+1}, \dots, x_{j-1}, \star, x_{j+1}, \dots, x_n]$  and comparison operator  $\theta \in \{=, <, >, <>, \leq, \geq\}$ :

$x \theta \text{ ANY}(xs) \equiv \exists i \in \{1, \dots, n\}: x \theta xs[i]$

$x \theta \text{ ALL}(xs) \equiv \forall i \in \{1, \dots, n\}: x \theta xs[i]$

$\text{array\_position}(xs, \star) \equiv i$  if  $\star$  not found: NULL

$\text{array\_positions}(xs, \star) \equiv \text{array}[i, j]$  if  $\star$  not found: array[]

$\text{array\_replace}(xs, \star, \blacklozenge) \equiv \text{array}[x_1, \dots, \underset{i}{\blacklozenge}, \dots, \underset{j}{\blacklozenge}, \dots, x_n]$

## 4 : A Bridge Between Arrays and Tables: `unnest` & `array_agg`

---

```
SELECT t.elem
FROM   unnest(array[x1, ..., xn]) AS t(elem)
```

$\underbrace{\hspace{10em}}_{\equiv XS}$

Table t

elem
x <sub>1</sub>
⋮
x <sub>n</sub>

```
SELECT array_agg(t.elem) AS xs
FROM   (VALUES (x1),
             ⋮
             (xn)) AS t(elem)
```

XS
{x <sub>1</sub> , ..., x <sub>n</sub> }

- `unnest(·)`: a *set-returning function*. More on that soon.
- ⚠ Preservation of order of the  $x_i$  is *not* guaranteed...

## Representing Order (Indices) As First-Class Values

---

```
SELECT t.*  
FROM   unnest(array[x1, ..., xn])  
       WITH ORDINALITY AS t(elem, idx)
```

recall ordered aggregates


```
SELECT array_agg(t.elem ORDER BY t.idx) AS xs  
FROM   (VALUES (x1, 1),  
           ⋮  
           (xn, n)) AS t(elem, idx)
```

≡

elem	idx
x <sub>1</sub>	1
⋮	⋮
x <sub>n</sub>	n

≡

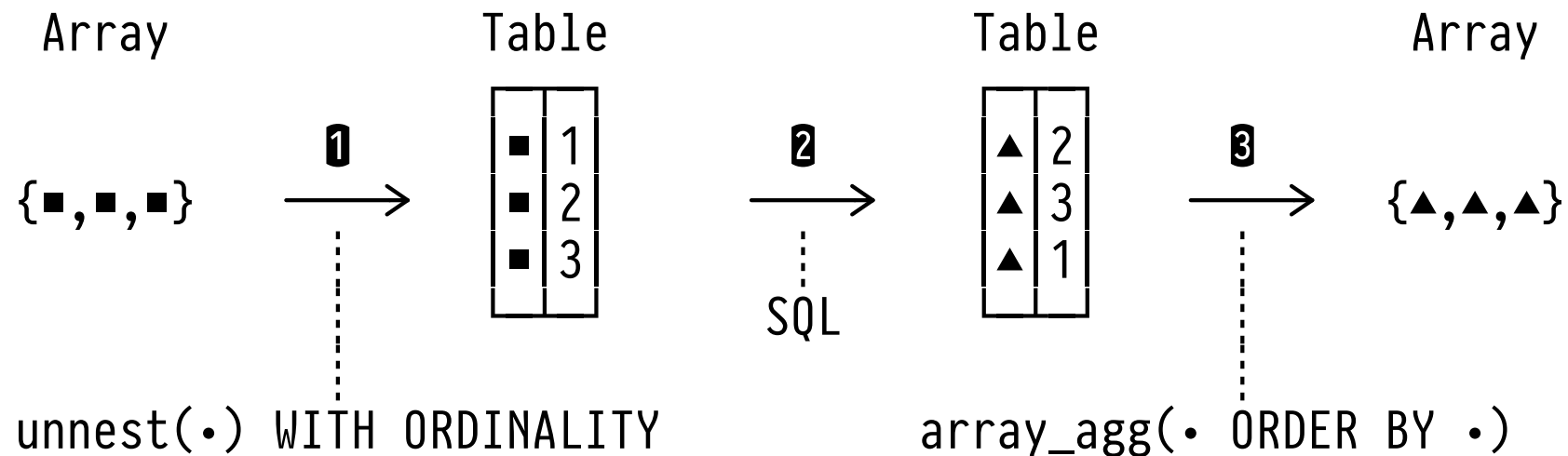
xs
{x <sub>1</sub> , ..., x <sub>n</sub> }

