

# CS 348 Lecture 6

## Recursion in SQL, Datalog and SQL Programming

Semih Salihoğlu

Jan 25<sup>th</sup>, 2022



UNIVERSITY OF  
**WATERLOO**



Data  
Systems  
Group

# Outline For Today

---

1. SQL Recursive Query Support
  - Recursion Motivation & FixedPoint Subroutine
  - WITH and WITH RECURSIVE Clauses
  - Monotonicity
  - Linear vs Non-Linear Recursion
  - Mutual Recursion
  - Important Note About Convergence of Recursive Queries
2. Datalog: A More Elegant Query Languages For Recursion
3. SQL Programming

# Strengths and Limitations of SQL So Far

---

## Strengths:

- Excellent fit for tasks using fundamental set operations:
  - projection, joins, filtering, grouping etc. and combinations
- Very high-level:
  - I. Declarative: abstracts users away from low-level computations
  - II. Physical data independence: abstracts away low-level storage

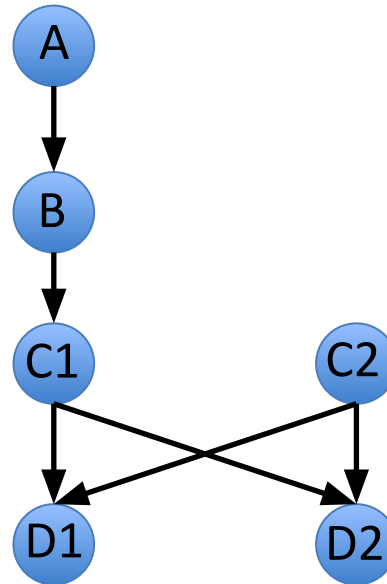
## Limitations:

- Is not Turing-complete
- More specifically: Cannot express recursive computations
- Historically: Recursion was an afterthought when standardizing SQL

# Motivating Example 1: Transitive Closure

- Ex: Given academic <(co-)supervisor, student> relationships:
  - Find all academic ancestors/descendants of an academic

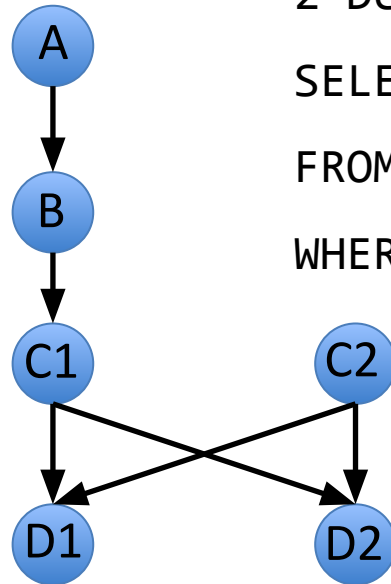
Advisor	
<u>supervisor</u>	<u>student</u>
C1	D1
C1	D2
C2	D1
C2	D2
B	C1
A	B



Ancestors	
<u>anc</u>	<u>desc</u>
...	...
...	...
...	...
...	...
...	...
...	...

# Motivating Example 1: Transitive Closure

Advisor	
<u>sup</u>	<u>stu</u>
C1	D1
C1	D2
C2	D1
C2	D2
B	C1
A	B



2 Degree Ancestors Query:

```
SELECT Adv1.sup AS anc, Adv2.stu as desc
FROM Advisor Adv1, Advisor Adv2
WHERE Adv1.stu = Adv2.sup
```

□ Can find ancestors at any fixed degree, e.g., 1<sup>st</sup>, 2<sup>nd</sup> or 4<sup>th</sup> degree

□ If max depth d is known: union all possible queries upto degree d:

(SELECT \* FROM Advisor) UNION

(SELECT Adv1.sup, Adv2.stu FROM Adv1,Adv2 WHERE Adv1.stu=Adv2.sup) UNION

... (SQL Query for d-degree ancestors)

□ But cannot express arbitrary depths

# Motivating Example 1: Transitive Closure

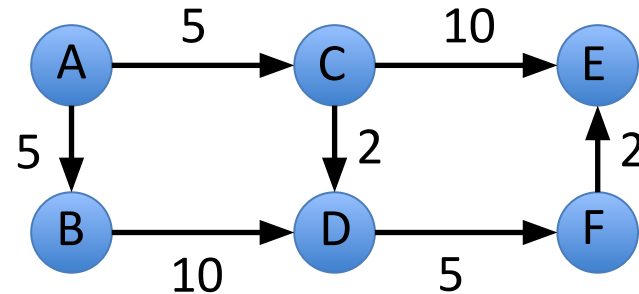
---

- Historical Fact: killer app of graph DBMSs before relational systems was the “parts explosion query” equivalent transitive closure
- Ask me offline if you want to hear more about this history!

# Motivating Example 2: Shortest Paths

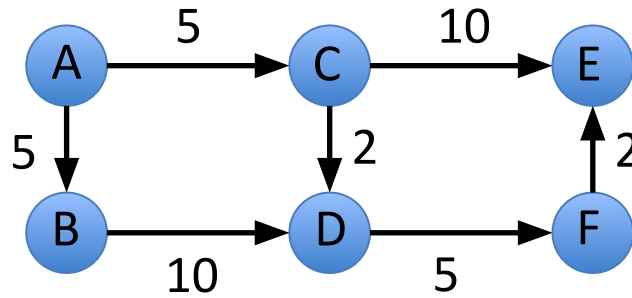
- Many other queries builds on top of transitive closure.
- Ex: Given flights <from, to, price> relationships:
  - Find cheapest paths from A to F

Flights		
<u>from</u>	<u>to</u>	<u>cost</u>
A	B	5
A	C	5
B	D	10
C	D	2
C	E	10
D	F	5
F	E	2



# Motivating Example 2: Shortest Paths

Flights		
<u>from</u>	<u>to</u>	<u>cost</u>
A	B	5
A	C	5
B	D	10
C	D	2
C	E	10
D	F	5
F	E	2



3-edge Paths Query:

```
SELECT F1.from, F3.to, F1.cost+F2.cost+F3.cost as cost
FROM Flights F1, Flights F2, Flights F3
WHERE F1.to=F2.from AND F2.to=F3.from
```

- Can find all (shortest) paths with any fixed number, e.g., k, edges
- If max depth d is known (\*and (directed) graph is acyclic\*)
  - i. Union all paths with up to d edges. Call this relation AllPaths:
  - ii. `SELECT from, to, min(cost) FROM AllPaths`
- But cannot express arbitrary depths



# Solution: Recursive “Fixed Point” Computations

- Transitive closure (TC) and all paths (or shortest paths which depend on all paths) are inherently recursive properties of graphs
- Example: TC of  $v$ : all nodes that  $v$  can directly or indirectly reach
- Computing them require a recursive computation subroutine:
  - High-level Recursive Subroutine for TC:

FixedPoint(fnc  $F$  w/  $T$  as input):

Important Questions:

*1. When does fp converge?*

*2. When is it unique?*

$$T_{\text{prev}} = \emptyset$$

$$T_{\text{new}} = F(T_{\text{prev}}) \quad // \text{ 1}^{\text{st}} \text{ degree ancestors}$$

while ( $T_{\text{prev}} \neq T_{\text{new}}$ ):

$$T_{\text{prev}} = T_{\text{new}}$$

$$T_{\text{new}} = F(T_{\text{prev}}) \quad // \text{ compute up to next-degree ancestors}$$

Equivalently: Compute  $T_0 = \emptyset$ ;  $T_1 = F(T_0)$ ;  $T_2 = F(T_1)$ ; ... until  $T_i = T_{i+1}$

# SQL WITH

- A convenient way to define sub-queries and temporary views

WITH **R1** AS **Q1**

- $R_i$  is the result of  $Q_i$

**R2** AS **Q2**

- $R_i$  visible to  $R_{i+1}, \dots, R_n$

...

- Can explicitly specify schema as

**Rn** AS **Qn**

- $R_1(\text{foo}, \text{bar})$  AS  $Q_1$  o.w inherits from  $Q$

$Q$  // a query that can use existing tables \*and\* **R1, ..., Rn**

- Ex:

```
WITH Deg2Anc AS (SELECT Adv1.sup AS anc, Adv2.stu as desc
                  FROM Advisor Adv1, Advisor Adv2
                  WHERE Adv1.stu = Adv2.sup)
```

```
Deg3Anc AS (....)
```

```
SELECT desc FROM (SELECT * FROM Deg2Anc UNION
```

# SQL WITH RECURSIVE

- WITH can be suffixed with RECURSIVE keyword

WITH RECURSIVE

R1 AS Q1 → Can reference R1

R2 AS Q2

...

Rn AS Qn

Q // a query that can use existing tables \*and\* R1, ..., Rn

- Semantics of “RECURSIVE T AS Q”: run FixedPoint subroutine

$$T_0 = \emptyset$$

$$T_1 = Q \text{ (but use } T_0 \text{ for } T)$$

$$T_2 = Q \text{ (but use } T_1 \text{ for } T)$$

until  $T = T$

*Note: In SQL standard RECURSIVE is bound to specific Ri. We will and some systems bound it to WITH, so all Ri.*

# TC: ATTEMPT 1

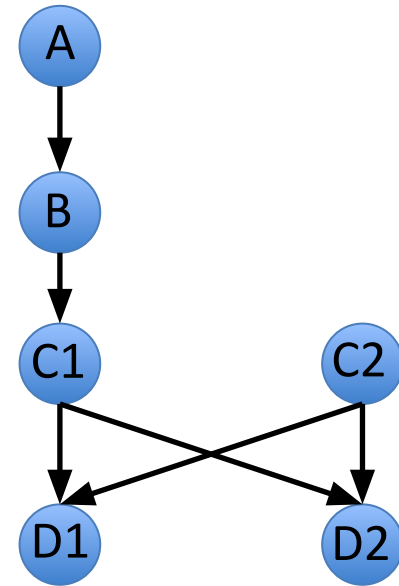
WITH RECURSIVE **Ancestors**(anc, desc) AS (

**SELECT** Ancestor.anc, Adv.stu

**FROM** Ancestor, Advisor

**WHERE** Ancestor.desc = Advisor.sup)

Advisor	
<u>sup</u>	<u>stu</u>
C1	D1
C1	D2
C2	D1
C2	D2
B	C1
A	B



- Problem? Ancestor starts as  $\emptyset$
- Common fix: UNION with a 2<sup>nd</sup> query that init's Ancestor to Advisor
- Common WITH RECURSIVE query template:

WITH RECURSIVE R AS ( e<sub>B</sub> **UNION** e<sub>R</sub> )

non-recursive  
“base” query

recursive query

# TC: ATTEMPT 2: Union w/ a “Base” Case

WITH RECURSIVE **Ancestors**(anc, desc) AS (

base query

{

SELECT sup as anc, stu as desc

FROM Advisor

}

UNION

→

*duplicate eliminating union*

recursive query

{

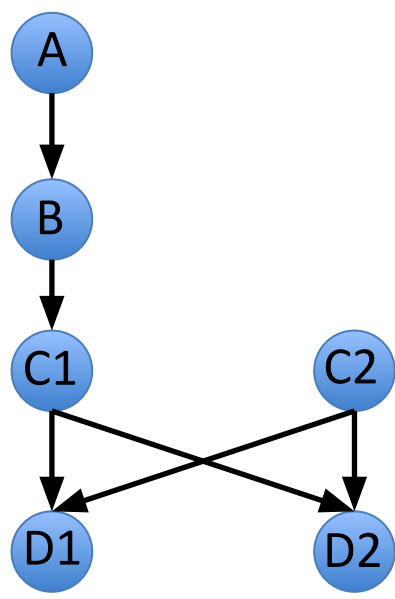
SELECT Ancestor.anc, Adv.stu

FROM Ancestor, Advisor

WHERE Ancestor.desc = Advisor.sup

}

Advisor	
<u>sup</u>	<u>stu</u>
C1	D1
C1	D2
C2	D1
C2	D2
B	C1
A	B



$Q_R = Anc_0 \bowtie Advisor$

Anc <sub>0</sub>	
<u>anc</u>	<u>desc</u>

Q <sub>B</sub>	
<u>anc</u>	<u>desc</u>
C1	D1
C1	D2
C2	D1
C2	D2
B	C1
A	B

Q <sub>R</sub>	
<u>anc</u>	<u>desc</u>

Anc <sub>1</sub>	
<u>anc</u>	<u>desc</u>
C1	D1
C1	D2
C2	D1
C2	D2
B	C1
A	B

Ancestor<sub>1</sub> =

U

=

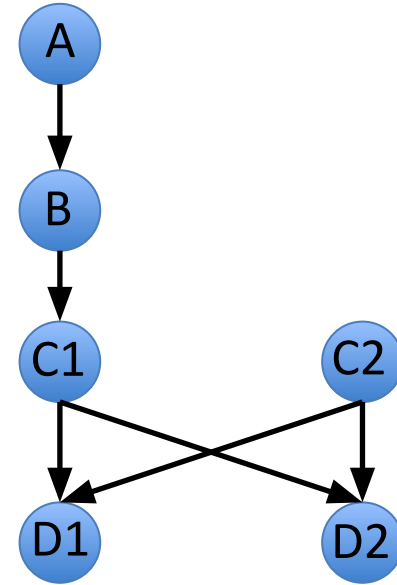
Is Anc<sub>1</sub> - Ans<sub>0</sub> = ∅?  
No: Repeat

