# Database Theory

Conjunctive Queries



#### Motivation

We've seen that we are limited in many things we want to do when it comes to powerful languages like FO/RA.

Let us instead study a restricted subclass of queries that lies at the core of important data retrieval tasks.

 $\{\ \bar{y}\ |\ \exists \bar{z}\ R_1(\bar{x}_1) \land R_2(\bar{x}_2) \land \cdots \land R_k(\bar{x}_k)\ \}$ 

We call queries of this form Conjunctive Queries (CQs).

That is, conjunctive queries are FO queries using only the connectives  $\exists$  and  $\land$ .

$$\pi_{\bar{y}}\left(R_1\bowtie R_2\bowtie\cdots\bowtie R_k\right)$$

(Assuming attributes for  $R_i$  are  $ar{x}_{i'}$  otherwise simply rename)

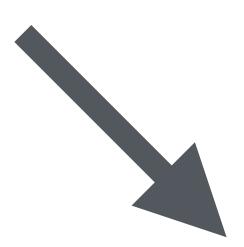
$$\left\{ \begin{array}{c|c} \bar{y} & \exists \bar{z} \ R_1(\bar{x}_1) \land R_2(\bar{x}_2) \land \cdots \land R_k(\bar{x}_k) \end{array} \right\}$$

$$\left\{ \begin{array}{c|c} \bar{y} & \exists \bar{z} \ R_1(\bar{x}_1) \land R_2(\bar{x}_2) \land \cdots \land R_k(\bar{x}_k) \end{array} \right\}$$

$$\pi_{\bar{y}} \left( R_1 \bowtie R_2 \bowtie \cdots \bowtie R_k \right)$$

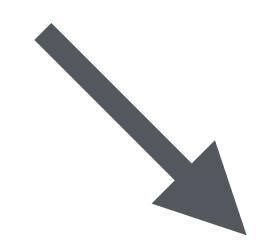
Conjunctive queries correspond so-called join queries in RA. That is, RA queries that only use projection, renaming, selection, and joins.

Recall our SQL example in the lecture on the relational model.



 $\{(c,s) \mid \exists sid, l, dob, active . Enrolled(c, WS24, sid) \land Course(c, WS24, l) \land Student(sid, s, dob, active)\}$ 

CQs cover the core part of most SQL queries!



 $\{(c,s) \mid \exists sid, l, dob, active . Enrolled(c, WS24, sid) \land Course(c, WS24, l) \land Student(sid, s, dob, active)\}$ 

◆ Conjunctive queries form the key part of most data retrieval tasks.

- → Join queries
- → Datalog rule bodies are CQs
- Basis for many other query languages,
   e.g., Conjunctive Regular Path Queries.

- ◆ Conjunctive queries form the key part of most data retrieval tasks.
- Optimising CQs can help to optimise the most expensive part of practical join evaluations

```
select min(ps_supplycost)
from

    partsupp,
    supplier,
    nation,
    region

where

    p_partkey = ps_partkey
    and s_suppkey = ps_suppkey
    and s_nationkey = n_nationkey
    and n_regionkey = r_regionkey
    and r_name = 'Asia'
```

Real systems will evaluate a CQ first and then evaluate the min aggregate on the result of the CQ.

- ◆ Conjunctive queries form the key part of most data retrieval tasks.
- Optimising CQs can help to optimise the most expensive part of practical join evaluations
- ◆ Complexity for results for CQs also give us lower bounds for more complex queries that often have CQs at their core.

```
select min(ps_supplycost)
from

    partsupp,
    supplier,
    nation,
    region

where

    p_partkey = ps_partkey
    and s_suppkey = ps_suppkey
    and s_nationkey = n_nationkey
    and n_regionkey = r_regionkey
    and r_name = 'Asia'
```

The query intuitively will be at least as hard to solve as the underlying CQ (without the min aggregate).

# Equivalence & Containment

# Query Containment

- lacktriangle For queries  $q_1,q_2$  we say that  $q_1$  is contained in  $q_2$  (in symbols,  $q_1\subseteq q_2$ ) if  $q_1(D)\subseteq q_2(D)$  for every database D.
- ◆ Equivalence of two queries  $q_1, q_2$  (in symbols,  $q_1 \equiv q_2$ ) is defined as  $q_1 \equiv q_2 \iff q_1 \subseteq q_2$  and  $q_2 \subseteq q_1$

See the Trakthenbrot exercise sheet!

Given  $q_1, q_2$  in RA, there is no algorithm to decide  $q_1 \subseteq q_2$ . But if  $q_1, q_2$  are CQs then the problem is decidable!

# Query Containment Example

Consider the following two queries

$$q_1 := \{(x, y) \mid \exists z \, R(y, x) \land R(x, z)\}$$
  
 $q_2 := \{(x, y) \mid R(y, x) \land R(x, y)\}$ 

Intuitively it seems clear that  $q_1$  describes a weaker requirement: x only needs to reach some z, whereas it needs to reach specifically y in  $q_2$ .

We would therefore suspect that  $q_2(D) \subseteq q_1(D)$  for all D. But how to prove it?

R

A	В
1	2
3	3
2	3

$$q_1(D) = \{ (2,1), (3,2), (3,3) \}$$
  
 $q_2(D) = \{ (3,3) \}$ 

#### The Tableau of a CQ

Basic Idea: we can represent a CQ as a database.

The tableau  $\mathsf{Tbl}(q)$  of a CQ q is the database where the tuples of relation R are all of the term lists that occur for R in the query (+ a relation for the output variables).

Consider the following query:

$$\{ (x,y) \mid \exists wz \ B(x,y) \land R(y,z) \land R(y,w) \land R(w,y) \}$$

We write  $\mathsf{Tbl}^*(q)$  for the tables without the special  $\mathsf{Out}$  relation for output variables.

Out

1	2
X	У

B

1	2
X	У

R

1	2
У	Z
У	W
W	У

# Homomorphisms

A homomorphism of two databases  $D_1, D_2$  is a function

 $h:Dom(D_1)\to Dom(D_2)$  such that:

$$(c_1, ..., c_n) \in R^{D_1} \implies (h(c_1), ..., h(c_n)) \in R^{D_2} \quad \forall R \in Rel$$

A homomorphism of two databases  $D_1, D_2$  is a function

 $h:Dom(D_1)\to Dom(D_2)$  such that:

$$(c_1, ..., c_n) \in R^{D_1} \implies (h(c_1), ..., h(c_n)) \in R^{D_2}$$

 $\forall R \in \text{Rel}$ 

Name	Age
Anna	54
Ben	26

**Parent** 

Ben

Child

Anna

R

Name	Age
David	85
Anna	40
Claire	34

Parent	Child
David	David
Anna	Claire

Parent	Child
David	David
Anna	Claire

A homomorphism of two databases  $D_1, D_2$  is a function

 $h:Dom(D_1)\to Dom(D_2)$  such that:

$$(c_1, ..., c_n) \in \mathbb{R}^{D_1} \implies (h(c_1), ..., h(c_n)) \in \mathbb{R}^{D_2}$$

 $\forall R \in \text{Rel}$ 

R

Name	Age
Anna	54
Ben	26

 $\begin{array}{c} Anna \mapsto Claire \\ Ben \mapsto Anna \\ 54 \mapsto 34 \\ 26 \mapsto 40 \end{array}$ 

R

Name	Age
David	85
Anna	40
Claire	34

5

Parent	Child
David	David
Anna	Claire

S

Parent	Child
Ben	Anna

May seem weird in meaning but the "structure" is preserved!

A homomorphism of two databases  $D_1, D_2$  is a function

 $h:Dom(D_1)\to Dom(D_2)$  such that:

$$(c_1, ..., c_n) \in \mathbb{R}^{D_1} \implies (h(c_1), ..., h(c_n)) \in \mathbb{R}^{D_2}$$

 $\forall R \in \text{Rel}$ 

Name	Age
Anna	54
Ben	26

 $Anna \mapsto David$  $Ben \mapsto David$  $54 \mapsto 85$  $26 \mapsto 85$ 

R

Name	Age
David	85
Anna	40
Claire	34

Nothing says the
function must be
injective

Child **Parent** David David Claire Anna

Parent	Child
Ben	Anna

A homomorphism of two databases  $D_1,D_2$  is a function

 $h:Dom(D_1)\to Dom(D_2)$  such that:

$$(c_1, ..., c_n) \in \mathbb{R}^{D_1} \implies (h(c_1), ..., h(c_n)) \in \mathbb{R}^{D_2}$$

 $\forall R \in \text{Rel}$ 

Name	Age
Anna	54
Ben	26

**Parent** 

Ben

Child

Anna

#### Not a hom!

Anna	$\mapsto$	Cla	ire
Ben	$\mapsto$	Cla	ire
	54	1 →	34
	26	$i \mapsto$	34

Name	Age
David	85
Anna	40
Claire	34

Child

David

Claire

at doon't	Parent
	David
ays work	Anna

But that doesn't	
always work	

# Theorem Let $q_1, q_2$ be CQs. Then $q_1 \subseteq q_2 \quad \text{if and only if} \quad \mathsf{Tbl}(q_2) \xrightarrow{hom} \mathsf{Tbl}(q_1)$

As a result we have an "easy" algorithm for deciding containment for CQs:

Careful!
Only holds for set semantics!

- 1. Compute **Tbl**(·) for both queries (trivial).
- 2. Check if there is a homomorphism (in NP).