- `<f>(…)` **WITH ORDINALITY** adds a trailing column () of ascending indices 1, 2, ... to the output of function `<f>`.

# A Relational Array Programming Pattern

---

Availability of `unnest(·)` and ordered `array_agg(·)` suggests a pattern for **relational array programming**:



- At ② use the full force of SQL, read/transform/generate elements and their positions at will.
- ①+③ constitute **overhead**: an RDBMS is *not* an array PL.

## 5 | Table-Generating Functions

---

What is the **type** of `unnest(·)`?

- `unnest(·)` establishes a bridge between arrays and SQL's tabular data model:<sup>2</sup>

```
unnest :: τ[] → SETOF τ
```

- In SQL, functions of type  $\tau_1 \rightarrow \text{SETOF } \tau_2$  are known as **set-returning** or **table(-generating) functions**. May be invoked wherever a query expects a table (`FROM` clause).
- Several built-in — may also be **defined by the user**.

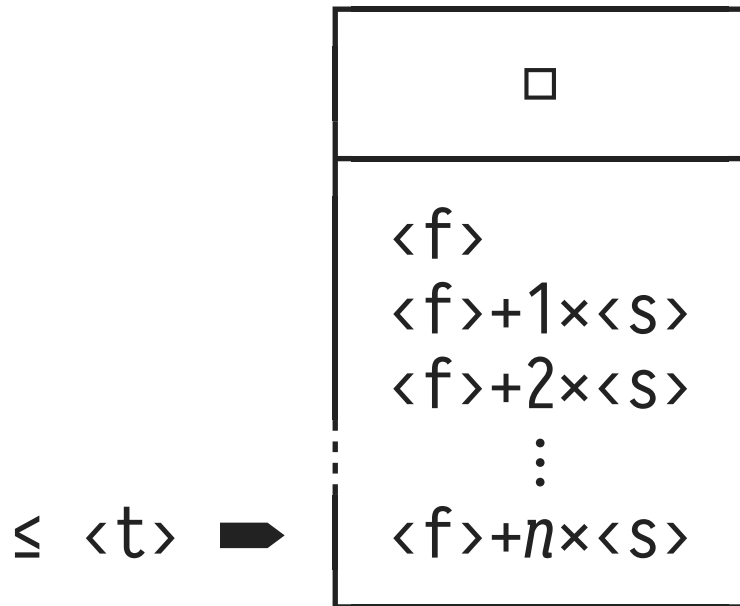
<sup>2</sup> Unfortunate naming again: `SETOF` should probably read `BAGOF` or `TABLE OF`.

## Series and Subscript Generators

---

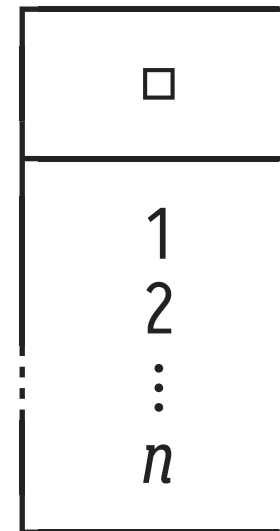
Built-in table-generating functions that generate **tables of consecutive numbers**:

`generate_series(<f>, <t>, <s>)`



`<s>`  $\equiv$  1, if absent  
`<f>`, `<t>`: numbers/timestamps

`generate_subscripts(<xs>, <d>)`



*n*  $\equiv$  `array_length(<xs>, <d>)`  
can also enumerate *n*, ..., 1

## Text Generators (Regular Expression Matching)

---

Use *regular expression*<sup>3</sup> `<re>` to extract **matched substrings** from `<t>` or **split** text `<t>` at defined positions:

1. `regexp_matches(<t>, <re>, 'g')`, yields SETOF `text[]`:  
Generates one array `xs` per match of `<re>` in `<t>`. Element `xs[i]` holds the **match** of the  $i^{\text{th}}$  *capture group* (in (...)).
2. `regexp_split_to_table(<t>, <re>)`, yields SETOF `text`:  
Uses the matches of `<re>` in `<t>` as *separators* to **split** `<t>`. Yields table of  $n+1$  rows if `<re>` matches  $n$  times.

<sup>3</sup> See [regexr.com](http://regexr.com) for tutorials and an interactive playground, for example.

## Breaking Bad: Parse a Chemical Formula (e.g., $C_6H_5O_7^{3-}$ )

---

```
SELECT t.match[1] AS element, -- } extract match details
       t.match[2] AS "# atoms", -- } from the (...)
       t.match[3] AS charge   -- } (capture groups)
FROM   regexp_matches(
       'C6H5O73-',
       '([A-Za-z]+)([0-9]*)([0-9]+[+-])?',
       'g')
AS t(match);
```

element	# atoms	charge
C	6	NULL
H	5	NULL
O	7	3-

} NULL ≡ no match  
}

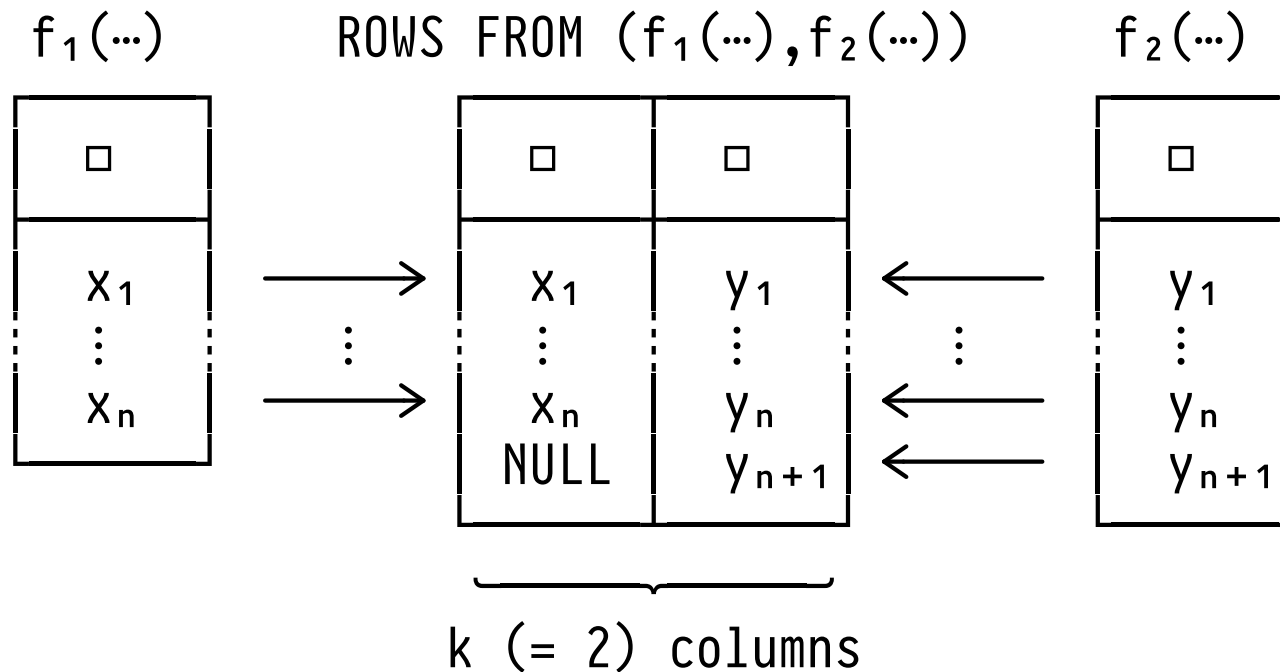


# Zippering Arrays and Table-Generating Functions

---

**Zip:** pair elements based on position (“ORDINALITY join”):

- Zippering table functions  $f_i$ : `ROWS FROM(f1(...),...,fk(...))`
- Zippering arrays  $xs_i$ : `unnest(xs1,xs2,...,xsk)`