# TC: ATTEMPT 2: Union w/ a “Base” Case

Anc <sub>1</sub>	
<u>anc</u>	<u>desc</u>
C1	D1
C1	D2
C2	D1
C2	D2
B	C1
A	B

→ All 1-degree  
ancestors

Advisor	
<u>sup</u>	<u>stu</u>
C1	D1
C1	D2
C2	D1
C2	D2
B	C1
A	B



$$Q_R = \text{Anc}_1 \bowtie \text{Advisor}$$

Ancestor<sub>2</sub> =

Q <sub>B</sub>	
<u>anc</u>	<u>desc</u>
C1	D1
C1	D2
C2	D1
C2	D2
B	C1
A	B

Q <sub>R</sub>	
<u>anc</u>	<u>desc</u>
A	C1
B	D1

U

=

Is Anc<sub>2</sub> - Ans<sub>1</sub> = ∅?  
No: Repeat

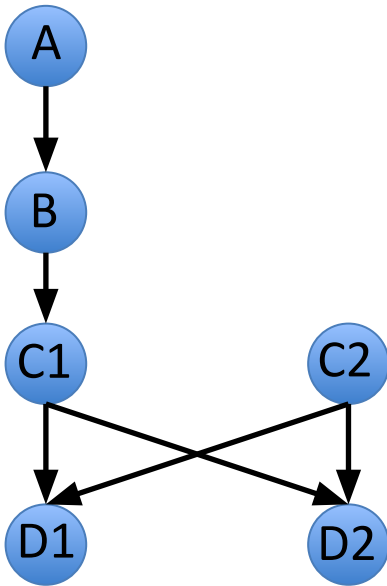
Anc <sub>2</sub>			
<u>anc</u>	<u>desc</u>	<u>anc</u>	<u>desc</u>
C1	D1	A	C1
C1	D2	B	D1
C2	D1		
C2	D2		
B	C1		
A	B		

# TC: ATTEMPT 2: Union w/ a “Base” Case

Anc <sub>2</sub>			
<u>anc</u>	<u>desc</u>	<u>anc</u>	<u>desc</u>
C1	D1	A	C1
C1	D2	B	D1
C2	D1		
C2	D2		
B	C1		
A	B		

→ All 1- and 2-degree ancestors

Advisor	
<u>sup</u>	<u>stu</u>
C1	D1
C1	D2
C2	D1
C2	D2
B	C1
A	B



$$Q_R = \text{Anc}_2 \bowtie \text{Advisor}$$

Ancestor<sub>3</sub> =

Q <sub>B</sub>	
<u>anc</u>	<u>desc</u>
C1	D1
C1	D2
C2	D1
C2	D2
B	C1
A	B

U

Q <sub>R</sub>	
<u>anc</u>	<u>desc</u>
A	C1
B	D1
A	D1
A	D2

=

Is  $\text{Anc}_3 - \text{Ans}_2 = \emptyset$ ?  
No: Repeat

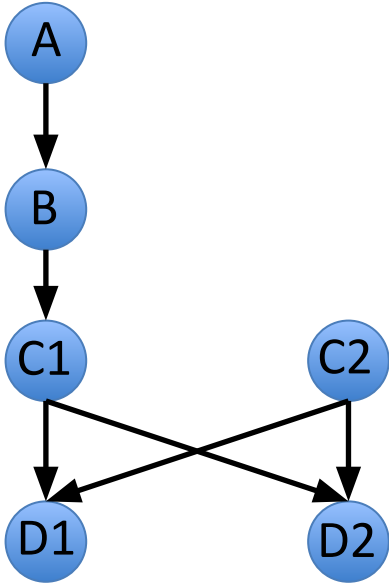
Anc <sub>3</sub>			
<u>anc</u>	<u>desc</u>	<u>anc</u>	<u>desc</u>
C1	D1	A	C1
C1	D2	B	D1
C2	D1	A	D1
C2	D2	A	D2
B	C1		
A	B		

# TC: ATTEMPT 2: Union w/ a “Base” Case

Anc <sub>3</sub>			
<u>anc</u>	<u>desc</u>	<u>anc</u>	<u>desc</u>
C1	D1	A	C1
C1	D2	B	D1
C2	D1	A	D1
C2	D2	A	D2
B	C1		
A	B		

All 1-, 2-, and 3-degree ancestors

Advisor	
<u>sup</u>	<u>stu</u>
C1	D1
C1	D2
C2	D1
C2	D2
B	C1
A	B



$$Q_R = Anc_3 \bowtie Advisor$$

Ancestor<sub>4</sub> =

Q <sub>B</sub>	
<u>anc</u>	<u>desc</u>
C1	D1
C1	D2
C2	D1
C2	D2
B	C1
A	B

Q<sub>R</sub>

<u>anc</u>	<u>desc</u>
A	C1
B	D1
A	D1
A	D2

=

Is Anc<sub>4</sub> − Ans<sub>3</sub> = ∅?

Yes: Stop

Anc <sub>4</sub>			
<u>anc</u>	<u>desc</u>	<u>anc</u>	<u>desc</u>
C1	D1	A	C1
C1	D2	B	D1
C2	D1	A	D1
C2	D2	A	D2
B	C1		
A	B		

Final Answer called the **fixed point** of Q.



# Some Comments

- Recall common WITH RECURSIVE query template:

WITH RECURSIVE R AS ( $Q_B$  UNION  $Q_R$ )

- Can use other queries/templates (e.g., multiple base cases)
  - But some restrictions apply (stay tuned)
- Note that fixed-point computation was very well-behaved in TC:
  - Computation converged:
    - In finite steps (and computed a finite relation)
    - No oscillations
- Question: Are there conditions that guarantee convergence to a unique fixed point of  $Q$ ?

# Monotonicity

---

- If we focus on core relational algebra foundation of SQL:
  - Select/project/cross product/join/union/set difference/intersection
  - Ignore group by and aggregations and arithmetic functions etc.
- Theorem: If a recursive  $Q$  is “monotone w.r.t to every relation it contains”, then  $Q$  has a unique and finite fixed point (i.e., the fixed point subroutine is guaranteed to converge)
- Definition:  $Q$  is monotone w.r.t  $R$  iff adding more tuples to  $R$  can not remove tuples from output of  $Q$  (but new tuples can appear)
  - i.e., if each  $t$  that used to be in the output of  $Q$  is guaranteed to remain in output if add more tuples to  $R$  (keeping all else same)

# Monotonicity

---

- Recall each core RA operator except set difference is monotone w.r.t their arguments
- E.g.:  $R \bowtie_p S$  is monotone w.r.t  $R$  and  $S$
- But:  $R - S$  is non-monotone w.r.t  $S$
- Therefore: Any  $Q$  that uses core relational algebraic operations and does not use set difference is monotone
  - =>  $Q$  will converge to a unique fixed point (if recursive)
- Note:  $Q$  can still be monotone even if it contains set difference. But not guaranteed to be.