Consider the following two queries

$$q_1 := \{(x, y) \mid \exists z \, R(y, x) \land R(x, z)\}$$

$$q_2 := \{(x,y) \mid \exists wu \ R(y,x) \land R(w,x) \land R(x,u)\}$$

 $Tbl(q_1)$ 

Out

1	2
X	У

R

1	2
У	X
X	Z

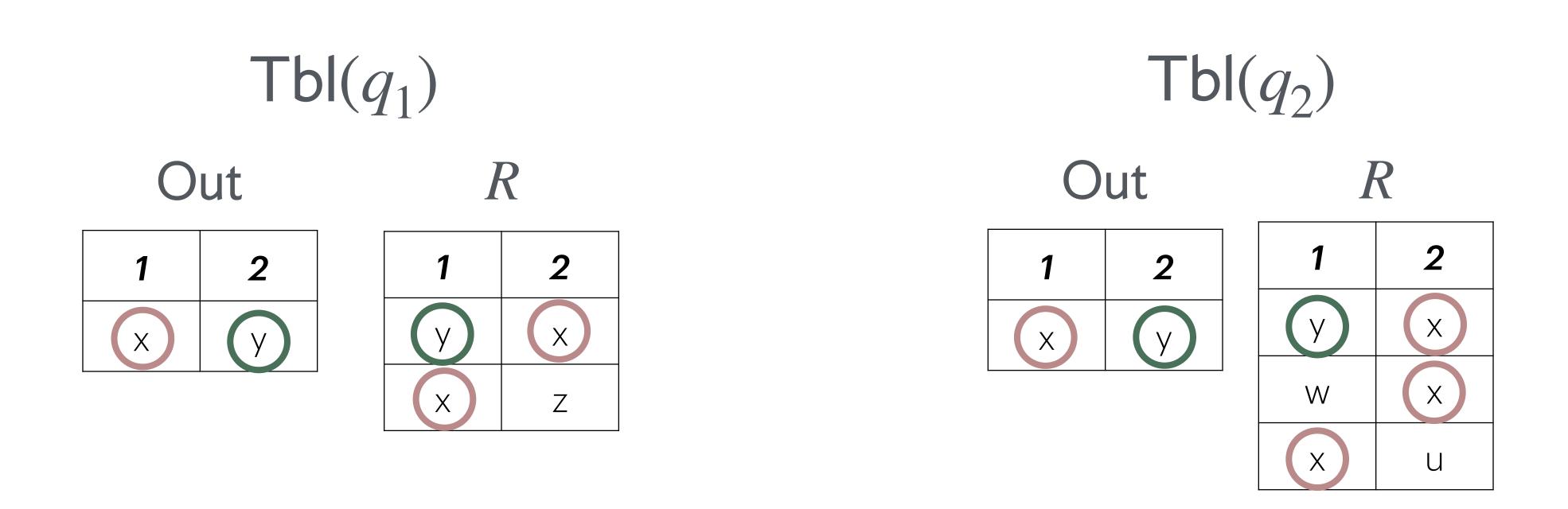
 $\mathsf{Tbl}(q_2)$ 

Out

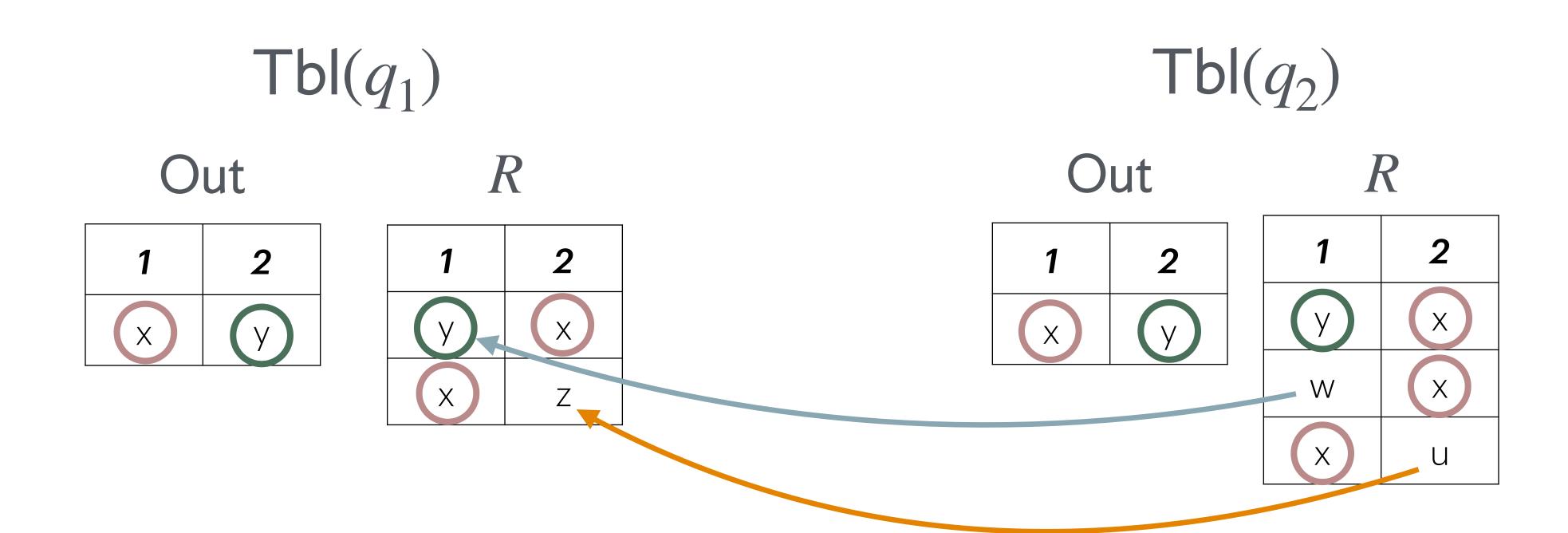
1	2
X	У

R

2
X
X
u

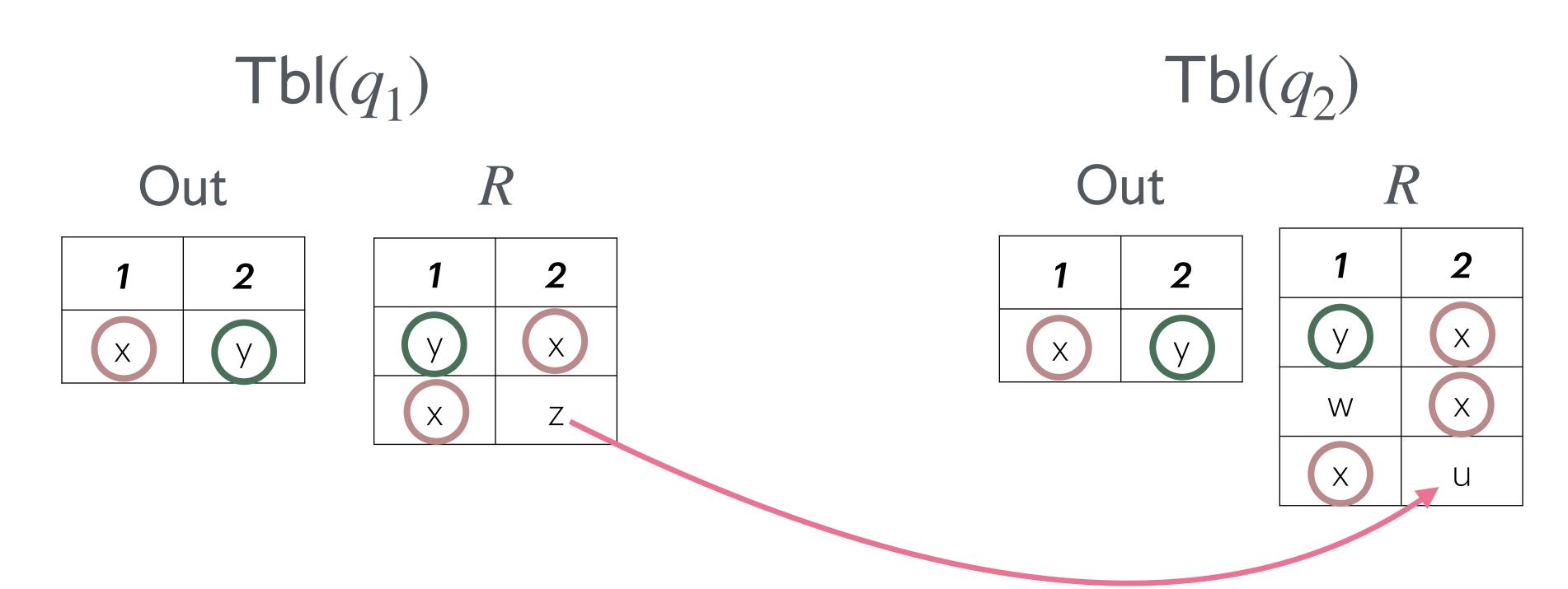


Since Out has only one tuple, a homomorphism h must necessarily have h(x) = x and h(y) = y.



Since  $\operatorname{Out}$  has only one tuple, a homomorphism h must necessarily have h(x) = x and h(y) = y.

With h(w) = y and h(u) = z we get a homomorphism  $\mathsf{Tbl}(q_2) \overset{hom}{\longrightarrow} \mathsf{Tbl}(q_1)$ .



Since Out has only one tuple, a homomorphism h must necessarily have h(x) = x and h(y) = y.

With h(z) = y and h(u) = z we get a homomorphism  $\mathsf{Tbl}(q_1) \xrightarrow{hom} \mathsf{Tbl}(q_2)$ .

Consider the following two queries

$$q_1 := \{(x, y) \mid \exists z R(y, x) \land R(x, z)\}$$

$$q_2 := \{(x, y) \mid \exists wu \ R(y, x) \land R(w, x) \land R(x, u)\}$$

We've seen that  $\mathsf{Tbl}(q_1) \xrightarrow{hom} \mathsf{Tbl}(q_2)$  and  $\mathsf{Tbl}(q_2) \xrightarrow{hom} \mathsf{Tbl}(q_1)$ . By the Homomorphism Theorem that means  $q_1 \subseteq q_2$  and  $q_2 \subseteq q_1$ , i.e.,  $q_1 \equiv q_2$ !



#### Important observation

Let q be a CQ with free variables  $\bar{y}$ . For any database D, we have  $\bar{c} \in q(D)$  if and only if there is a homomorphism h from  $\mathsf{Tbl}^*(q)$  to D such that  $h(\bar{y}) = \bar{c}$ .

$$q_2 := \{(x,y) \mid \exists wu \ R(y,x) \land R(w,x) \land R(x,u)\}$$

 $(a,b) \in q_2(D)$  means that there is an interpretation I such that

1. 
$$R(I(y), I(x)) \in D$$
,  $R(I(w), I(x)) \in D$ , and  $R(I(x), I(u)) \in D$ .

2. and 
$$I(x) = a$$
 and  $I(y) = b$ .

So I is precisely a homomorphism from  $\mathsf{TbI}^*(q_2)$  to D!

If  $q_1 \subseteq q_2$ , then  $\mathsf{Tbl}(q_2) \xrightarrow{hom} \mathsf{Tbl}(q_1)$ 

By assumption  $q_1(D) \subseteq q_2(D)$ . So take  $\mathsf{Tbl}^*(q_1)$  as database D. Let  $\bar{y}$  be the free variables of  $q_1$ . We have that  $(x,y) \in q_1(\mathsf{Tbl}^*(q_1))$  since we can just map every variable to itself.

 $\mathsf{Tbl}^*(q_1) R$ 

1	2
У	X
X	Z

$$h(x) = x, h(y) = y, h(z) = z$$

So 
$$(x, y) \in q_1(\mathsf{Tbl}^*(q_1))$$

 $\mathsf{Tbl}^*(q_1) \; R$ 

1	2
У	X
X	Z

Example queries from before

$$q_1 := \{(x, y) \mid \exists z \, R(y, x) \land R(x, z)\}$$

$$q_2 := \{(x,y) \mid \exists wu \ R(y,x) \land R(w,x) \land R(x,u)$$

If  $q_1 \subseteq q_2$ , then  $\mathsf{Tbl}(q_2) \xrightarrow{hom} \mathsf{Tbl}(q_1)$ 



By assumption  $q_1(D) \subseteq q_2(D)$ . So take  $\mathsf{Tbl}^*(q_1)$  as database D. Let  $\bar{y}$  be the free variables of  $q_1$ . We have that  $(x,y) \in q_1(\mathsf{Tbl}^*(q_1))$  since we can just map every variable to itself.

Then also  $(x, y) \in q_2(\mathsf{Tbl}^*(q_1))$ . By the key observation:

$$\bar{y} \in q_2(\mathsf{Tbl}^*(q_1)) \iff \mathsf{Tbl}^*(q_2) \xrightarrow{hom} \mathsf{Tbl}^*(q_1)$$

That is, the tuple in the **Out** relation of  $\mathsf{Tbl}(q_2)$  maps into a tuple of  $\mathsf{Out}$  of  $\mathsf{Tbl}(q_1)$ . Furthermore, that homomorphism maps the free variables of  $q_2$  to the free variables of  $q_1$ .

We then have  $\mathsf{Tbl}(q_2) \xrightarrow{hom} \mathsf{Tbl}(q_1)$ .

Example queries from before

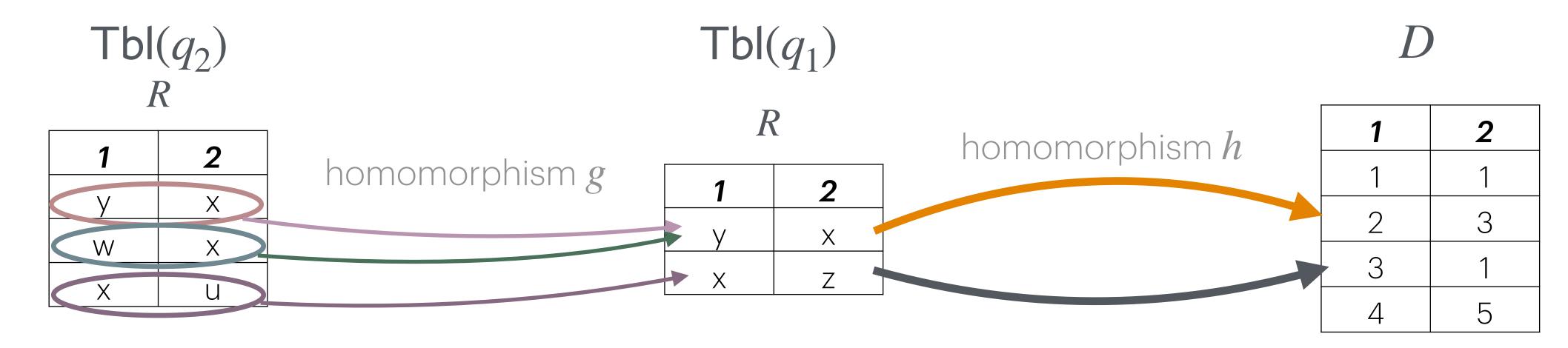
$$q_1 := \{(x, y) \mid \exists z \, R(y, x) \land R(x, z)\}$$

$$q_2 := \{(x, y) \mid \exists wu \, R(y, x) \land R(w, x) \land R(x, u)\}$$

If 
$$\mathsf{Tbl}(q_2) \xrightarrow{hom} \mathsf{Tbl}(q_1)$$
, then  $q_1 \subseteq q_2$ 

If  $\bar{c}$  is an answer of  $q_1$  on some database D then there is a homomorphism  $h: \mathsf{Tbl}^*(q_1) \to D$  that maps the output variables of  $q_1$  to  $\bar{c}$ .

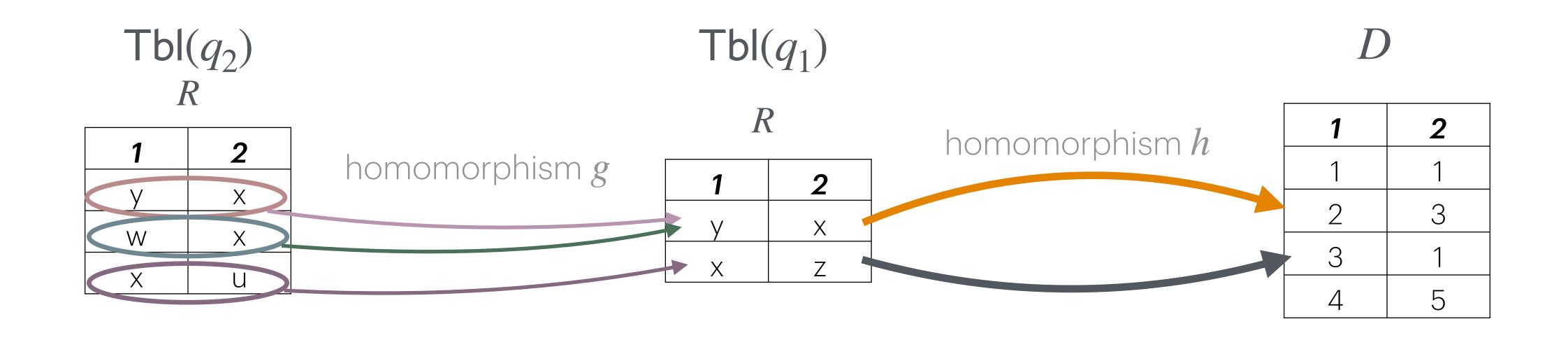
Let g be the homomorphism from  $\mathsf{Tbl}(q_2)$  to  $\mathsf{Tbl}(q_1)$ . It is not hard to see that  $h \circ g$  is homomorphism from  $\mathsf{Tbl}^*(q_2)$  to D that also maps the output of variables of  $q_2$  to  $\bar{c}$ .

