

# Advanced SQL

---

## 05 — Recursion

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

## Computational Limits of SQL

---

SQL has grown to be an **expressive data-oriented language**. Intentionally, it has *not* been designed as a general-purpose programming language:

1. *SQL does not loop forever:*

Any SQL query is expected to **terminate**, regardless of the size/contents of the input tables.

2. *SQL can be **evaluated efficiently**:*

A SQL query over table **T** of  $c$  columns and  $r$  rows can be evaluated in  $O(r^c)$  space and time.<sup>1</sup>


<sup>1</sup> SQL cannot compute the set of all subsets of rows in **T** which requires  $O(2^r)$  space, for example.

## A Giant Step for SQL

---

The addition of **recursion** to SQL changes everything:

**Expressiveness** SQL becomes a **Turing-complete language** and thus a general-purpose PL (albeit with a particular flavor).

**Efficiency**  **No longer** are queries guaranteed to **terminate** or to be **evaluated with polynomial effort**.

Like a pact with the  — but the payoff is magnificent...

## Recursion in SQL: WITH RECURSIVE

---

### Recursive common table expression (CTE):

#### WITH RECURSIVE

```
<T1>(<C11>, ..., <C1,k1>>) AS (
  <q1> ),
  ⋮
<Tn>(<Cn1>, ..., <Cn,kn>>) AS (
  <qn> )
<q>
```

Queries <q<sub>j</sub>> may refer **to all** <T<sub>i</sub>>

<q> may refer **to all** <T<sub>i</sub>>

- In particular, any <q<sub>j</sub>> may refer to itself (⊙)! Mutual references are OK, too. (Think `letrec` in FP.)
- Typically, final query <q> performs post-processing only.

## Shape of a Self-Referential Query

---

### WITH RECURSIVE

```
T(c1, ..., ck) AS (  
  q0 -- common schema of q0 and q∅(·)  
      -- base case query, evaluated once  
  
  UNION [ ALL ] -- either UNION or UNION ALL  
  
  q∅(T) -- recursive query refers to T itself,  
        -- evaluated repeatedly  
)  
q(T) -- final post-processing query
```

- Semantics in a nutshell:

$$q(q_{\emptyset}(\dots q_{\emptyset}(q_{\emptyset}(q_0))\dots)) \cup \dots \cup q_{\emptyset}(q_{\emptyset}(q_0)) \cup q_{\emptyset}(q_0) \cup q_0$$

repeated evaluation of  $q_{\emptyset}$  (when to stop?)

## Semantics of a Self-Referential Query (**UNION** Variant)

---

Iterative and recursive semantics—both are equivalent:

<pre>iterate(<math>q\vartheta</math>, <math>q_0</math>):   <math>r \leftarrow q_0</math>   <math>t \leftarrow r</math>   <b>while</b> <math>t \neq \emptyset</math>     <math>t \leftarrow q\vartheta(t) \setminus r</math>     <math>r \leftarrow r \uplus t</math>   <b>return</b> <math>r</math></pre>	<pre>recurse(<math>q\vartheta</math>, <math>r</math>):   <b>if</b> <math>r \neq \emptyset</math> <b>then</b>     <b>return</b> <math>r \uplus \text{recurse}(q\vartheta, q\vartheta(r) \setminus r)</math>   <b>else</b>     <b>return</b> <math>\emptyset</math></pre>
---	---

- Invoke the recursive variant with `recurse( $q\vartheta$ ,  $q_0$ )`.
- $\uplus$  denotes disjoint set union,  $\setminus$  denotes set difference.
- `$q\vartheta(\cdot)$`  evaluated over **the *new* rows found in the last iteration/recursive call**. Exit if there were no new rows.