# Why Does Monotonicity Guarantee A Unique Fixed Point For A Recursive Query?

□ Proof Sketch: Recall fixed point subroutine:

$T_0 = \emptyset$ ;  $T_1 = Q$  (but use  $T_0$  for  $T$ );  $T_2 = Q$  (but use  $T_1$  for  $T$ )

...

□ Note we are assuming we are focusing on core RA:

□ Each value in a column of  $T_i$  is from a value from base relation

Anc <sub>4</sub>				Advisor	
<u>anc</u>	<u>desc</u>	<u>anc</u>	<u>desc</u>	<u>sup</u>	<u>stu</u>
C1	D1	A	C1	C1	D1
C1	D2	B	D1	C1	D2
C2	D1	A	D1	C2	D1
C2	D2	A	D2	C2	D2
B	C1			B	C1
A	B			A	B

□ In any relation, no matter what its schema is, has a finite maximum size.

□ B/c  $Q$  is monotone (specifically w.r.t to  $T$ ):

$T_1 \subset T_2 \subset T_3 \subset \dots$  (must stop b/c finiteness)

i.e.  $T_1 \subset T_2 \subset \dots T_k = T_{k+1}$  (and fp stops)

# Example Non-Monotone Recursive Query 1

WITH RECURSIVE  $T(x)$  AS (

SELECT  $x$  FROM  $R$

UNION

SELECT  $\text{sum}(x)$  as  $x$  FROM  $T$ )

R
$x$
1
2

$T_0$	$T_1$	$T_2$	$T_3$	$T_4$
$x$	$x$	$x$	$x$	$x$
	1	1	1	1
	2	2	2	2
		<del>3</del>	<del>3</del>	9

... would never converge

□  $Q$  is non-monotone b/c as we added 3, 3 got deleted, as we added 6, 6 got deleted etc.

□ That's why aggr. not allowed in recursive queries in SQL standard .

A 2<sup>nd</sup> example after we cover “mutual recursion” (stay tuned).

# Linear vs Non-linear Recursion

---

□ Recall  $Q_R$  in transitive closure:

```
SELECT Ancestor.anc, Adv.stu  
      FROM Ancestor, Advisor  
      WHERE Ancestor.desc = Advisor.sup
```

Has 1 reference to itself Ancestor: Called *linear recursion*

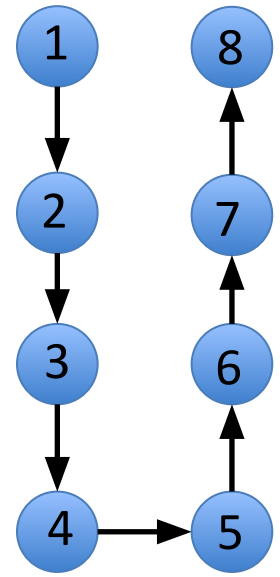
Can have > 1 reference to Ancestor, called *non-linear recursion*

# Non-linear Recursive Computation of Ancestors

WITH RECURSIVE **Ancestors(anc, desc)** AS (

```
SELECT sup as anc, stu as desc
FROM Advisor
UNION
SELECT Anc1.anc, Anct.desc
FROM Ancestor Anc1, Ancestor Anc2
WHERE Ac1.desc = Anc2.anc)
```

Advisor	
<u>sup</u>	<u>stu</u>
1	2
2	3
3	4
4	5
5	6
6	7
7	8



Anc <sub>0</sub>	
<u>anc</u>	<u>desc</u>

Q <sub>B</sub>	
<u>anc</u>	<u>desc</u>
1	2
2	3
3	4
4	5
5	6
6	7
7	8

Q <sub>R</sub>	
<u>anc</u>	<u>desc</u>

Anc <sub>1</sub>	
<u>anc</u>	<u>desc</u>
1	2
2	3
3	4
4	5
5	6
6	7
7	8

Ancestor<sub>1</sub> =                      U                      =

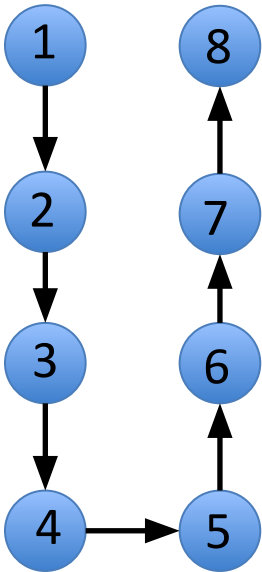
Is Anc<sub>1</sub> - Ans<sub>0</sub> = ∅?  
No: Repeat

# Non-linear Recursive Computation of Ancestors

Anc <sub>1</sub>	
<u>anc</u>	<u>desc</u>
1	2
2	3
3	4
4	5
5	6
6	7
7	8

→ All 1-degree  
ancestors

Advisor	
<u>sup</u>	<u>stu</u>
1	2
2	3
3	4
4	5
5	6
6	7
7	8



Q <sub>B</sub>	
<u>anc</u>	<u>desc</u>
1	2
2	3
3	4
4	5
5	6
6	7
7	8

∪

Q <sub>R</sub>	
<u>anc</u>	<u>desc</u>
1	3
2	4
3	5
4	6
5	7
6	8

=

Anc <sub>2</sub>			
<u>anc</u>	<u>desc</u>	<u>anc</u>	<u>desc</u>
1	2	1	3
2	3	2	4
3	4	3	5
4	5	4	6
5	6	5	7
6	7	6	8
7	8		

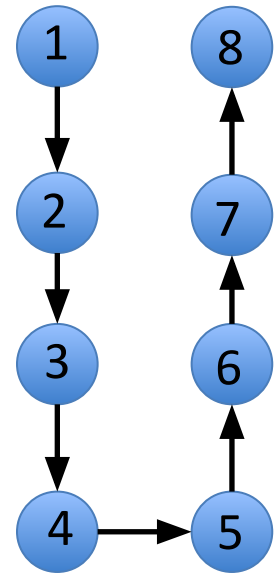


# Non-linear Recursive Computation of Ancestors

Anc <sub>2</sub>			
<u>anc</u>	<u>desc</u>	<u>anc</u>	<u>desc</u>
1	2	1	3
2	3	2	4
3	4	3	5
4	5	4	6
5	6	5	7
6	7	6	8
7	8		