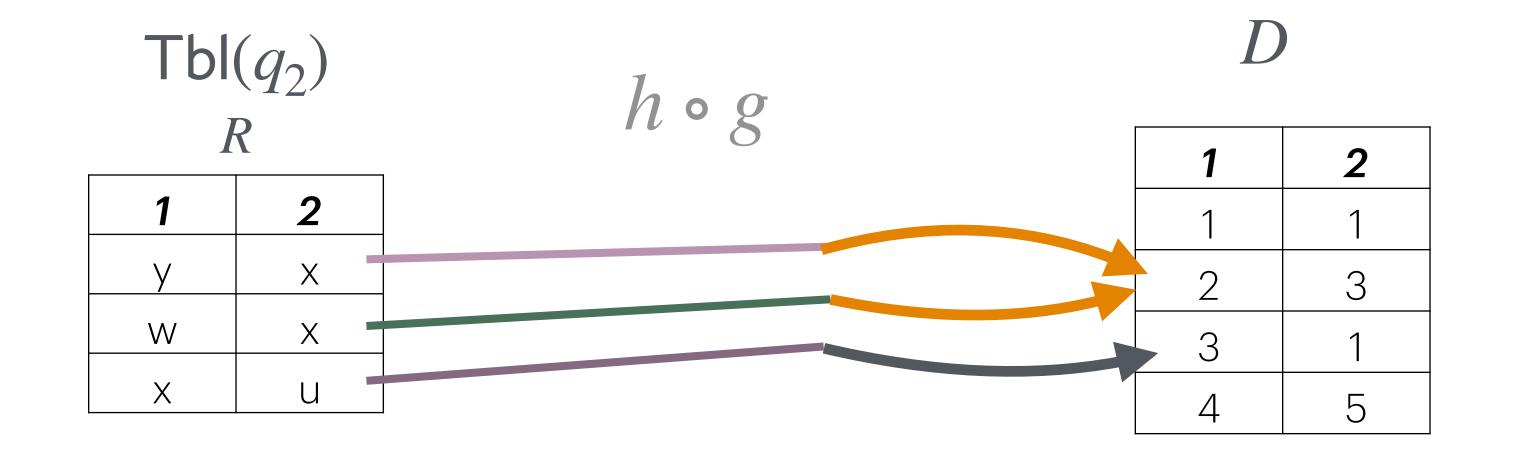


Example queries from before

$$q_1 := \{(x, y) \mid \exists z \, R(y, x) \land R(x, z)\}$$

$$q_2 := \{(x,y) \mid \exists wu \ R(y,x) \land R(w,x) \land R(x,u)\}$$





#### Still a homomorphism!

Output variables map to the same values in D!

Example queries from before

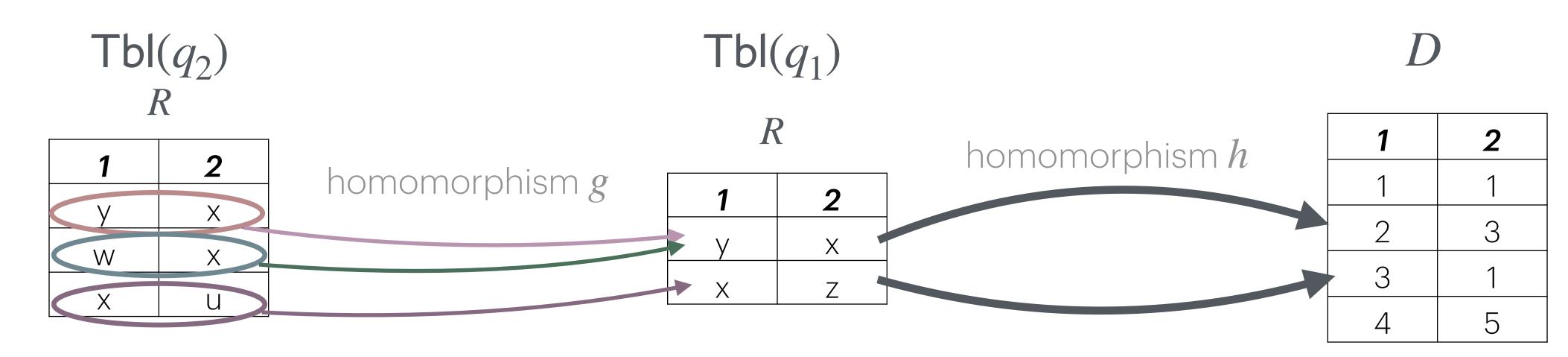
$$q_1 := \{(x, y) \mid \exists z \, R(y, x) \land R(x, z)\}$$

$$q_2 := \{(x, y) \mid \exists wu \, R(y, x) \land R(w, x) \land R(x, u)\}$$

If 
$$\mathsf{Tbl}(q_2) \xrightarrow{hom} \mathsf{Tbl}(q_1)$$
, then  $q_1 \subseteq q_2$ 

If  $\bar{c}$  is an answer of  $q_1$  on some database that maps the output variables of  $q_1$  to  $\bar{c}$ .

Let g be the homomorphism from  $\mathsf{Tbl}(q_2)$  to  $\mathsf{Tbl}(q_1)$ . It is not hard to see that  $h \circ g$  is homomorphism from  $\mathsf{Tbl}^*(q_2)$  to D that also maps the output of variables of  $q_2$  to  $\bar{c}$ .



#### Time for a break.

Turn to your neighbour: briefly discuss the lecture Stretch, grab water, reset

# Query Minimisation

## 

#### Goal:

Given a CQ q, we want the equivalent CQ  $q^\prime$  with the least amount of atoms.

Formally, a CQ q is minimal if there does not exist a CQ  $q^\prime$  such that:

- a)  $q' \equiv q$
- b) q' has fewer atoms (=terms in the conjunction) than q

We would like to replace a CQ with its minimal equivalent CQ before evaluating it.

How do we find this minimal equivalent CQ?

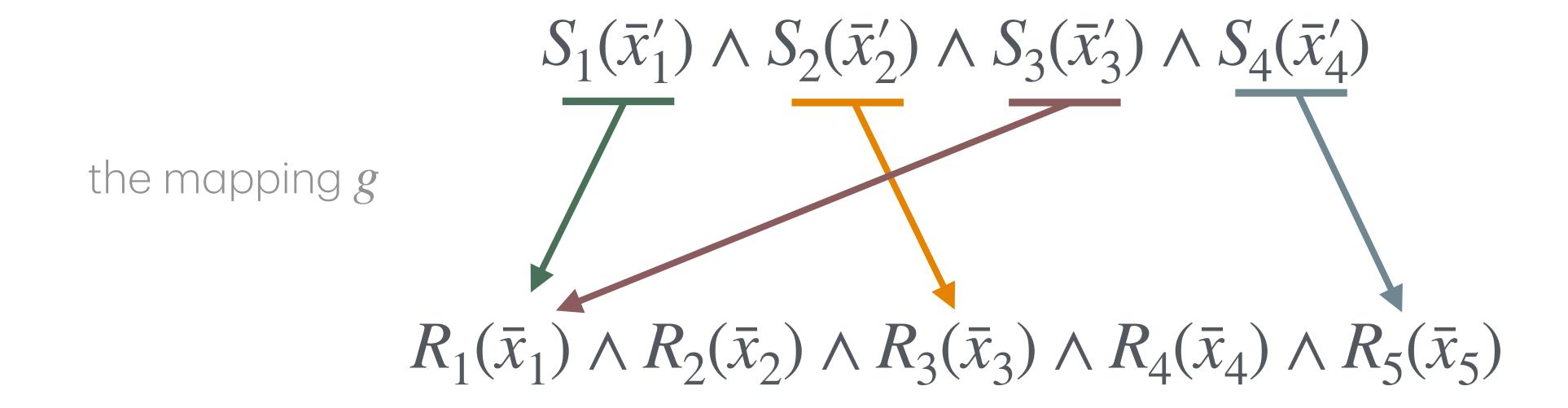
To minimise CQ q, it is enough to check only those queries obtained by deleting atoms from q!

Assume CQ  $q = \{ \ \bar{y} \ | \ \exists \bar{z} \ R_1(\bar{x}_1) \land R_2(\bar{x}_2) \land \cdots \land R_k(\bar{x}_k) \ \}.$  Furthermore, assume q has an equivalent CQ  $q' = \{ \ \bar{y}' \ | \ \exists \bar{z}' \ S_1(\bar{x}_1') \land S_2(\bar{x}_2') \land \cdots \land S_j(\bar{x}_j') \ \} \text{ with } j < k.$ 

By the Homomorphism Theorem there are homomorphisms:

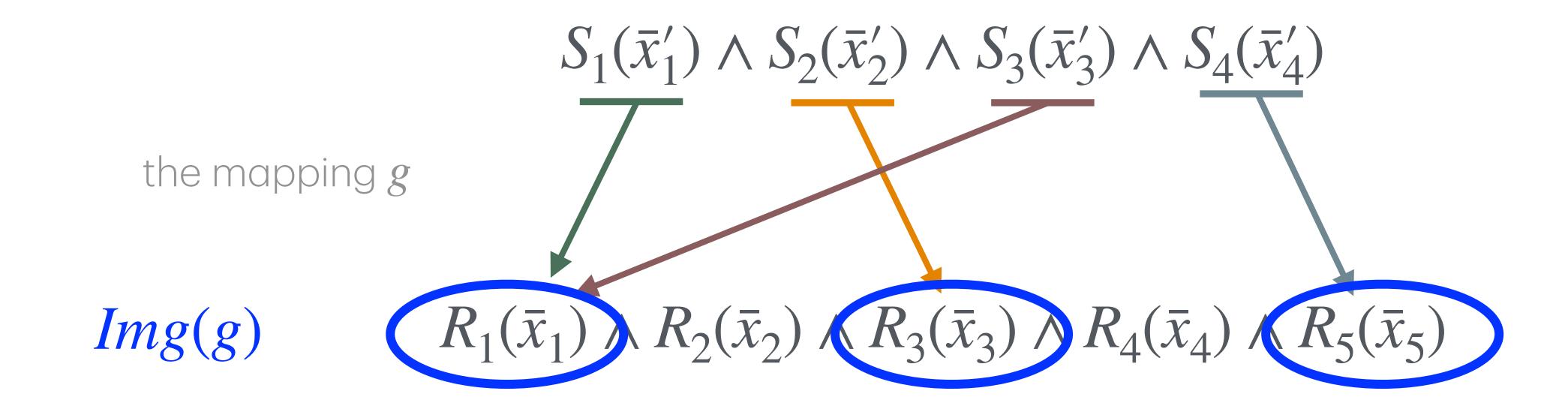
$$f: \mathsf{Tbl}(q) \to \mathsf{Tbl}(q')$$
 and  $g: \mathsf{Tbl}(q') \to \mathsf{Tbl}(q)$ 

We have that  $g: \mathrm{Tbl}(q') \to \mathrm{Tbl}(q)$  maps every  $S_{\alpha}(\bar{x}'_{\alpha}) \in \mathrm{Tbl}(q')$  into some  $R_{i_{\alpha}}(\bar{x}_{i_{\alpha}}) \in \mathrm{Tbl}(q)$  with  $S_{\alpha} = R_{i_{\alpha}}$  and  $\bar{x}_{i_{\alpha}} = h(\bar{x}'_{\alpha})$ .