## 6 | User-Defined SQL Functions (UDFs)

---

The body of a **user-defined SQL function (UDFs)** evaluates  $n \geq 1$  arbitrary SQL statements in sequence:


```
CREATE FUNCTION <f>(<X1> τ1,...,<Xk> τk) RETURNS τ AS
$$
    <q1>;      --
    <q2>;      -- } evaluate the <qi> in order,
    ⋮           -- ⋮ <qn> defines the result
    <qn>;      -- ]
$$
LANGUAGE SQL [IMMUTABLE];
```

all <q<sub>i</sub>> are read-only ⇒ <f> is free of side effects

- UDF <f> is stored persistently. Remove via **DROP FUNCTION**.

## UDF Types

---

- UDF  $\langle f \rangle$  is  $k$ -ary with type  $\tau_1 \times \dots \times \tau_k \rightarrow \tau$ .
  - **Argument types**  $\tau_i$  must be **atomic** or **row types**.
  - **Overloading** allowed as long as  $(\langle f \rangle, \tau_1, \dots, \tau_k)$  is unique.
  - Limited form of **polymorphism**: any  $\tau_i$  and  $\tau$  may be **anyelement/anyarray/anyenum/anyrange**.
    -  If **any...** occurs more than once in the function signature, *all* occurrences denote the *same* type:

```
f1 :: anyelement × anyelement → boolean
f2 :: anyarray × integer → anyelement
```

elem =

## UDFs Can Return Tables

---

A UDF  $\langle f \rangle :: \tau_1 \times \dots \times \tau_k \rightarrow \tau$  may be of **two flavors**:

### Regular vs. Table-generating UDFs

	<b>atomic <math>\tau</math></b>	<b><math>\tau \equiv \text{SETOF } \tau'</math></b>
If $\langle q_n \rangle^4$ returns no rows, If $\langle q_n \rangle$ returns rows, May be invoked	returns <b>NULL</b> returns the first row wherever <b><math>v::\tau</math></b> is used	returns empty table returns all rows in the <b>FROM</b> clause

- A UDF may invoke **INSERT/DELETE/UPDATE** statements in  $\langle q_i \rangle$  and thus incur side-effects. (Hmm, UDF...<sup>4</sup>)
  - No **IMMUTABLE** option — use **VOLATILE** instead.
  - Use  **$\tau \equiv \text{void}$**  if  $\langle f \rangle$  is all about side-effects or consider adding **... RETURNING  $\langle e_1 \rangle, \dots, \langle e_m \rangle$**  if  $i = n$ .

<sup>4</sup> Recall:  $\langle f \rangle$ 's body evaluates queries  $\langle q_1 \rangle, \dots, \langle q_n \rangle$  (in this order).



## Example UDF: Map Unicode Subscripts

---

Map subscript symbol '0', ..., '9' to its value in {0, ..., 9}:

```
CREATE FUNCTION subscript(s text) RETURNS int AS
$$
  SELECT subs.value::int - 1
  FROM   unnest(array['0', '1', '2', ..., '9'])
        WITH ORDINALITY AS subs(sym, value)
  WHERE  subs.sym = s
$$
LANGUAGE SQL IMMUTABLE;
```

- This is a UDF with atomic return type: yields **NULL** if **s** does not denote a valid subscript.