→ All 1 and 2-degree ancestors

Advisor	
<u>sup</u>	<u>stu</u>
1	2
2	3
3	4
4	5
5	6
6	7
7	8



Q <sub>B</sub>	
<u>anc</u>	<u>desc</u>
1	2
2	3
3	4
4	5
5	6
6	7
7	8

U

Q <sub>R</sub>					
<u>anc</u>	<u>desc</u>	<u>anc</u>	<u>desc</u>	<u>anc</u>	<u>desc</u>
1	3	2	5	4	8
2	4	3	6		
3	5	4	7		
4	6	5	8		
5	7	1	5		
6	8	2	6		
1	4	3	7		

=

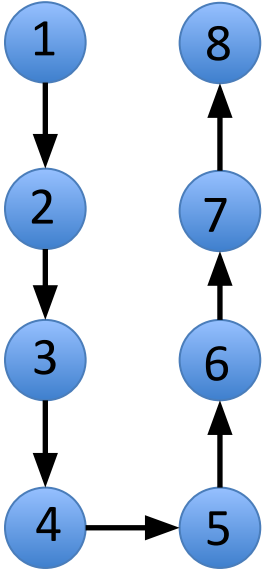
Anc <sub>3</sub>							
<u>a</u>	<u>d</u>	<u>a</u>	<u>d</u>	<u>a</u>	<u>d</u>	<u>a</u>	<u>d</u>
1	2	1	3	2	5	4	8
2	3	2	4	3	6		
3	4	3	5	4	7		
4	5	4	6	5	8		
5	6	5	7	1	5		
6	7	6	8	2	6		
7	8	1	4	3	7		

# Non-linear Recursive Computation of Ancestors

Anc <sub>3</sub>							
<u>a</u>	<u>d</u>	<u>a</u>	<u>d</u>	<u>a</u>	<u>d</u>	<u>a</u>	<u>d</u>
1	2	1	3	2	5	4	8
2	3	2	4	3	6		
3	4	3	5	4	7		
4	5	4	6	5	8		
5	6	5	7	1	5		
6	7	6	8	2	6		
7	8	1	4	3	7		

All 1, 2, 3, and 4-degree  
ancestors

Advisor	
<u>sup</u>	<u>stu</u>
1	2
2	3
3	4
4	5
5	6
6	7
7	8



Q <sub>B</sub>	
<u>anc</u>	<u>desc</u>
1	2
2	3
3	4
4	5
5	6
6	7
7	8

U

Q <sub>R</sub>					
<u>a</u>	<u>d</u>	<u>a</u>	<u>d</u>	<u>a</u>	<u>d</u>
1	3	2	5	4	8
2	4	3	6	1	6
3	5	4	7	2	7
4	6	5	8	3	8
5	7	1	5	1	7
6	8	2	6	2	8
1	4	3	7	1	8

=

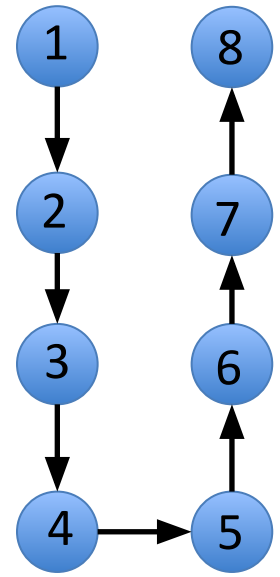
Anc <sub>4</sub>							
<u>a</u>	<u>d</u>	<u>a</u>	<u>d</u>	<u>a</u>	<u>d</u>	<u>a</u>	<u>d</u>
1	2	1	3	2	5	4	8
2	3	2	4	3	6	1	6
3	4	3	5	4	7	2	7
4	5	4	6	5	8	3	8
5	6	5	7	1	5	1	7
6	7	6	8	2	6	2	8
7	8	1	4	3	7	1	8

# Non-linear Recursive Computation of Ancestors

Anc <sub>4</sub>							
<u>a</u>	<u>d</u>	<u>a</u>	<u>d</u>	<u>a</u>	<u>d</u>	<u>a</u>	<u>d</u>
1	2	1	3	2	5	4	8
2	3	2	4	3	6	1	6
3	4	3	5	4	7	2	7
4	5	4	6	5	8	3	8
5	6	5	7	1	5	1	7
6	7	6	8	2	6	2	8
7	8	1	4	3	7	1	8

All 1, ... 8-degree  
ancestors

Advisor	
<u>sup</u>	<u>stu</u>
1	2
2	3
3	4
4	5
5	6
6	7
7	8



Q <sub>B</sub>	
<u>anc</u>	<u>desc</u>
1	2
2	3
3	4
4	5
5	6
6	7
7	8

U

Q <sub>R</sub>					
<u>a</u>	<u>d</u>	<u>a</u>	<u>d</u>	<u>a</u>	<u>d</u>
1	3	2	5	4	8
2	4	3	6	1	6
3	5	4	7	2	7
4	6	5	8	3	8
5	7	1	5	1	7
6	8	2	6	2	8
1	4	3	7	1	8

=

Anc <sub>5</sub>							
<u>a</u>	<u>d</u>	<u>a</u>	<u>d</u>	<u>a</u>	<u>d</u>	<u>a</u>	<u>d</u>
1	2	1	3	2	5	4	8
2	3	2	4	3	6	1	6
3	4	3	5	4	7	2	7
4	5	4	6	5	8	3	8
5	6	5	7	1	5	1	7
6	7	6	8	2	6	2	8
7	8	1	4	3	7	1	8

Is  $Anc_5 - Ans_4 = \emptyset$ ?  
Yes: Stop  
Fixed Point  
Linear recursion  
would take 8  
steps

# Linear vs Non-linear Recursion

---

- For tc-like computations:
  - Linear recursion:
    - Takes \*linear\* # iterations in the depth of the relationships
    - But each iteration might perform less work b/c joins are between smaller tables
  - Non-linear recursion:
    - Takes logarithmic # iterations in the same depth
    - But each iteration performs more work
- SQL standard requires/allows linear recursion for performance reasons

# Mutual Recursion

- Each  $Q_i$  in our examples so far referred to itself.
- We can have the following “mutually recursive” set of queries

WITH RECURSIVE

RECURSIVE  $R_1$  AS  $Q_1$   $\longrightarrow$  e.g. references  $R_2$

RECURSIVE  $R_2$  AS  $Q_2$   $\longrightarrow$  e.g. references  $R_3$

RECURSIVE  $R_3$  AS  $Q_3$   $\longrightarrow$  e.g. references  $R_1$

...

- No query alone is recursive but  $Q_1, Q_2, Q_3$  together is recursive
- So they need to be executed “in tandem” until fixed point.