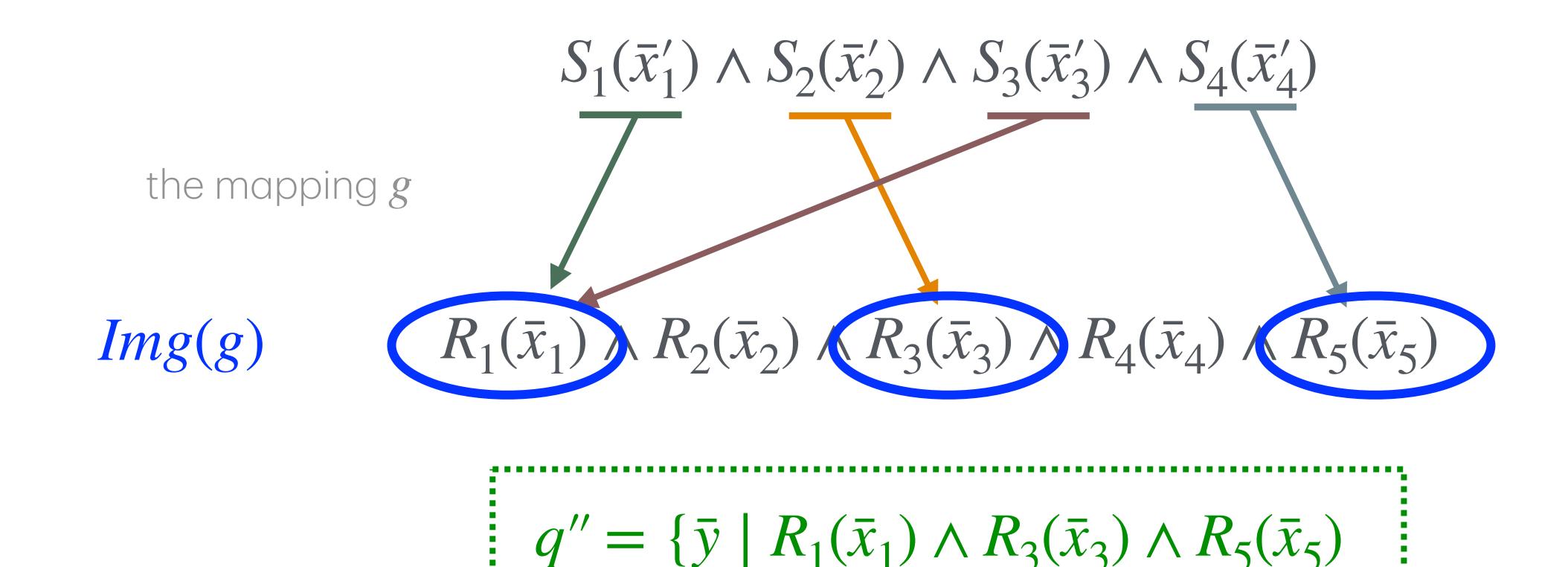


We have that  $g: \mathrm{Tbl}(q') \to \mathrm{Tbl}(q)$  maps every  $S_{\alpha}(\bar{x}'_{\alpha}) \in \mathrm{Tbl}(q')$  into some  $R_{i_{\alpha}}(\bar{x}_{i_{\alpha}}) \in \mathrm{Tbl}(q)$  with  $S_{\alpha} = R_{i_{\alpha}}$  and  $\bar{x}_{i_{\alpha}} = h(\bar{x}'_{\alpha})$ .

Let  $Img(g) = \{R_{i_1}(\bar{x}_{i_1}), R_{i_2}(\bar{x}_{i_2}), \dots, R_{i_\ell}(\bar{x}_{i_\ell})\}$  be the set of all such images of the mapping g applied to the terms of q' and observe that  $|Img(g)| \leq j < k$ .



Let us define the query  $q'' = \{ \bar{y} \mid \exists \bar{z} \; R_{i_1}(\bar{x}_{i_1}) \land R_{i_2}(\bar{x}_{i_2}) \land \cdots \land R_{i_\ell}(\bar{x}_{i_\ell}) \}$  consisting of the terms in Img(g). We see that q'' can be obtained by simply deleting some terms from q.



We use the Homomorphism Theorem to show that q'' is also equivalent to q:

♦ There is a trivial homomorphism  $\mathsf{Tbl}(q'') \to \mathsf{Tbl}(q)$ : simply map every variable to itself.

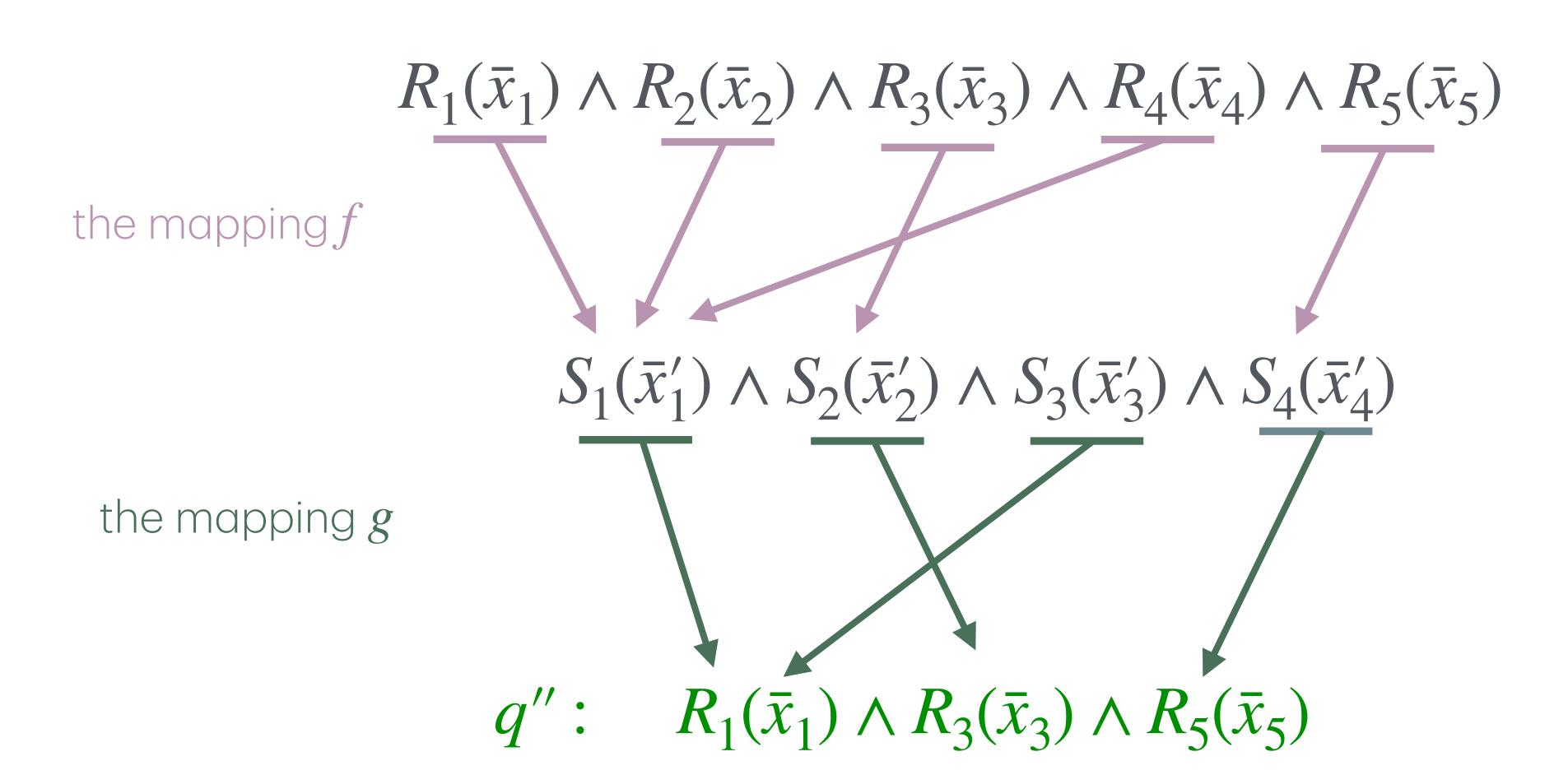
$$q'': R_{1}(\bar{x}_{1}) \wedge R_{3}(\bar{x}_{3}) \wedge R_{5}(\bar{x}_{5})$$

$$q: R_{1}(\bar{x}_{1}) \wedge R_{2}(\bar{x}_{2}) \wedge R_{3}(\bar{x}_{3}) \wedge R_{4}(\bar{x}_{4}) \wedge R_{5}(\bar{x}_{5})$$

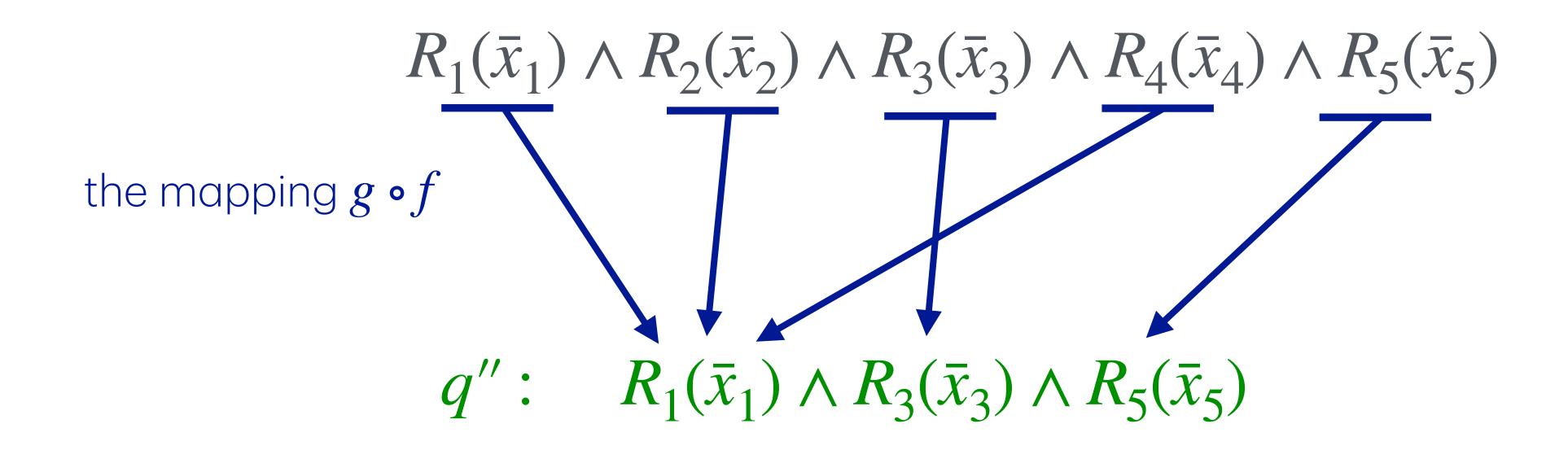
We use the Homomorphism Theorem to show that q'' is also equivalent to q:

- ♦ There is a trivial homomorphism  $\mathsf{Tbl}(q'') \to \mathsf{Tbl}(q)$ : simply map every variable to itself.
- $lacktriangledown g \circ f$  i.e., the function composition  $g(f(\,\cdot\,))$  is a homomorphism  $\mathsf{Tbl}(q) \to \mathsf{Tbl}(q'')$ :
  - function f maps every  $R_i(ar{x}_i)$  in q into an  $S_{lpha}(ar{x}_lpha)$  in q' by definition,
  - function g maps every  $S_{lpha}(ar{x}_{lpha})$  in q' into an atom of q'' by construction.

- ♦  $g \circ f$  i.e., the function composition  $g(f(\cdot))$  is a homomorphism  $\mathsf{Tbl}(q) \to \mathsf{Tbl}(q'')$ :
  - function f maps every  $R_i(ar{x}_i)$  in q into an  $S_{lpha}(ar{x}_lpha)$  in q' by definition,
  - -function g maps every  $S_{lpha}(ar{x}_{lpha})$  in q' into an atom of q'' by construction.



- ♦  $g \circ f$  i.e., the function composition  $g(f(\cdot))$  is a homomorphism  $\mathsf{Tbl}(q) \to \mathsf{Tbl}(q'')$ :
  - function f maps every  $R_i(ar{x}_i)$  in q into an  $S_{lpha}(ar{x}_lpha)$  in q' by definition,
  - function g maps every  $S_{lpha}(ar{x}_{lpha})$  in q' into an atom of q'' by construction.



#### Lemma

Assume CQ  $q=\{\bar{y}\mid \exists \bar{z}\; R_1(\bar{x}_1) \land R_2(\bar{x}_2) \land \cdots \land R_k(\bar{x}_k)\}.$ 

Furthermore, assume q has a semantically equivalent CQ

$$q' = \{ \bar{y'} \mid \exists \bar{z'} S_1(\bar{x'}_1) \land S_1(\bar{x'}_1) \land \cdots \land S_j(\bar{x'}_j) \}.$$

Then q is also semantically equivalent to a CQ q''that can obtained by deleting atoms from q.

Interesting consequence: there is always a unique minimal equivalent query. We call this minimal equivalent subquery of q the  $\underline{core}$  of q.

# An Algorithm for Minimisation

### **Algorithm 1:** Query Minimisation **Input:** Conjunctive query q**Result:** A minimal conjunctive query q' with $q \equiv q$ $q' \leftarrow q$ repeat for Term $R(\bar{x})$ in q' do q''q' without $R(\bar{x})$ if there is a homomorphism $\mathsf{Tbl}(q') \to \mathsf{Tbl}(q'')$ then $q' \leftarrow q''$ end end until no change to q'

return q'

In plain text

Delete terms from the CQ as long as there is still a homomorphism to the query after deletion.

Once this is no longer possible, the minimum is reached.