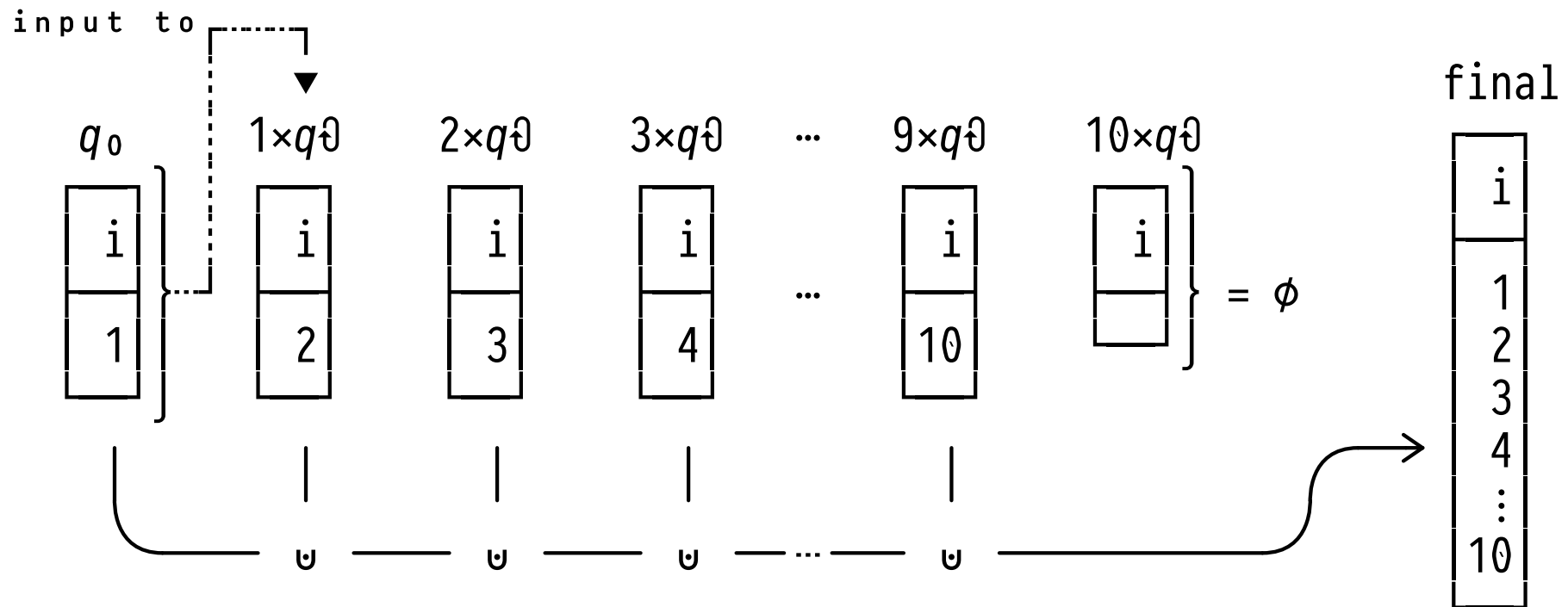


## Y A Home-Made `generate_series()`

---

- Assume `<from> = 1`, `<to> = 10`:

New rows in table **series** after evaluation of...