# Mutual Recursion Example

□ Table *Natural* (*n*) contains 1

□ Even/Odd numbers < 100

WITH RECURSIVE

```
    Even(n) AS (SELECT n FROM Natural
                WHERE n = ANY(SELECT n+1 FROM Odd) AND n < 100),
    Odd(n) AS (
        (SELECT n FROM Natural WHERE n = 1)
        UNION
        (SELECT n FROM Natural
         WHERE n = ANY(SELECT n+1 FROM Even) AND n < 100)
```

Even<sub>0</sub> = ∅,      Odd<sub>0</sub> = ∅

Even<sub>1</sub> = ∅,      Odd<sub>1</sub> = {1}

Even<sub>2</sub> = {2},    Odd<sub>2</sub> = {1}

Even<sub>3</sub> = {2},    Odd<sub>3</sub> = {1, 3}

Even<sub>4</sub> = {2, 4}, Odd<sub>4</sub> = {1, 3}

Even<sub>5</sub> = {2, 4}, Odd<sub>5</sub> = {1, 3, 5}

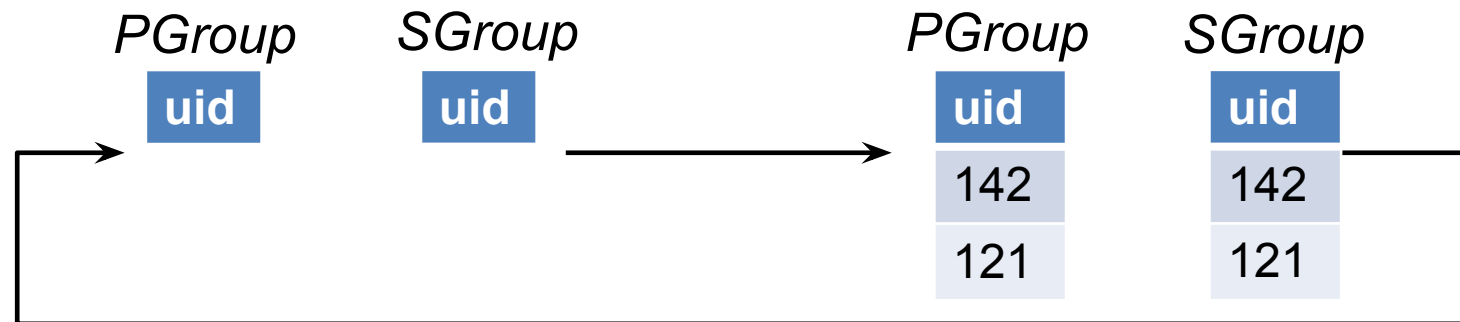
...

# Example Non-Monotone Recursive Query 2: Set Difference

```
WITH RECURSIVE PGroup(uid) AS
  (SELECT uid FROM User
   AND uid NOT IN (SELECT uid FROM SGroup)),
RECURSIVE SGroup(uid) AS
  (SELECT uid FROM User
   AND uid NOT IN (SELECT uid FROM PGroup))
```

uid	name	age
142	Bart	10
121	Allison	8

*MINUS can replace with AND uid NOT IN. In general negated sub-queries or MINUS in recursive parts are not allowed.*



□ Q is non-monotone b/c recall set diff. is nonmonotone w.r.t 2<sup>nd</sup> arg<sub>31</sub>

# Important Note On Monotonicity/Convergence

---

- In practice: DBMSs will not/cannot check for monotonicity and may allow much more than SQL standard: arithmetic, aggregations.
- Nor will they detect oscillations
- You can write non-converging code. Systems will often run a max # iterations (e.g., 100) and error
- SQL compiler will not error for these errors. This is on the user!
- Be careful with recursive queries: Know your query & database!



# Example Non-convergence Based On Database

□ Consider this All Paths query Q:

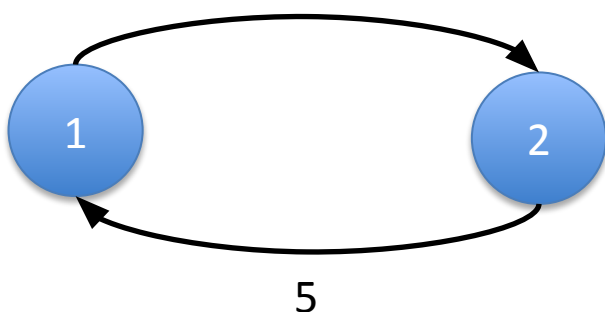
```
WITH RECURSIVE AllPaths(s, d, cost) AS
  (SELECT s, d, cost FROM Edges)
  UNION
  (SELECT AllPaths.s, Edges.d, AllPaths.cost+Edges.cost
   FROM AllPaths, Edges
   WHERE AllPaths.d = Edges.s)
```

□ If Edges { (1, 2, 10) } => All Paths: { (1, 2, 10) }

□ Keep Q the same but add one more tuple (2, 1, 5) to Edges

□ Now there are infinitely many (1, 2) and (2, 1) paths:

(1, 2, 10), (1, 2, 25), (1, 2, 40) etc..



Systems will allow this query!

# Summary of SQL Recursion

---

- Recursion did not exist from 1986-1999 in SQL Standard

- General Syntax: `WITH RECURSIVE`

`R1 AS Q1`

`R2 AS Q2`

`...`

`Rn AS Qn`

- Basic functionality: linear recursion

- Extended functionality: non-linear and mutual recursion

- Unsafe recursive queries: non-monotone (may not converge) queries or query is monotone but output relation's size is infinite (e.g., due to use of arithmetic)

- Personal opinion: Recursive computations are not elegant in SQL. 34

# Datalog: Logic-based DB Query Language with Recursion As a First Class Citizen

- A QL based on logical rules of the form: **Head** := **Body**
- A DB consists of a set of “base relations” (called “extensional” db)

Likes(person, foodItem)

Sells(restaurant, foodItem, cost)

Frequents(person, restaurant)

- Ex Rule: Happy(p) := Likes(p, f), Frequents(p, r), Sells(r, f, c)

Head

Body: conjunction/AND of “subgoals”

- For simplicity: assume head, subgoals can be relation names (called predicates) with arguments that can be variables or constants.
- Datalog allows other predicates: e.g.,  $c < 20$