# CQ Minimisation Example

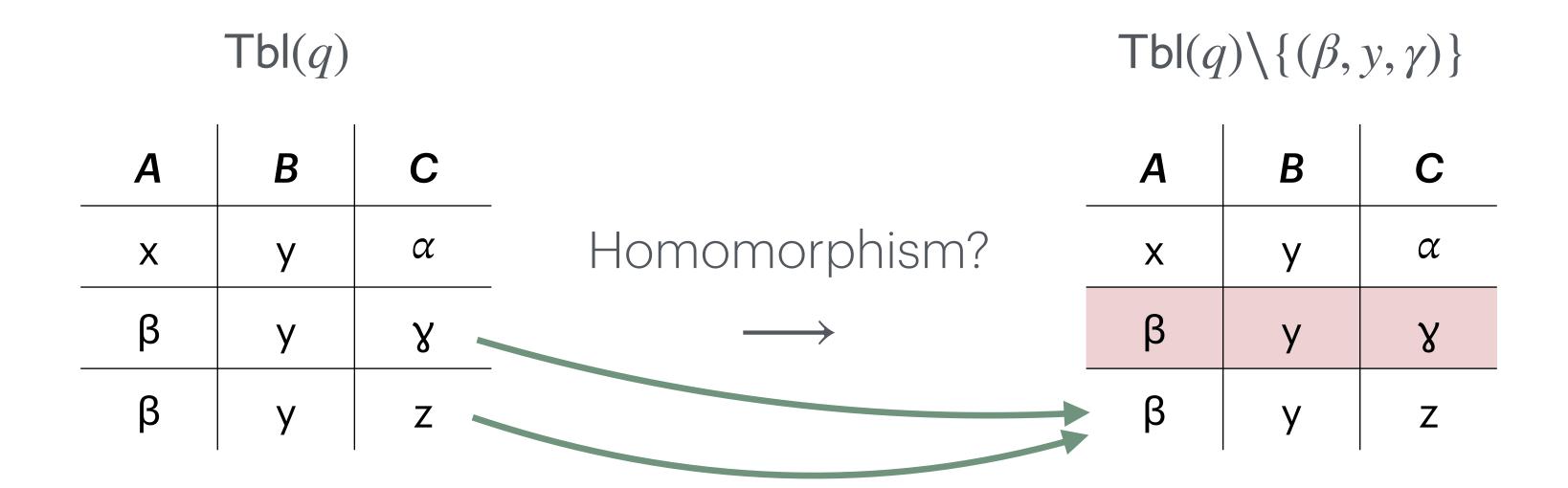
$$q = \{(x, y, z) \mid \exists \alpha \beta \gamma \ R(x, y, \alpha) \land R(\beta, y, \gamma) \land R(\beta, y, z) \}$$

Tbl(q)				$Tbl(q) \setminus \{(x, y, \alpha)\}$		
A	В	С		A	В	С
X	У	α	Homomorphism?	X	У	α
β	У	γ	$\longrightarrow$	β	У	γ
β	У	Z		β	У	Z

No, because h(x) = x, the first row can't be mapped into the right-hand tableau.

# CQ Minimisation Example

$$q = \{(x, y, z) \mid \exists \alpha \beta \gamma \ R(x, y, \alpha) \land R(\beta, y, \gamma) \land R(\beta, y, z) \}$$



Yes!  $x, y, z, \beta$  map to themselves and  $h(\gamma) = z$ 

# CQ Minimisation Example

$$q' = \{(x, y, z) \mid \exists \alpha \beta \ R(x, y, \alpha) \land R(\beta, y, z) \}$$

 $\mathsf{Tbl}(q')$ 

A	В	С
X	У	α
β	У	Z

Homomorphism?



Both times no.

Hence, q' is minimal!

A	В	С
X	У	α
β	y	Z

or

A	В	С
X	У	α
β	У	Z

# Complexity of CQs

# Data vs Combined Complexity

Recall, in classical computational complexity theory we study algorithm behavior (time/space/etc.) relative to the size of the input.

In query answering problems there are different variants of this problem:

## Eval(q, D)

Input size is the sum of the query size and database size

Matches natural settings such as a DBMS, where queries and data come from user and are arbitrary.

## q-Eval(D)

Input size is only the database!

Motivated by the fact that queries are usually much smaller than the databases.

# Data vs Combined Complexity

Recall, in classical computational complexity theory we study algorithm behavior (time/space/etc.) relative to the size of the input.

In query answering problems there are different variants of this problem:

## **Combined Complexity**

Input size is the sum of the query size and database size

Matches natural settings such as a DBMS, where queries and data come from user and are arbitrary.

## Data Complexity

Input size is only the database!

Motivated by the fact that queries are usually much smaller than the databases.

## Our Focus Now

#### : CQ-EVAL

Input: Conjunctive query q, database D (of same schema):

 $\vdots \text{ Output: } q(D) \neq \emptyset$ 

Recall, this corresponds to combined complexity.

# NP-Membership

$$\{ \bar{y} \mid \exists \bar{z} \ R_1(\bar{x}_1) \land R_2(\bar{x}_2) \land \cdots \land R_k(\bar{x}_k) \}$$

When is  $q(D) \neq \emptyset$ ?

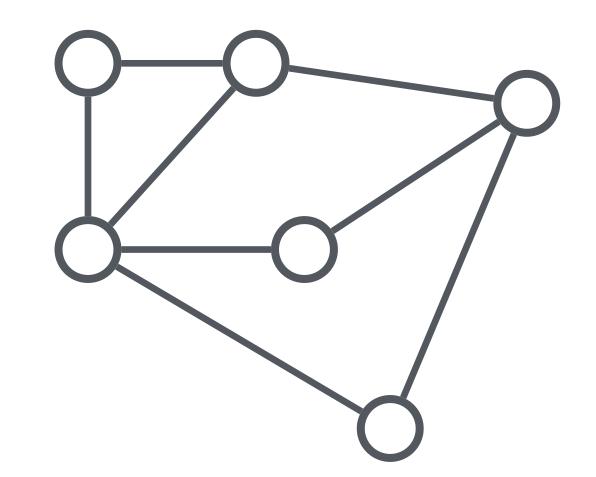
If there is any homomorphism from  $\mathsf{Tbl}^*(q)$  to D.

NP-membership is straightforward: guess and check a homomorphism.

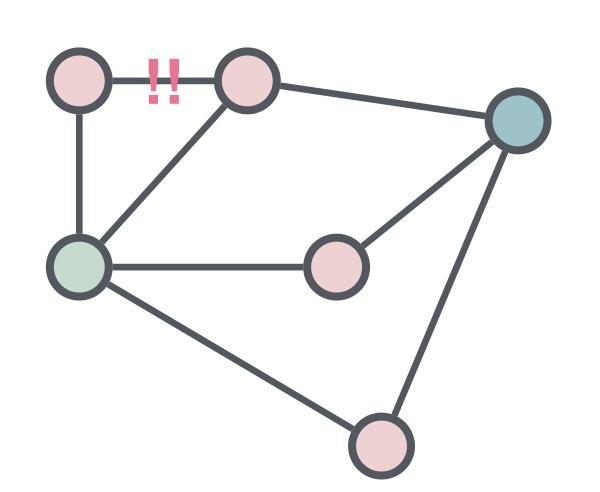
## NP-Hardness

- There is an easy reduction from 3-Colourability.
- lacktriangle 3-Colourability takes a graph G as input and decides whether G is 3-colourable.

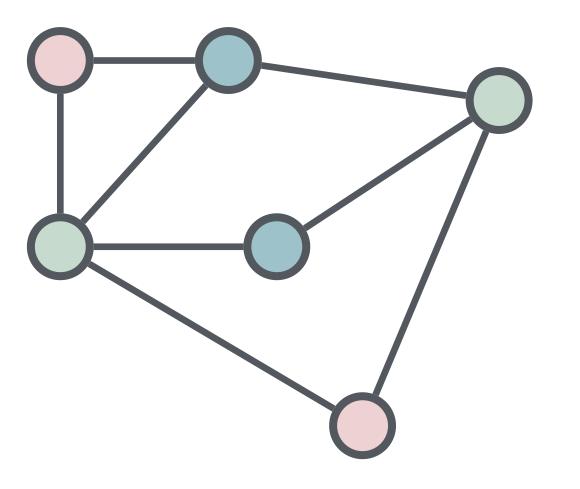
That is, can we color the vertices of G with red, green, and blue such that no edge is between two vertices of the same colour?



Not a 3-colouring

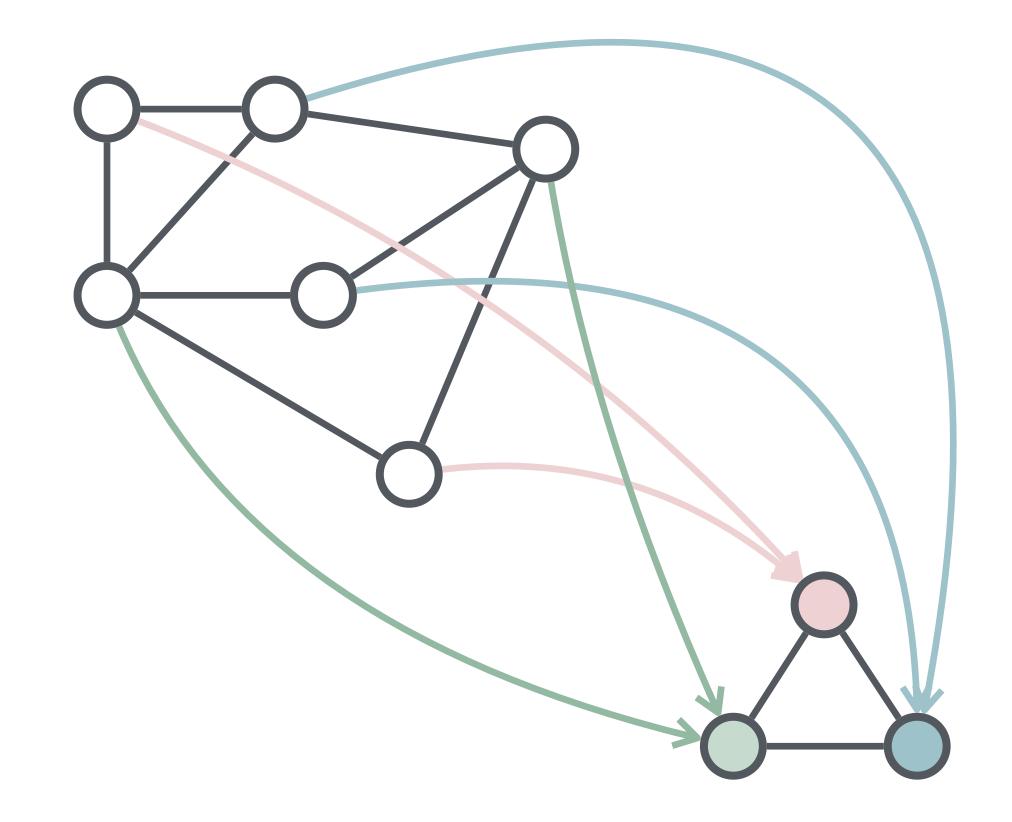


Valid 3-colouring



## NP-Hardness

- ◆ 3-Colourability is equivalent to having a homomorphism into the triangle graph.
- ◆ The three nodes of the triangle intuitively represent the three colours.
- Note that if there is an edge between v and u, then v, u can't be mapped to the same vertex, i.e., adjacent vertices can't be mapped to the same colour.



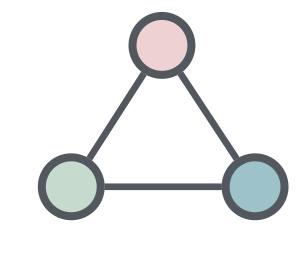
## NP-Hardness

This homomorphism into the triangle can be trivially expressed as a conjunctive query.

- lacktriangle Take an input for **3-Colourability**, i.e., a graph G.
- lacktriangle Create a database with relation E for the triangle:
- ◆ Encode the graph as a conjunctive query:

$$q = \{ () \mid \exists \bar{v} \land \bigwedge E(v_i, v_j) \land E(v_j, v_i) \}$$
$$\{v_i, v_j\} \in E(G)$$

A	В
red	green
green	red
red	blue
blue	red
green	blue
blue	green



lacktriangle There is a homomorphism  $\mathsf{Tbl}^*(q) \to D$  if and only if G is 3-colourable.

# Complexity of CQ-Eval

#### Theorem

CQ-Eval is NP-complete in combined complexity. Moreover, the NP-hardness holds already for schemas with a single binary relation symbol.

# Complexity of CQ Containment

Recall the Homomorphism Theorem:

$$q_1 \subseteq q_2 \iff \mathsf{Tbl}(q_2) \stackrel{hom}{\longrightarrow} \mathsf{Tbl}(q_1)$$

Same reduction applies here too: check whether the query  $q_1$  that represents the triangle is contained in the query  $q_2$  that represents graph G.

#### Theorem

Deciding CQ Containment is NP-complete.

# Complexity of CQ Minimisation

#### Theorem

Checking whether a query q is minimal is co-NP-complete.

**Intuition:** We need to check whether there are *no homomorphisms* into any query obtained by deleting atoms.

# Structure of CQs

## Structure?

The number of joins is not the only important factor in how difficult it is to evaluate a CQ.

$$q_1 = \{ (x) \mid R(x, y, z) \land S(z, w, u) \land R(u, v, x) \}$$

Which is easier to evaluate 1?

$$q_2 = \{ (x) \mid R(x, y, z) \land S(z, w, u) \land R(u, v, x)$$

$$\land T(x, y, z, w, u) \}$$

 $q_2$  can be solved in time linear in the database  $q_1$  cannot (under standard assumptions from fine-grained complexity theory).