## Semantics of a Self-Referential Query (**UNION ALL** Variant)

---

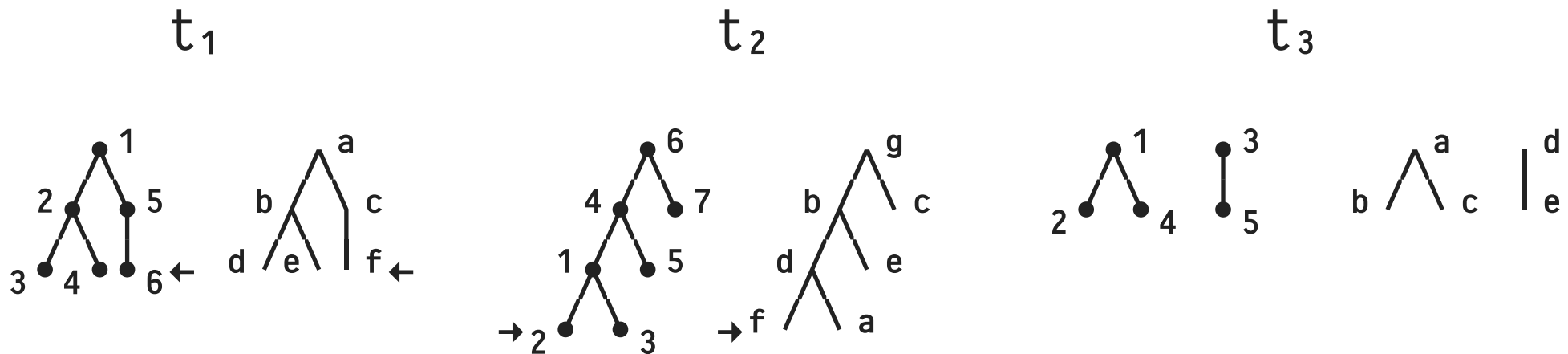
With **UNION ALL**, recursive query  $q\theta$  sees **all** rows added in the last iteration/recursive call:

```
iteratea11( $q\theta$ ,  $q_0$ ): | recursea11( $q\theta$ ,  $r$ ):  
   $r \leftarrow q_0$  |   if  $r \neq \phi$  then  
   $t \leftarrow r$  |   | return  $r \uplus$  recursea11( $q\theta$ ,  $q\theta(r)$ )  
  while  $t \neq \phi$  |   else  
  |  $t \leftarrow q\theta(t)$  |   | return  $\phi$   
  |  $r \leftarrow r \uplus t$   
  return  $r$  |
```

- Invoke the recursive variant via  $\text{recurse}^{a11}(q\theta, q_0)$ .
- $\uplus$  denotes bag (multiset) union.
-  Need to be extra careful to control termination.

# 1 : Y Traverse the Paths from Nodes 'f' to their Root

---



Array-based tree encoding (parent of node  $n \equiv \text{parents}[n]$ ):

## Trees

tree	parents ( $\square \equiv \text{NULL}$ )	labels
$t_1$	$\{\square, 1, 2, 2, 1, 5\}$	$\{'a', 'b', 'd', 'e', 'c', 'f'\}$
$t_2$	$\{4, 1, 1, 6, 5, \square, 6\}$	$\{'d', 'f', 'a', 'b', 'e', 'g', 'c'\}$
$t_3$	$\{\square, 1, \square, 1, 3\}$	$\{'a', 'b', 'd', 'c', 'e'\}$
	1 2 3 4 5 6 7	1 2 3 4 5 6 7 ← node

## Y Traverse the Paths from Nodes 'f' to their Root

---

**WITH RECURSIVE**

paths(tree, node) AS (

SELECT t.tree, array\_position(t.labels, 'f') AS node

FROM Trees AS t

**UNION**

SELECT t.tree, t.parents[p.node] AS node

FROM paths AS p,

Trees AS t

WHERE p.tree = t.tree

)