## Example UDF: Issue Unique ID, Write Protocol

---

Generate ID of the form '<prefix>###' and log time of issue:

```
CREATE FUNCTION new_ID(prefix text) RETURNS text AS
$$
  INSERT INTO issue(id,"when") VALUES
    (DEFAULT, 'now'::timestamp)
  RETURNING prefix || id::text          -- id: just generated
$$
LANGUAGE SQL VOLATILE;                -- function is side-effecting
```

Table `issue`

<code>id :: serial</code>	<code>when</code>
⋮	⋮
42	2017-05-17 14:25:36.928441
⋮	⋮

## Example Table-Generating UDF: Flatten a 2D-Array

---

Unnest 2D array `xss` in *column-major order*:<sup>5</sup>

```
CREATE OR REPLACE FUNCTION unnest2(xss anyarray)
  RETURNS SETOF anyelement AS
$$
SELECT xss[i][j]
FROM   generate_subscripts(xss,1) _(i),
       generate_subscripts(xss,2) __(j)
ORDER BY j, i -- return elements in column-major order
$$
LANGUAGE SQL IMMUTABLE;
```

-  Intended type is `unnest2 :: τ[][] → SETOF τ`.

<sup>5</sup> Built-in function `unnest(·)` can flatten  $n$ -dimensional arrays in row-major order.



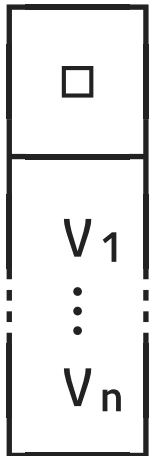
# Table-Generating UDFs: Returning Typed Rows

---

Assume a table-generating UDF  $\langle f \rangle :: \dots \rightarrow \tau$ .

If  $\tau \equiv$

SETOF  $\tau'$   
 $\tau'$  atomic



$v_i :: \tau'$

SETOF  $\tau'$   
 $\tau' \equiv (c_1 :: \tau_1, \dots, c_m :: \tau_m)$

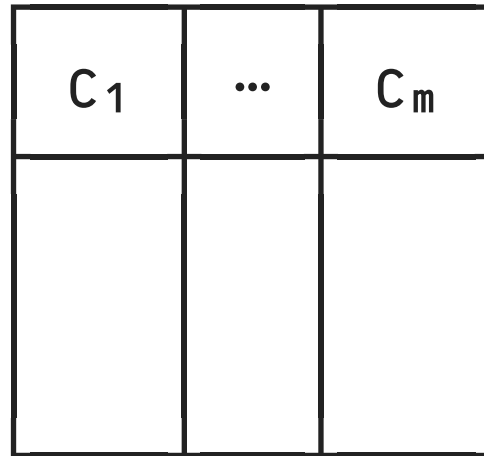
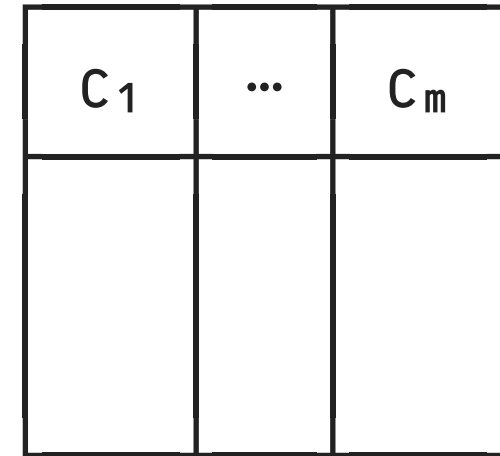


TABLE (c<sub>1</sub>  $\tau_1, \dots, c_m$   $\tau_m$ )