## Structure

To understand why, we need to understand how the structure of CQs is connected to their evaluation.

But what is the structure of a CQ?

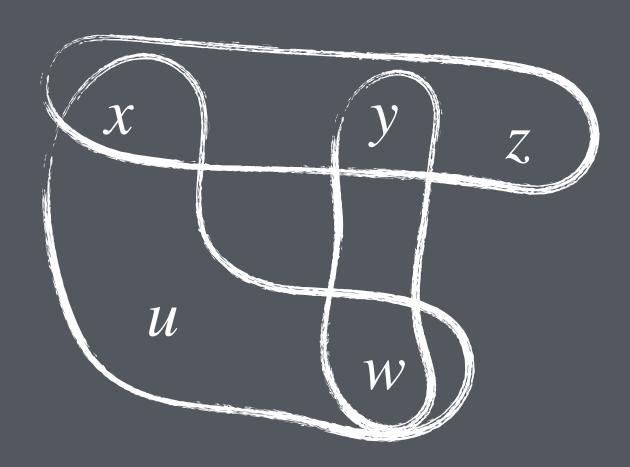
# Hypergraphs

- lacktriangle A hypergraph H is a pair  $(V_H, E_H)$ .
- $lacktriangledown V_H$  is the set of vertices, just as in a graph.
- ♦ The set of hyperedges  $E_H \subseteq 2^{V(H)}$  is some set of sets of vertices.
- ◆ Graphs are a special case of hypergraphs where every hyperedge has size exactly 2.

#### Example hypergraph H

$$V_H = \{x, y, z, u, w\}$$

$$E_H = \{\{x, y, z\}, \{y, u\}, \{x, u, w\}\}$$



# CQs as Hypergraphs

$$q = \{ \bar{y} \mid \exists \bar{z} R_1(\bar{x}_1) \land R_2(\bar{x}_2) \land \cdots \land R_k(\bar{x}_k) \}$$

We capture the structure of q as hypergraph H in the following way:

- $lacktriangledow V_H = vars(q)$ , i.e., we have a vertex for each variable in the query.
- lacktriangle For every atom  $R(x_1,\ldots,x_{\#R})$  we add the hyperedge  $\{x_1,\ldots,x_{\#R}\}$  to  $E_H$ .

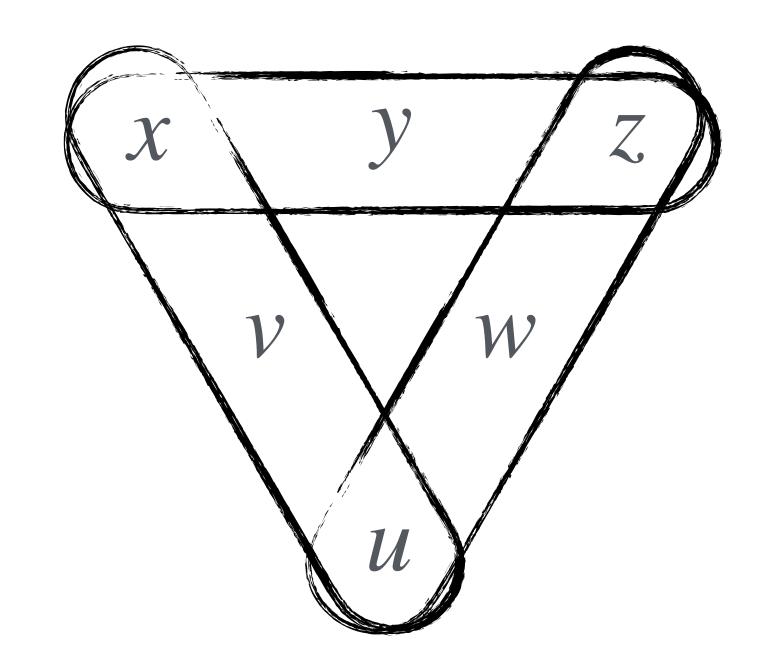
# Example

$$q_1 = \{ (x) \mid R(x, y, z) \land S(z, w, u) \land R(u, v, x) \}$$

$$V_H = vars(q_1)$$

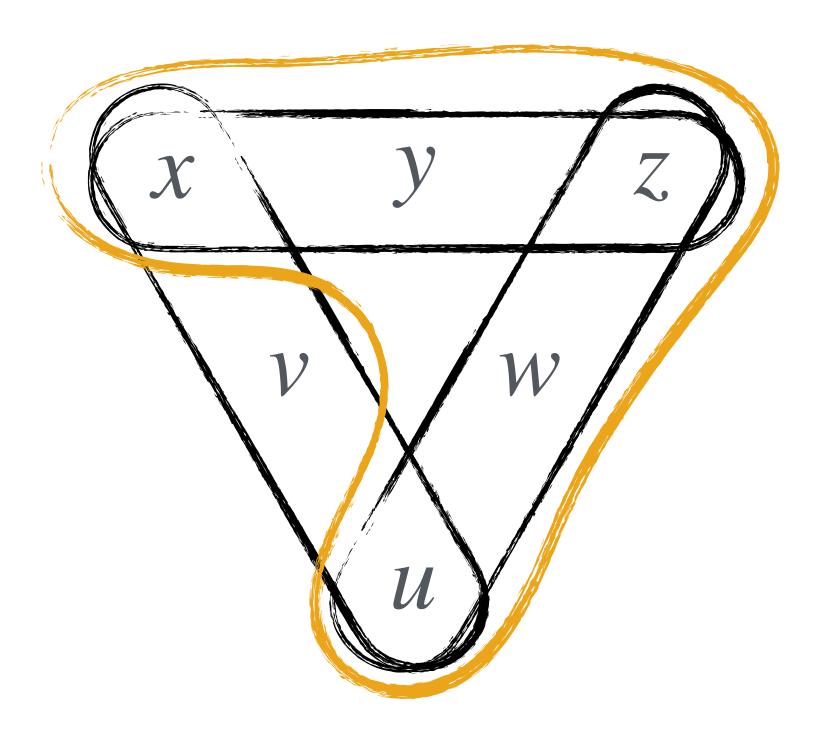
Add edges for each atom.

R(x, y, z) becomes edge  $\{x, y, z\}$ , and so on.



# Example

$$q_2 = \{ (x) \mid R(x, y, z) \land S(z, w, u) \land R(u, v, x)$$
$$\land T(x, y, z, w, u) \}$$



# When is the structure of a CQ good for evaluation?

# Join Trees & Yannakakis Algorithm

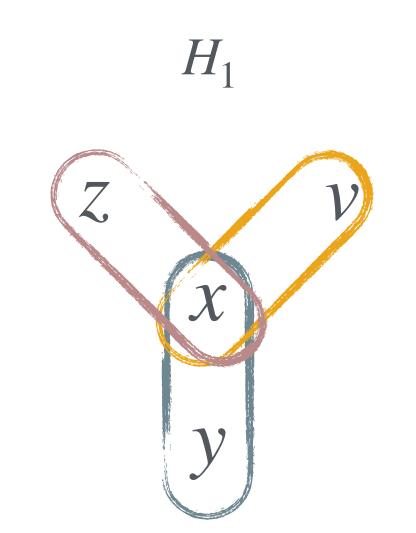
# Acyclic Hypergraphs

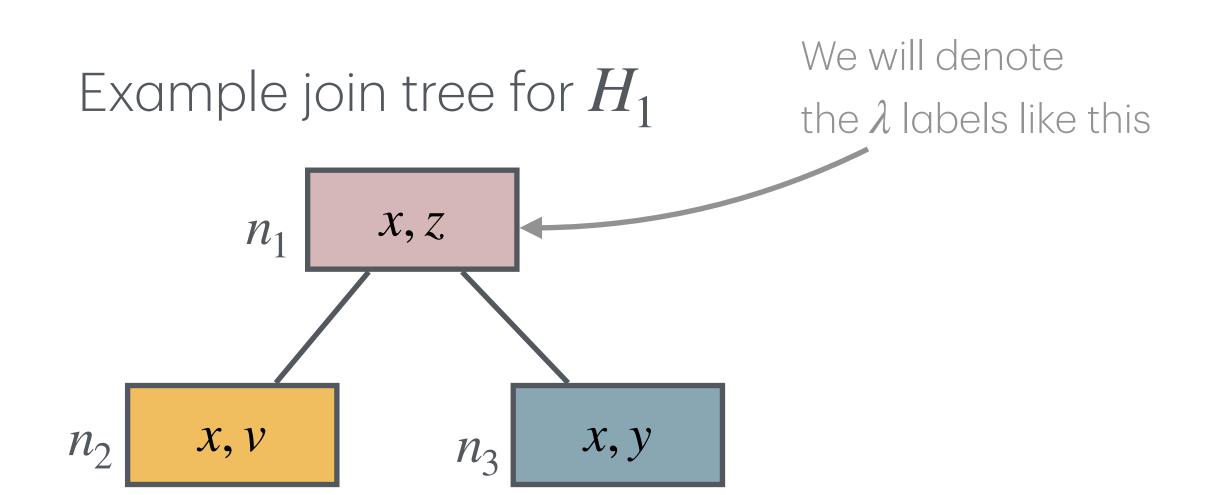
A join tree for a hypergraph H is a tree T with a labelling function  $\lambda:V(T)\to E(H)$  that labels every node of the tree with a an edge of H such that:

- i) For every  $e \in E(H)$  there is exactly one node n of the tree T such that  $\lambda(n) = e$
- ii) For every vertex  $v \in V(H)$ , the nodes for which  $v \in \lambda(n)$  are connected. (connectedness condition)

We say that a hypergraph H is lpha-acyclic if it has a join tree.

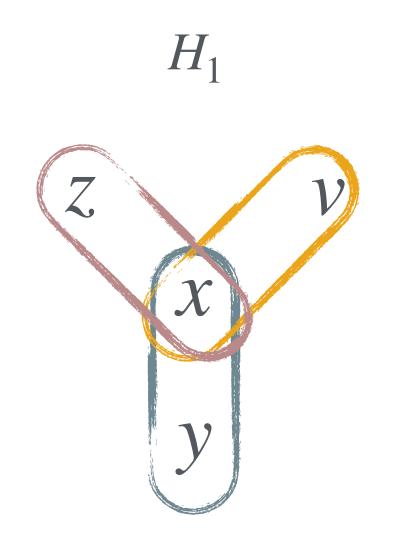
# Example

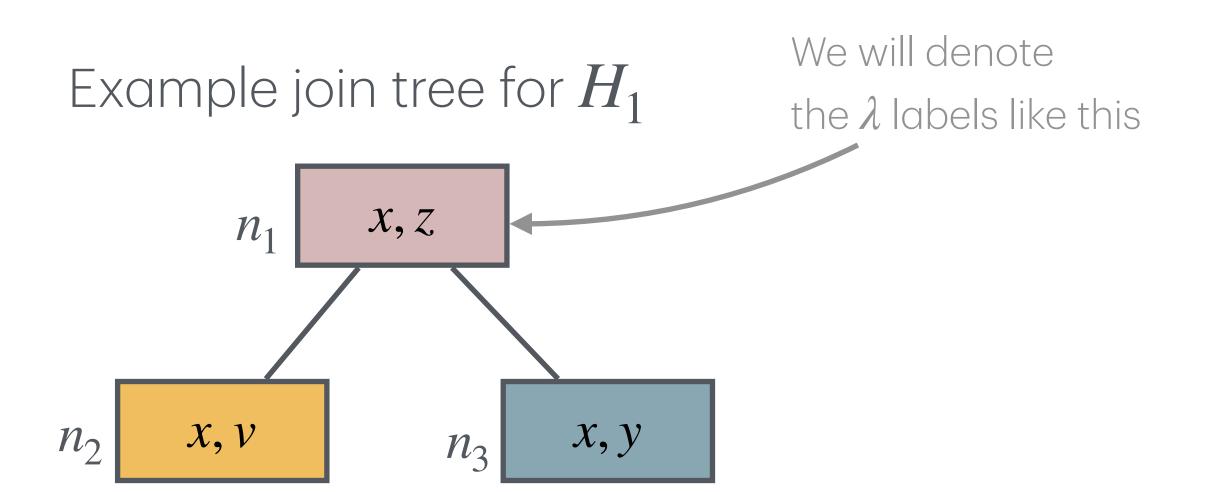




Let's check the two conditions:

- ◆ Every edge is mapped to some label
- ◆ Connectedness condition?