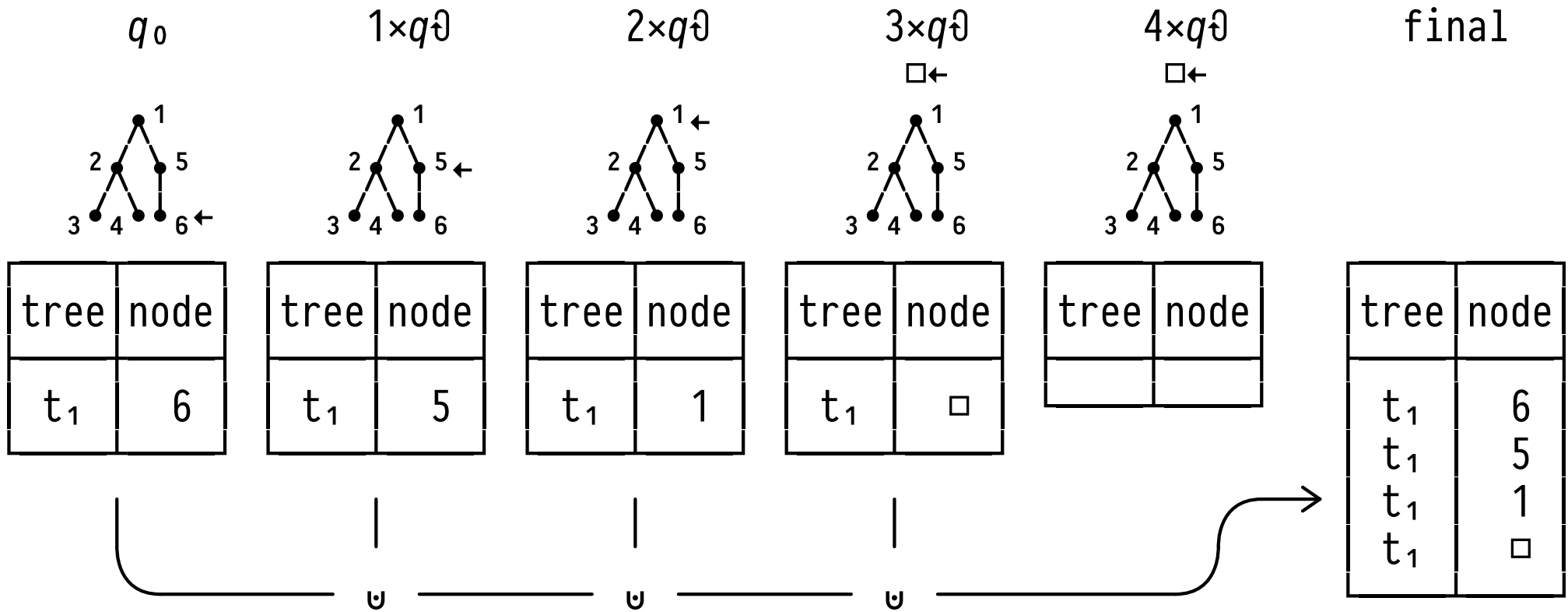
**TABLE** paths

$(t, n) \in \text{paths} \iff$  node  $n$  lies on path from 'f' to  $t$ 's root

# Y A Trace of the Path in Tree $t_1$

---

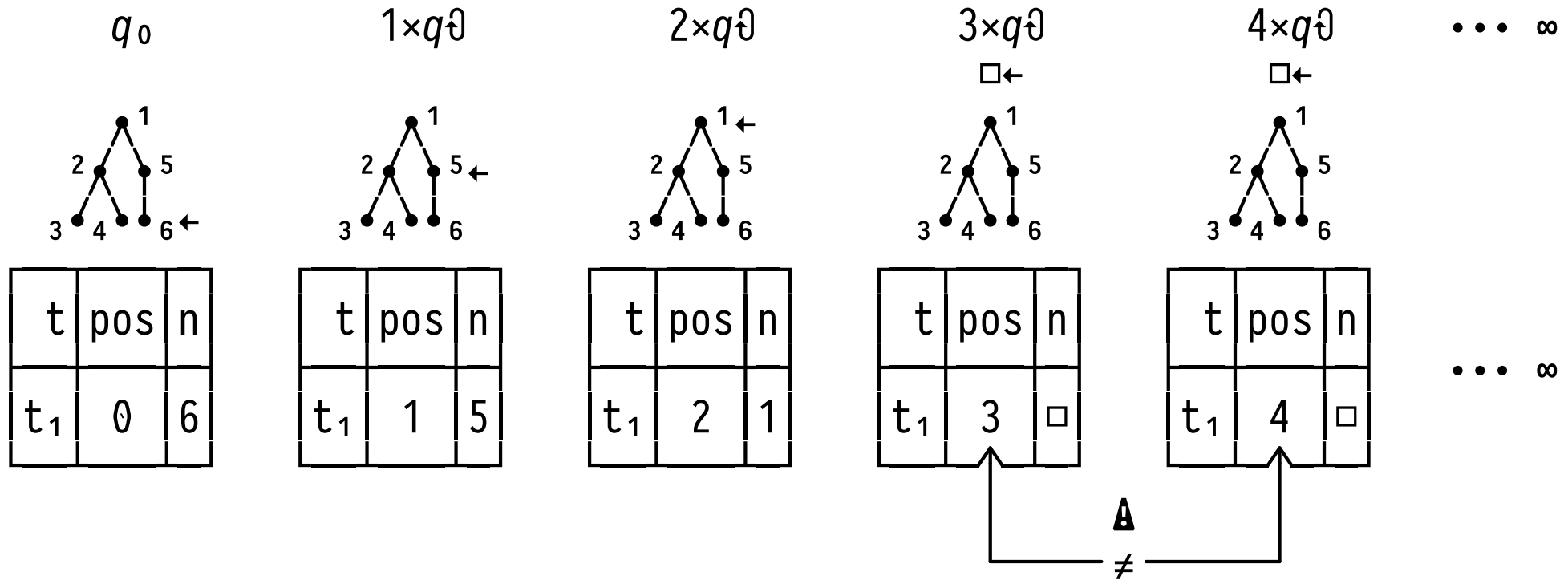
New rows produced by...