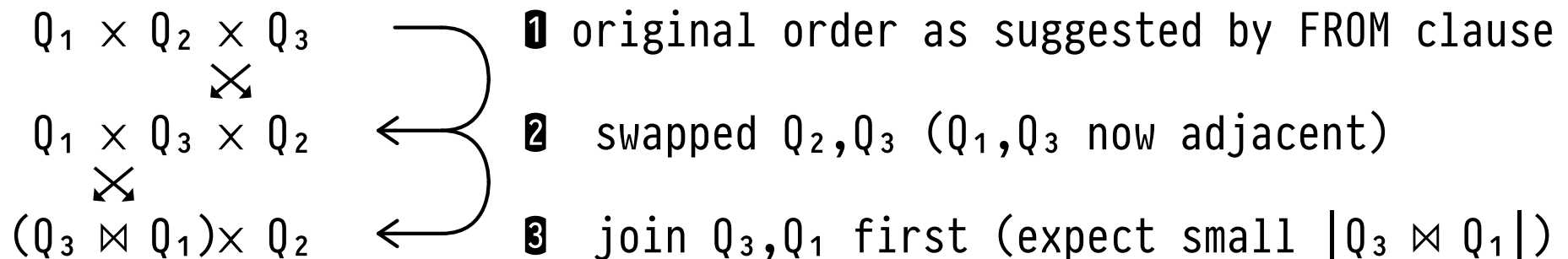
equivalent, but do not need named row type  $\tau'$

## 7 : ', ' in the FROM Clause and Row Variable References

---

```
SELECT ...  
FROM   Q1 AS t1, Q2 AS t2, Q3 AS t3 -- ti<j not free in Qj
```

- Q: Why is  $t_{i<j}$  not usable in  $Q_j$ ?
- A: “... the ', ' in FROM is commutative and associative...”.  
Query optimization might rearrange the  $Q_j$ :



## But Dependent Iteration in **FROM** is Useful...

---

Recall (find largest label in each tree  $t_1$ ):

```
SELECT t1.tree, MAX(t2.label) AS "largest label"
--           Q1                Q2
--           └───┬──────────┬──────────┘
FROM   Trees AS t1, unnest(t1.labels) AS t2(label)
GROUP BY t1.tree;
           ▲
           ⚡
```

- **Dependent iteration** (here:  $Q_2$  depends on  $t_1$  defined in  $Q_1$ ) has its uses and admits intuitive query formulation.
- $\Rightarrow$  Exception: the arguments of table-generating functions may refer to row variables defined earlier (like  $t_1$ ).

## LATERAL:<sup>6</sup> Dependent Iteration for Everyone

---

Prefix  $Q_j$  with **LATERAL** in the **FROM** clause to announce dependent iteration:

```
SELECT ...  
FROM   Q1 AS t1, ..., LATERAL Qj AS tj, ...  
                        ▲  
                        may refer to t1, ..., tj-1
```

- Works for *any* table-valued SQL expression  $Q_j$ , subqueries in (...) in particular.
  - Good style: be explicit and use **LATERAL** even with table functions.

<sup>6</sup> Lateral /'læt(ə)rəl/ a. [Latin lateralis]: *sideways*

## LATERAL: SQL's for each-Loop

---

LATERAL admits the formulation of **nested-loops** computation:

```
SELECT e
FROM   Q1 AS t1, LATERAL Q2 AS t2, LATERAL Q3 AS t3
```

is evaluated just like this nested loop:

```
for t1 in Q1
  for t2 in Q2(t1)
    for t3 in Q3(t1,t2)
      return e(t1,t2,t3)
```

- Convenient, intuitive, and perfectly OK.  
But much like hand-cuffs for the query optimizer. ⚠

## LATERAL Example: Find the Top *n* Rows Among a Peer Group

---

Which are **the three tallest** two- and four-legged dinosaurs?

```
SELECT locomotion.legs, tallest.species, tallest.height
FROM   (VALUES (2), (4)) AS locomotion(legs),
       LATERAL (SELECT d.*
                FROM   dinosaurs AS d
                WHERE  d.legs = locomotion.legs ←
                ORDER BY d.height DESC
                LIMIT 3) AS tallest
```

legs	species	height
2	Tyrannosaurus	7
2	Ceratosaurus	4
2	Spinosaurus	2.4
4	Supersaurus	10
4	Brachiosaurus	7.6
4	Diplodocus	3.6