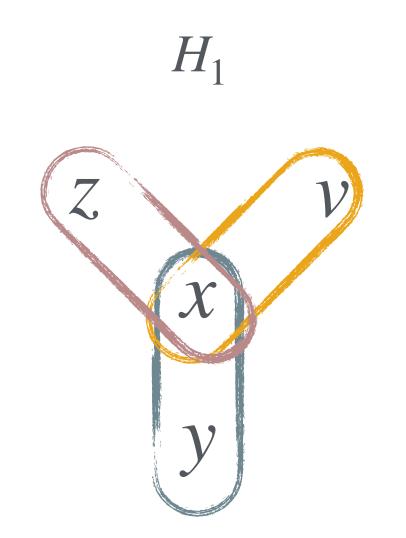


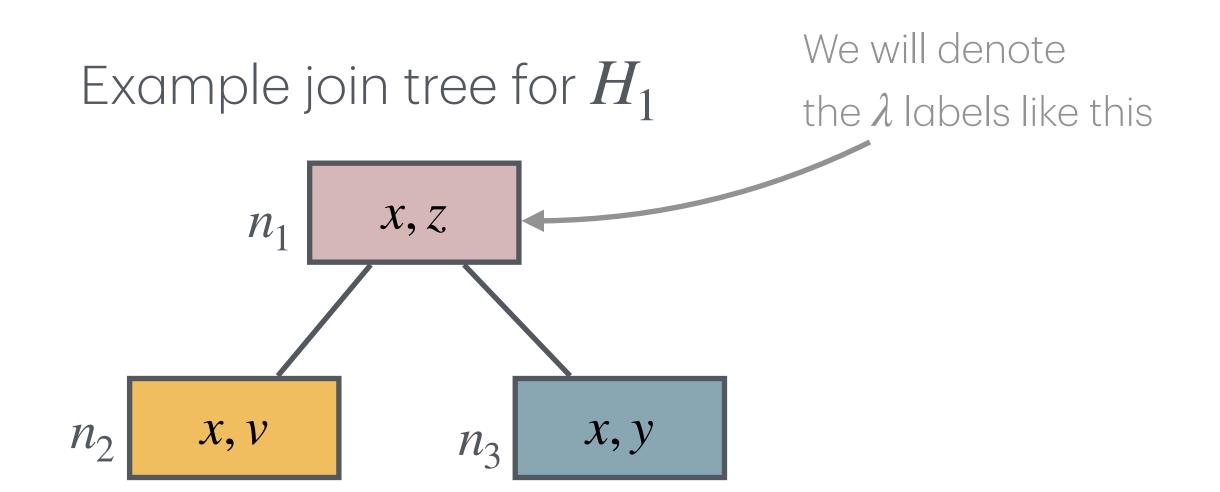
Let's check the two conditions:

- ◆ Every edge is mapped to some label
- ◆ Connectedness condition?

Nodes that contain x







Let's check the two conditions:

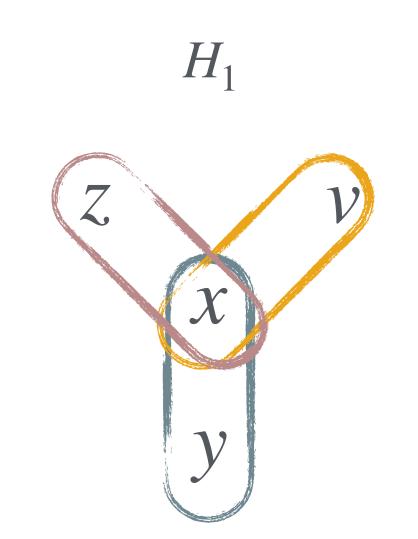
Every edge is mapped to some label

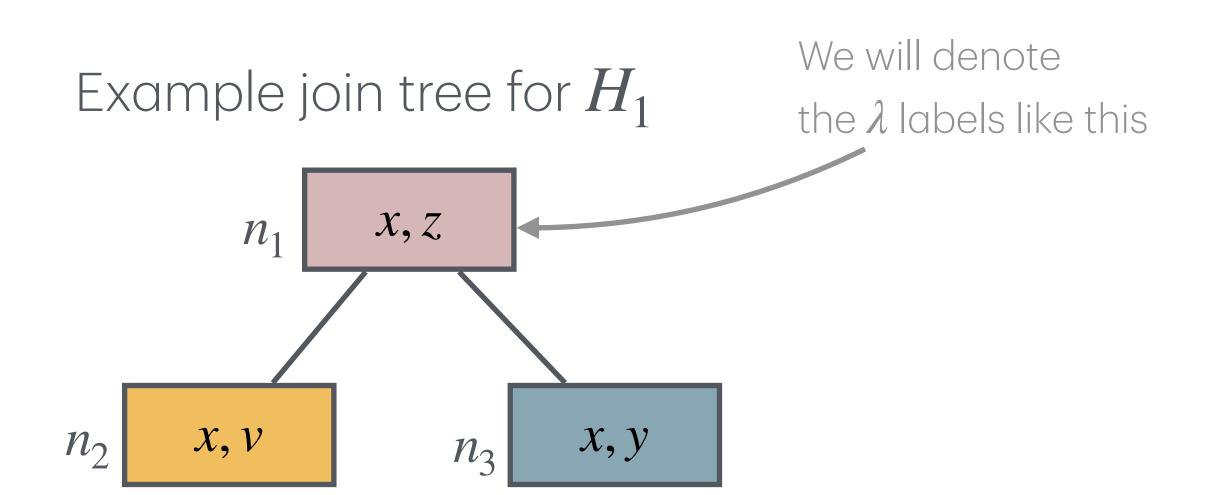
◆ Connectedness condition?

Nodes that contain y

Connected  $n_3$ 

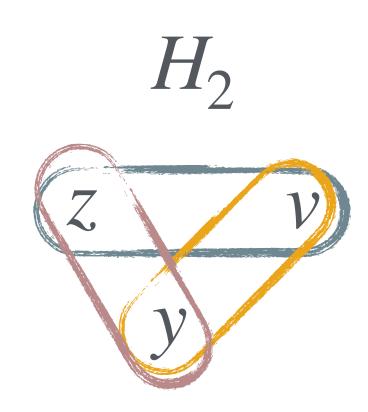
Same for v and z.



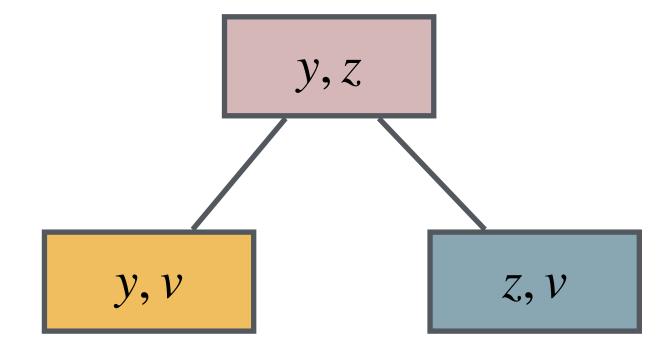


Let's check the two conditions:

- ◆ Every edge is mapped to some label
- ◆ Connectedness condition ✓

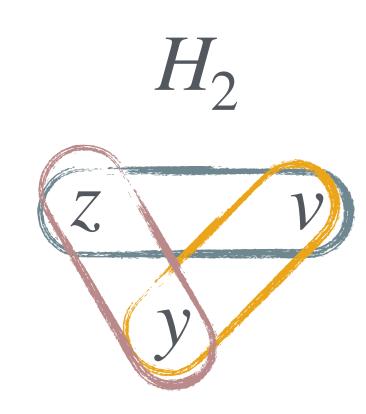


 $H_2$  has no join tree!

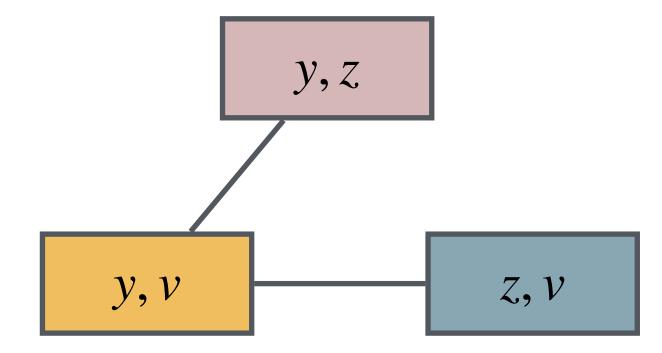


v is in the blue and yellow node but not in the red.

Violates connectedness condition!

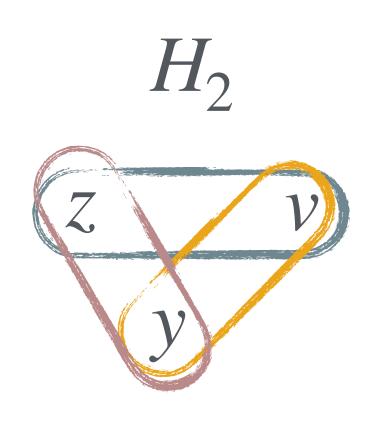


 $H_2$  has no join tree!

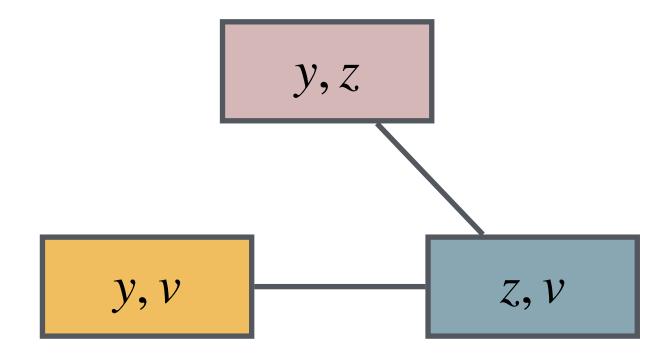


z is in the blue and red node but not in the yellow.

Violates connectedness condition!



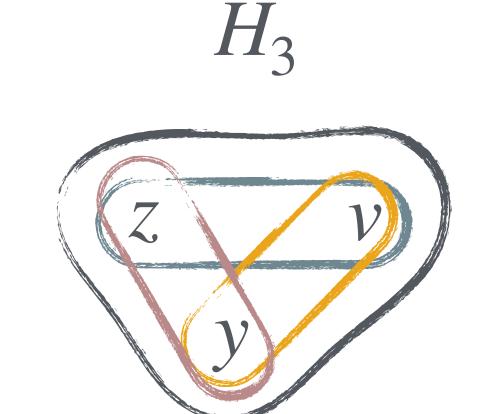
 $H_2$  has no join tree!



This covers all possible trees on three nodes. We see that none of them is a join tree!

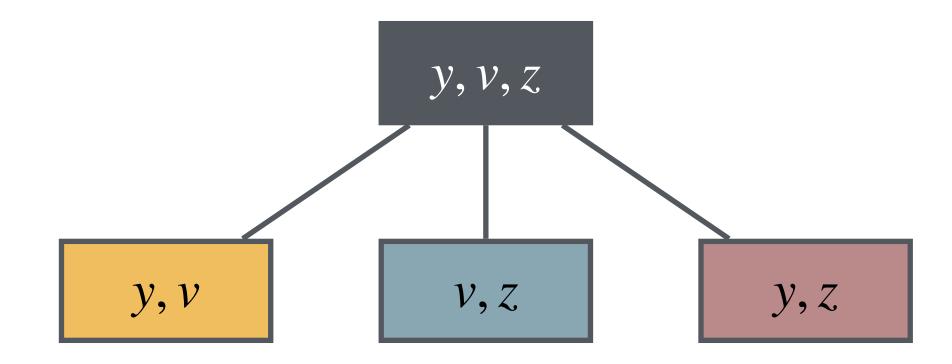
y is in the yellow and red node but not in the blue.

Violates connectedness condition!



That is, the triangle + the edge  $\{y, v, z\}$ 

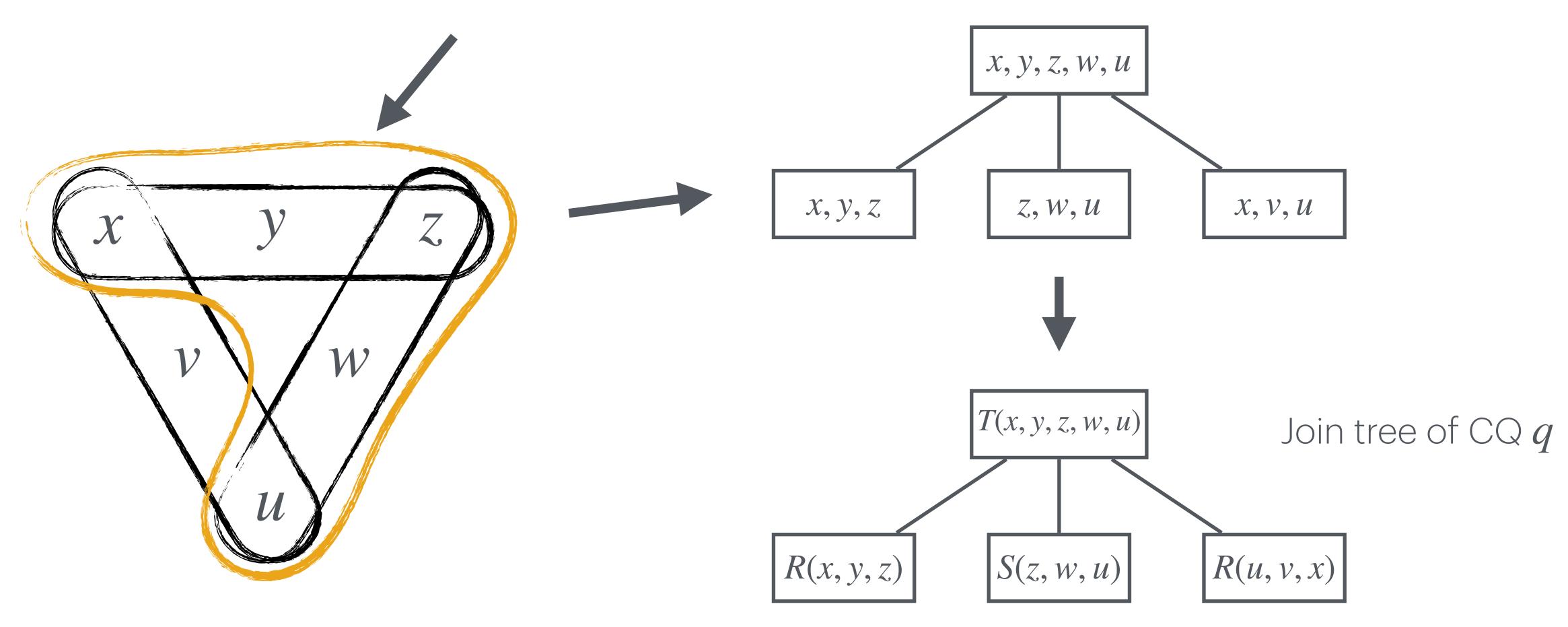
 $H_3$  has a join tree!