- $4 \times q \emptyset$  yields no new rows (recall:  $t.\text{parents}[\text{NULL}] \equiv \text{NULL}$ ).

# Y Ordered Path in Tree $t_1$ (New Rows Trace)

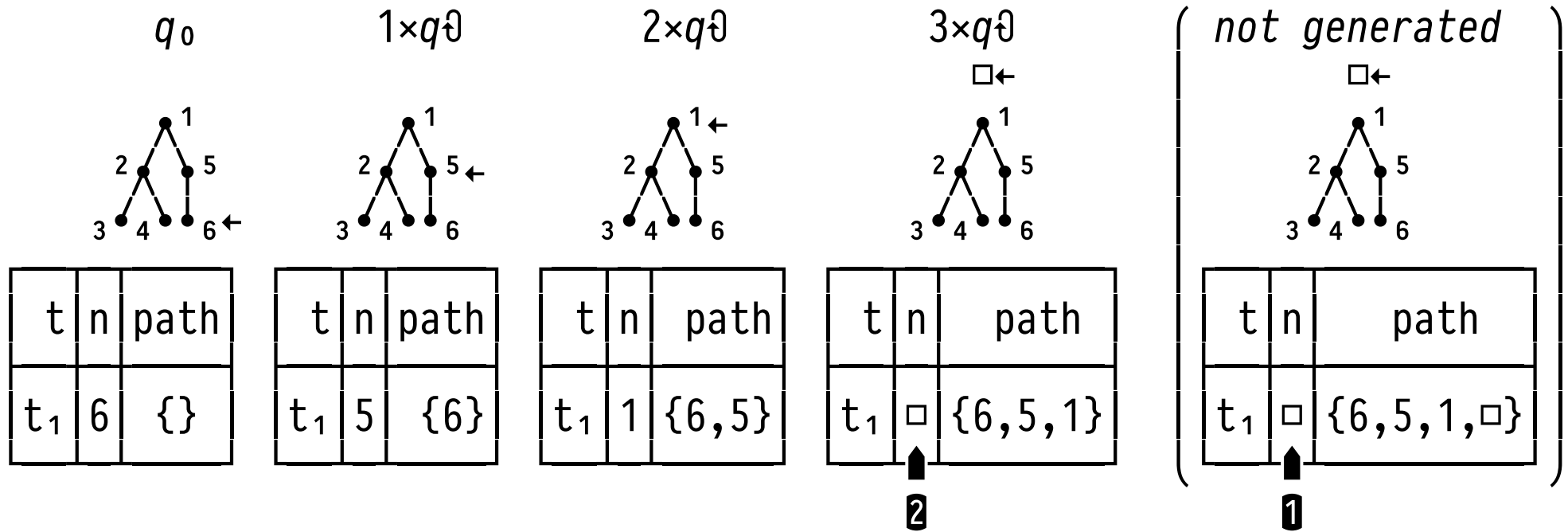
---



The (non-)generation of new rows to ensure termination is the user's responsibility — a common source of .

# Y Path as Array in Tree $t_1$ (New Rows Trace)

---



- ① Ensure termination (enforce  $\emptyset$ ): filter on  $n \neq \square$  in  $q\emptyset$ .
- ② Post-process: keep rows of last iteration ( $n = \square$ ) only.