# Semantics of Datalog Rules

---

- $\text{Happy}(p) := \text{Likes}(p, f), \text{Frequents}(p, r), \text{Sells}(r, f, c)$
- Natural join on common variables:
  - Any  $p$  s.t. "  $\exists$  a food  $f$  & rest  $r$  |  $p$  likes  $f$  &  $p$  frequents  $r$  &  $r$  sells  $f$ "  
is happy
  - In RA:  $\Pi_{\text{person}} (\text{Likes} \bowtie \text{Frequents} \bowtie \text{Sells})$
- Equality filters on constants:
  - $\text{Happy}(p) := \text{Likes}(p, f), \text{Frequents}(p, r), \text{Sells}(r, f, 20)$
  - Any  $p$  s.t. "  $\exists$  a food  $f$  & rest  $r$  |  $p$  likes  $f$  &  $p$  frequents  $r$  &  $r$  sells  $f$   
&  $x$  costs 20 CAD" is happy
- Note: also declarative

# More “Datalog Program” Examples

---

- There can be multiple rules with the same head predicate

`Happy(p) := Likes(p, f), Frequents(p, r), Sells(r, f, c)`

`Happy(p) := Likes(p, “Chocolate Cake”)`

`Happy(p) := Frequents(p, r), Frequents(“Karim”, r)`

- Meaning: Any p s.t:

1) “  $\exists$  a food f & rest r | p likes f & p frequents r & r sells f ” OR

2) “p likes Chocolate Cake” OR

3) “  $\exists$  a restaurant r | both Karim and p frequent r ”

is happy!

- Advantage: Arbitrary recursive, non-recursive, mutually recursive

rules can be just written down as logical “derivation” rules

# More Elegant Recursive Programs

---

## Example 1: Transitive Closure:

`Ancestor(a, d) := Advisor(a, d)`

`Ancestor(a, d) := Ancestor(a, b), Advisor(b, d)`

## Example 2: Shortest Paths:

`AllPaths(a, d, cost) := Edge(a, d, cost)`

`AllPaths(a, d, totalCost) := AllPaths(a, k, cost1), Edge(k, d, cost2),  
totalCost = cost1 + cost2`

`ShortestPaths(a, d, min(cost)) := AllPaths(a, k, cost)`

- Can be done in SQL WITH RECURSIVE but don't need to think about any recursive execution.
- Syntax forces one to focus on logical derivation rules for relations.

# Very Strong and Beautiful Result

- Given a Datalog program that satisfy some properties (specifically some monotonicity and finiteness rules as before):

$R_1 := \text{body 1 (possibly recursive)}$

$R_2 := \text{body 2 (possibly recursive)}$

...

$R_7 := \text{body 7 (possibly recursive)}$

...

$R_k := \text{body 1000 (possibly recursive)}$

- Apply rules in arbitrary order to generate new tuples and one always converges to same unique fixed-point => i.e., the order of execution does not matter

- If you want: run  $R_1 := \text{body 1}$  500 times if it keeps producing new tuples; then run  $R_2 := \text{body 2}$ , then  $R_j$ , then  $R_1$  again etc.

# Last Comments On Datalog

---

- Several DBMSs, e.g., recent LogicBlox or LinkedIn's core graph DBMS, adopts Datalog as a query language instead of SQL
- Better fit for apps requiring recursion and logical inference rules (e.g., in knowledge management and traditional AI applications)  
$$\text{Sibling}(x, y) := \text{BioParent}(z, x), \text{BioParent}(z, y), x \neq y$$
- Many cool applications have been developed on Datalog: (e.g., declarative distributed network programming)
  - See Peter Alvaro's work from UC Santa Cruz
- Has been the foundation for many seminal theoretical results



# Programming Applications W/ SQL

---

- Challenge of using SQL on a real app:
  - Not intended for general-purpose computation
- Solutions
  - Augment SQL with constructs from general-purpose programming languages
    - E.g.: SQL/PSM
  - Use SQL together with general-purpose programming languages: many possibilities
    - Through an API, e.g., Python psycopg2
    - Embedded SQL, e.g., in C
  - SQL generating approaches: Web Programming Frameworks (e.g., Django)

# 1) Augmenting SQL: SQL/PSM

---

- An ISO standard to extend SQL to an advanced prog. lang.
  - Control flow, exception handling, etc.
- Several systems adopt SQL/PSM partially (e.g. MySQL, PostgreSQL)
- PSM = **P**ersistent **S**tored **M**odules
- **CREATE PROCEDURE** *proc\_name*(*param\_decls*)  
    *local\_decls*  
    *proc\_body*;
- **CREATE FUNCTION** *func\_name*(*param\_decls*)  
    **RETURNS** *return\_type*  
    *local\_decls*  
    *func\_body*;
- **CALL** *proc\_name*(*params*);
- Inside procedure body:  
    **SET** *variable* = **CALL** *func\_name*(*params*);

# SQL/PSM Example

```
CREATE FUNCTION SetMaxPop(IN newMaxPop FLOAT)
RETURNS INT
```

```
-- Enforce newMaxPop; return # rows modified.
```

```
BEGIN
```

```
DECLARE rowsUpdated INT DEFAULT 0;
```

```
DECLARE thisPop FLOAT;
```

```
-- A cursor to range over all users:
```

```
DECLARE userCursor CURSOR FOR
```

```
  SELECT pop FROM User
```

```
FOR UPDATE;
```

```
-- Set a flag upon "not found" exception:
```

```
DECLARE noMoreRows INT DEFAULT 0;
```

```
DECLARE CONTINUE HANDLER FOR NOT FOUND
```

```
  SET noMoreRows = 1;
```

```
... (see next slide) ...
```

```
  RETURN rowsUpdated;
```

```
END
```

Declare  
local  
variables

# SQL/PSM Example

```
-- Fetch the first result row:
OPEN userCursor;
FETCH FROM userCursor INTO thisPop;
-- Loop over all result rows:
WHILE noMoreRows <> 1 DO
  IF thisPop > newMaxPop THEN
    -- Enforce newMaxPop:
    UPDATE User SET pop = newMaxPop
    WHERE CURRENT OF userCursor;
    -- Update count:
    SET rowsUpdated = rowsUpdated + 1;
  END IF;
  -- Fetch the next result row:
  FETCH FROM userCursor INTO thisPop;
END WHILE;
CLOSE userCursor;
```

Function  
body

# Other SQL/PSM Features

---

- Assignment using scalar query results
  - SELECT INTO
- Other loop constructs
  - FOR, REPEAT UNTIL, LOOP
- Flow control
  - GOTO
- Exceptions
  - SIGNAL, RESIGNAL
- ...
- For more PostgreSQL-specific information, look for “PL/pgSQL” in PostgreSQL documentation

## 2) Working with SQL through an API

---

- E.g.: Python psycopg2, JDBC, ODBC (C/C++/VB)
  - Based on the SQL/CLI (Call-Level Interface) standard
- The application program sends SQL commands to the DBMS at runtime. Gets back a “cursor” that can iterate over results.
- Results are converted to objects in the application program.  
Often you use a cursor to loop through result tuples.