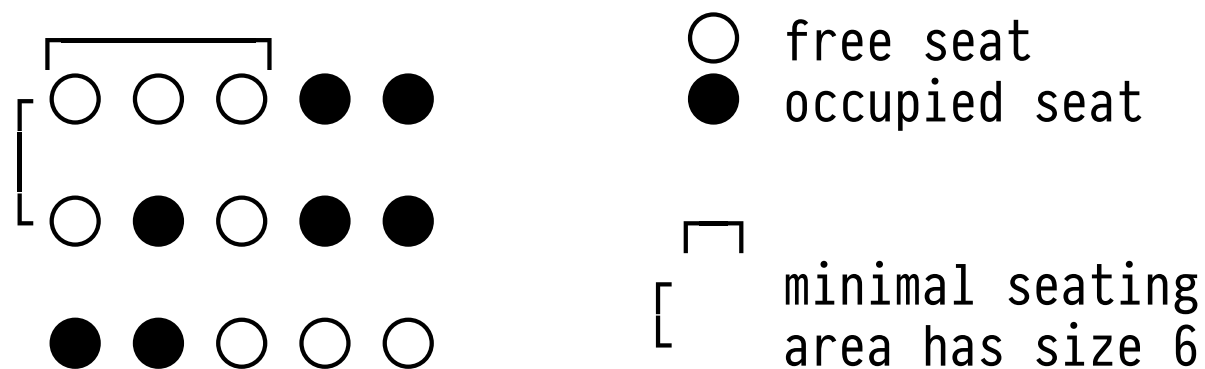
## 8 : Y ACM ICPC: Finding Seats

---

ACM ICPC Task **Finding Seats** (South American Regionals, 2007)

“ $K$  friends go to the movies but they are late for tickets. To sit all nearby, they are looking for  $K$  free seats such that the rectangle containing these seats has minimal area.”

- Assume  $K = 5$ :



## Y Finding Seats: Parse the ICPC Input Format

---

- Typical ICPC character-based input format:

...XX <sup>C<sub>R</sub></sup>	·	free seat
.X.XX <sup>C<sub>R</sub></sup>	X	occupied seat
XX...	<sup>C<sub>R</sub></sup>	new line

- **Parse into table** making seat position/status explicit:

Table `seats`

<u>row</u>	<u>col</u>	<u>taken?</u>
1	1	false
1	2	false
1	3	false
1	4	true
⋮	⋮	⋮
3	5	false



## Y Finding Seats: Parse the ICPC Input Format (Table `seats`)

---

```
\set cinema '...XX\n.X.XX\nXX...'

SELECT row.pos, col.pos, col.x = 'X' AS "taken?"
FROM   -- rows
       unnest(string_to_array(:cinema, '\n'))
       WITH ORDINALITY AS row(xs, pos),
       -- columns
       LATERAL unnest(string_to_array(row.xs, NULL))
              WITH ORDINALITY AS col(x, pos)
```

- `string_to_array(:cinema, '\n')` yields an array of three row strings: `{'...XX', '.X.XX', 'XX...'}`.
- `string_to_array(row.xs, NULL)` splits string `row.xs` into an array of individual characters (= seats).

## Y Finding Seats: A Problem Solution (Generate and Test)

---

- **Query Plan:** 

1. Determine the extent (*rows* × *cols*) of the cinema seating plan.
2. **Generate all** possible north-west (*nw*) and south-east (*se*) corners of rectangular seating areas:
  - For each such 「*nw,se*」 rectangle, scan its seats and **test** whether the number of free seats is  $\geq K$ .
  - If so, record *nw* together with the rectangle's *width/height*.
3. Among these rectangles with sufficient seating space, select those with minimal area.

## Y Finding Seats: Generating All Possible Rectangles

---

Generate all 「nw,se」 corners for rectangles up to maximum size *rows* × *cols*:

```
SELECT ROW(row_nw, col_nw) AS nw,  
       ROW(row_se, col_se) AS se  
FROM   generate_series(1, rows)           AS row_nw,  
       generate_series(1, cols)          AS col_nw,  
       LATERAL generate_series(row_nw, rows) AS row_se,  
       LATERAL generate_series(col_nw, cols) AS col_se
```

Generates  $\binom{rows}{\sum_{r=1}^{} r} \times \binom{cols}{\sum_{c=1}^{} c}$  rectangles  $\Rightarrow$  test/filter early!