# Example API: Python psycopg2

```
import psycopg2
conn = psycopg2.connect(dbname='beers')
cur = conn.cursor()
# list all drinkers:
cur.execute('SELECT * FROM Drinker')
for drinker, address in cur:
    print(drinker + ' lives at ' + address)
# print menu for bars whose name contains "a":
cur.execute('SELECT * FROM Serves WHERE bar LIKE %s',
            ('%a%',))
for bar, beer, price in cur:
    print('{} serves {} at ${:,.2f}'.format(bar, beer, price))
cur.close()
conn.close()
```

You can iterate over cur  
one tuple at a time

Placeholder for  
query parameter

Tuple of parameter values,  
one for each %s

# More psycpg2 Examples

```
# “commit” each change immediately—need to set this option just once at  
the start of the session
```

```
conn.set_session(autocommit=True)
```

```
# ...
```

```
bar = input('Enter the bar to update: ').strip()
```

```
beer = input('Enter the beer to update: ').strip()
```

```
price = float(input('Enter the new price: '))
```

```
try:
```

```
    cur.execute("""
```

```
        UPDATE Serves
```

```
        SET price = %s
```

```
        WHERE bar = %s AND beer = %s""", (price, bar, beer))
```

```
    if cur.rowcount != 1:
```

```
        print('{} row(s) updated: correct bar/beer?\'
```

```
            .format(cur.rowcount))
```

```
except Exception as e:
```

```
    print(e)
```

Perform passing,  
semantic analysis,  
optimization,  
compilation, and finally  
execution



# More psycopg2 Examples

```
....  
while true:  
# Input bar, beer, price...  
    cur.execute("""  
        UPDATE Serves  
        SET price = %s  
        WHERE bar = %s AND beer = %s""", (price, bar, beer))  
....  
# Check result...
```

Perform passing,  
semantic analysis,  
optimization,  
compilation, and finally  
execution

Execute many times  
Can we reduce this overhead?

# Prepared Statements: Example

```
cur.execute("""          # Prepare once (in SQL) Prepare only once
    PREPARE update_price AS      # Name the prepared plan,
    UPDATE Serves
    SET price = $1              # and note the $1, $2, ... notation for
    WHERE bar = $2 AND beer = $3""") # parameter placeholders.
```

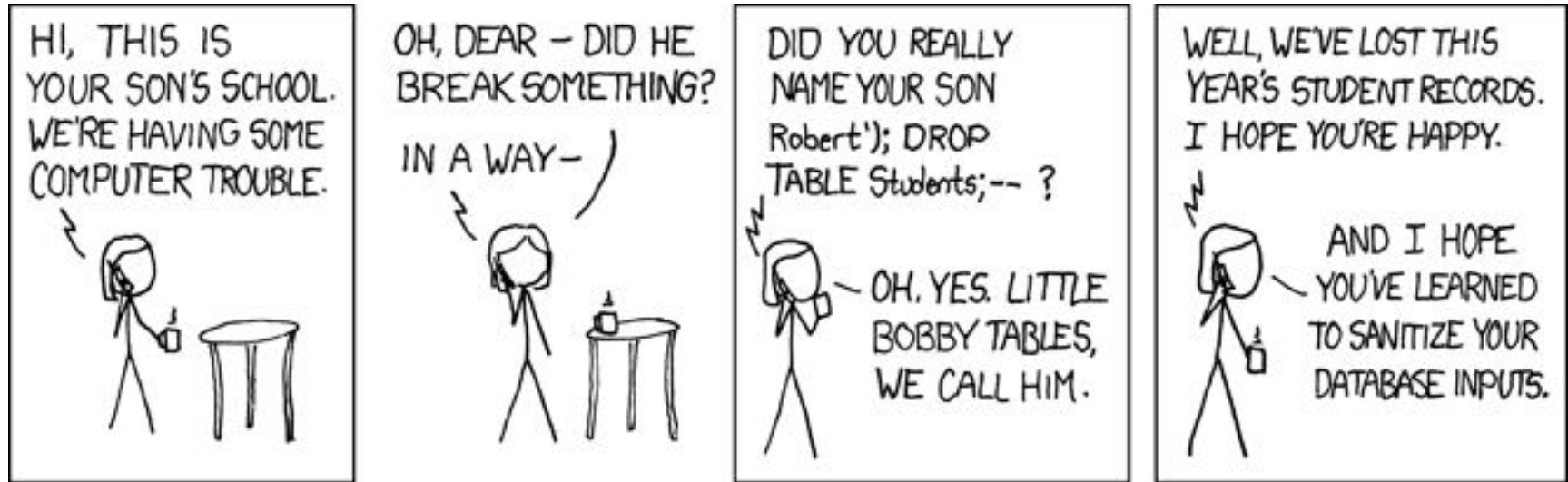
while true:

# Input bar, beer, price...

```
cur.execute('
    EXECUTE update_price(%s, %s, %s)',\ # Execute many times.
    (price, bar, beer))....
```

# Check result...

# Watch Out For SQL Injection Attacks!



<http://xkcd.com/327/>

- The school probably had something like:

```
cur.execute("SELECT * FROM Students " +  
  \  
  "WHERE (name = " + name + ")")
```

where **name** is a string input by user

- Called an SQL injection attack. Most APIs have ways to sanitize inputs.

# Augmenting SQL vs. Programming Through an API

---

## □ Pros of augmenting SQL:

- More processing features for DBMS
- More application logic can be pushed closer to data

## □ Cons of augmenting SQL:

- SQL is already too big
- Complicate optimization and make it impossible to guarantee safety

### 3) “Embedding” SQL in a host language

---

- Can be thought of as the opposite of SQL/PSM
- Extends a host language, e.g., C or Java, with SQL-based features
- Can compile host language together with SQL statements and catch SQL errors during *application compilation time*

# 4) Web Programming Frameworks

- A web development “framework” e.g., Django or Ruby on Rails
- Very frequent approach to web apps that need a DB
- For most parts, no explicitly writing SQL is needed:
- Example: Django Web App Programming:
  - Define “Models”: python objects and only do oo programming
  - Models will be backed up with Relations in an RDBMS
- E.g.: a Person class/object with first and lastName:

```
from django.db import models

class Person(models.Model):
    f_name = models.CharField(max_len=30)
    l_name = models.CharField(max_len=30)
```

```
CREATE TABLE myapp_person (
    "id" serial NOT NULL PRIMARY KEY,
    "f_name" varchar(30) NOT NULL,
    "l_name" varchar(30) NOT NULL );
```

- Would lead the “framework” (not the user) to generate the following SQL code somewhere in the web application files:

Thank You