 $H_3$  is lpha-acyclic even though it contains a cyclic subhypergraph!

#### Back to Queries

$$q_2 = \{ (x) \mid R(x, y, z) \land S(z, w, u) \land R(u, v, x)$$
$$\land T(x, y, z, w, u) \}$$



# We can use join trees to design efficient algorithms for CQ evaluation.

#### The Semi-join Operation X

$$R \bowtie S := \pi_{\operatorname{attr}(R)}(R \bowtie S)$$

Instead of creating the combined tuples as in a join, a semi-join only keeps the rows in  ${\it R}$  that have a join partner in  ${\it S}$ 

Α	В	C
3	2	5
6	7	9
5	5	5



В	С	D
1	9	1
2	5	3
5	5	1
1	3	2

A	В	C	
3	2	5	
6		9	
5	5	5	

#### The Semi-join Operation X

$$R \bowtie S := \pi_{\operatorname{attr}(R)}(R \bowtie S)$$

Instead of creating the combined tuples as in a join, a semi-join only keeps the rows in  $\emph{R}$  that have a join partner in  $\emph{S}$ 

Unlike a join, a semi-join cannot create larger relations:

$$|R \bowtie S| \leq |R|$$

In fact, we can compute any semi-join in  $O(n \log n)$  time, where n = |R| + |S|.

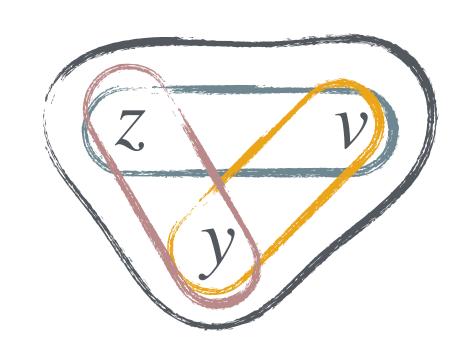
Show this in a theory exercise!

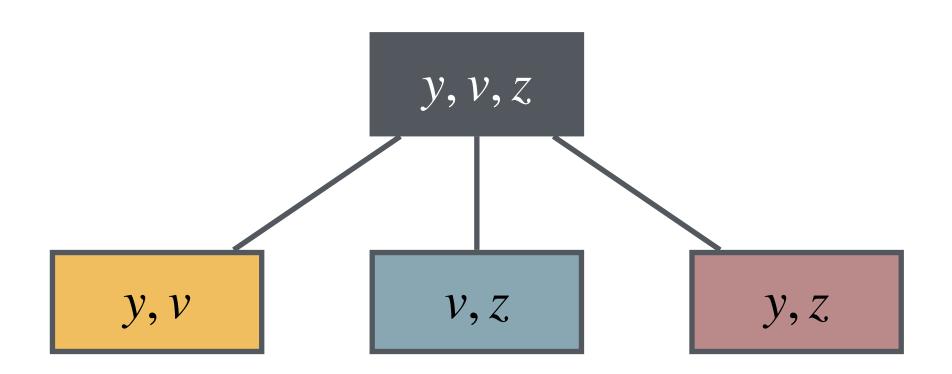
# Yannakakis' Algorithm

Let  $(T, \lambda)$  be a join tree of a CQ q and root it arbitrarily. Given a database D we can decide  $q(D) \neq \emptyset$  as follows:

- 1. Assign to each node labeled with atom  $A(\bar{x})$  the relation  $A^D$  (and rename columns according to the variables in the respective atom). We write rel(N) for the relation associated to node n.
- 2. In a bottom-up traversal: update rel(n) to  $\left(rel(n) \ltimes rel(c_1)\right) \ltimes \cdots \ltimes rel(c_\ell)$  where  $c_1, \ldots, c_\ell$  are the children of n
- 3. At the end, for the root r we have  $rel(r) \neq \emptyset$  if and only if  $q(D) \neq \emptyset$ .

$$R(v,z) \wedge S(z,y) \wedge R(y,v) \wedge T(v,y,z)$$

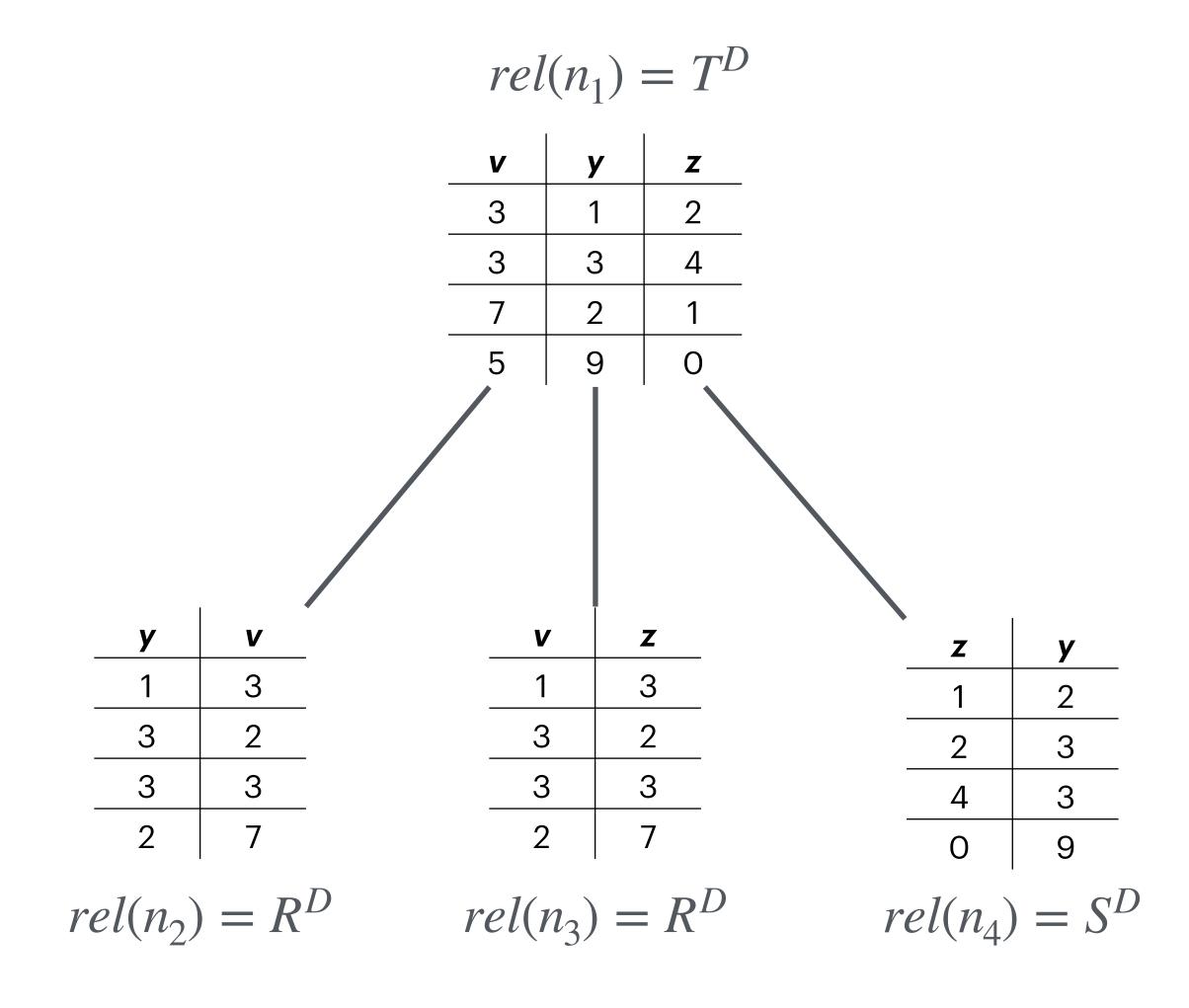




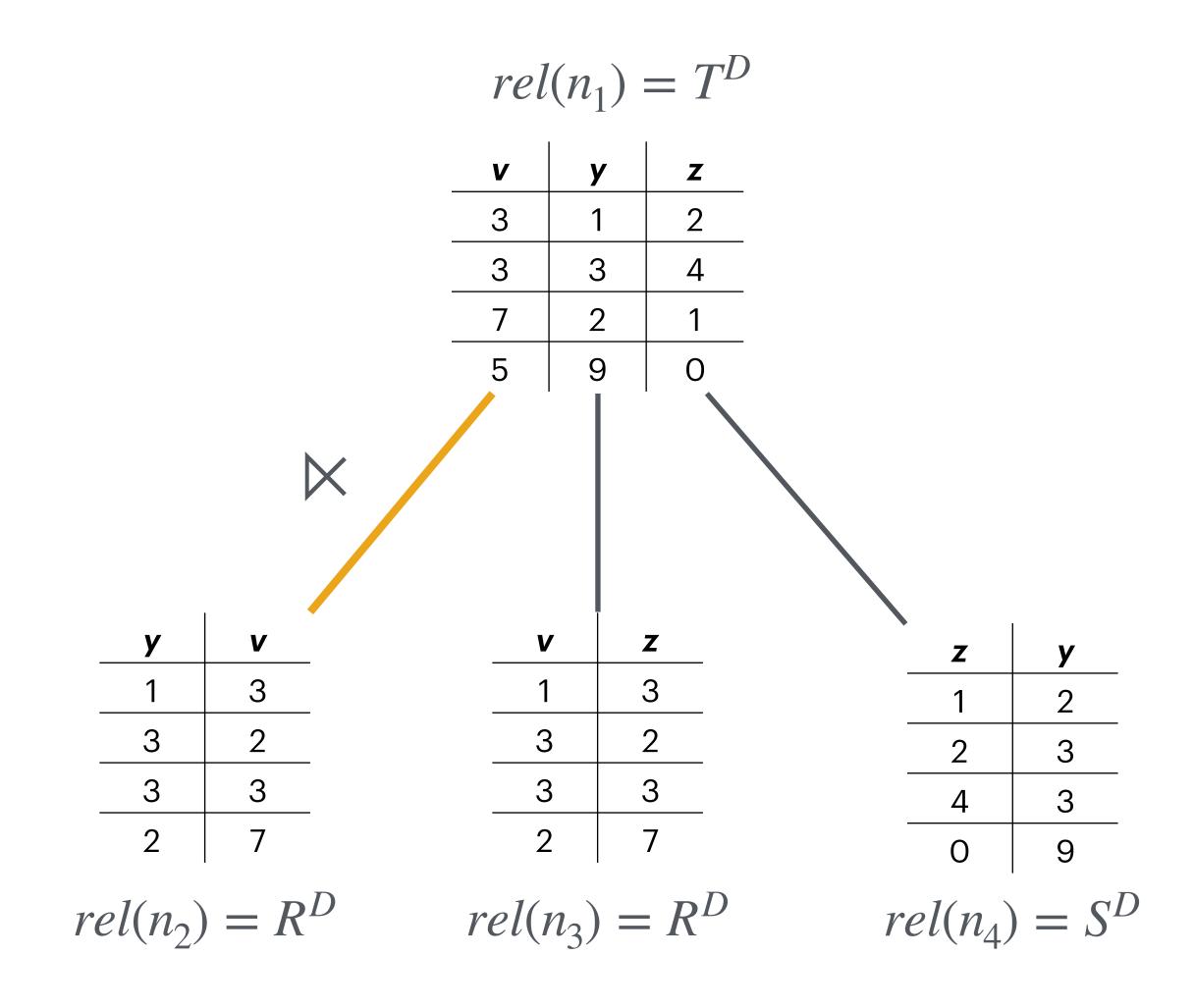
 $R(v,z) \wedge S(z,y) \wedge R(y,v) \wedge T(v,y,z)$ 

1. Assign to each node labeled with atom  $A(\bar{x})$  the relation  $A^D$  (and rename columns according to the variables in the respective atom).

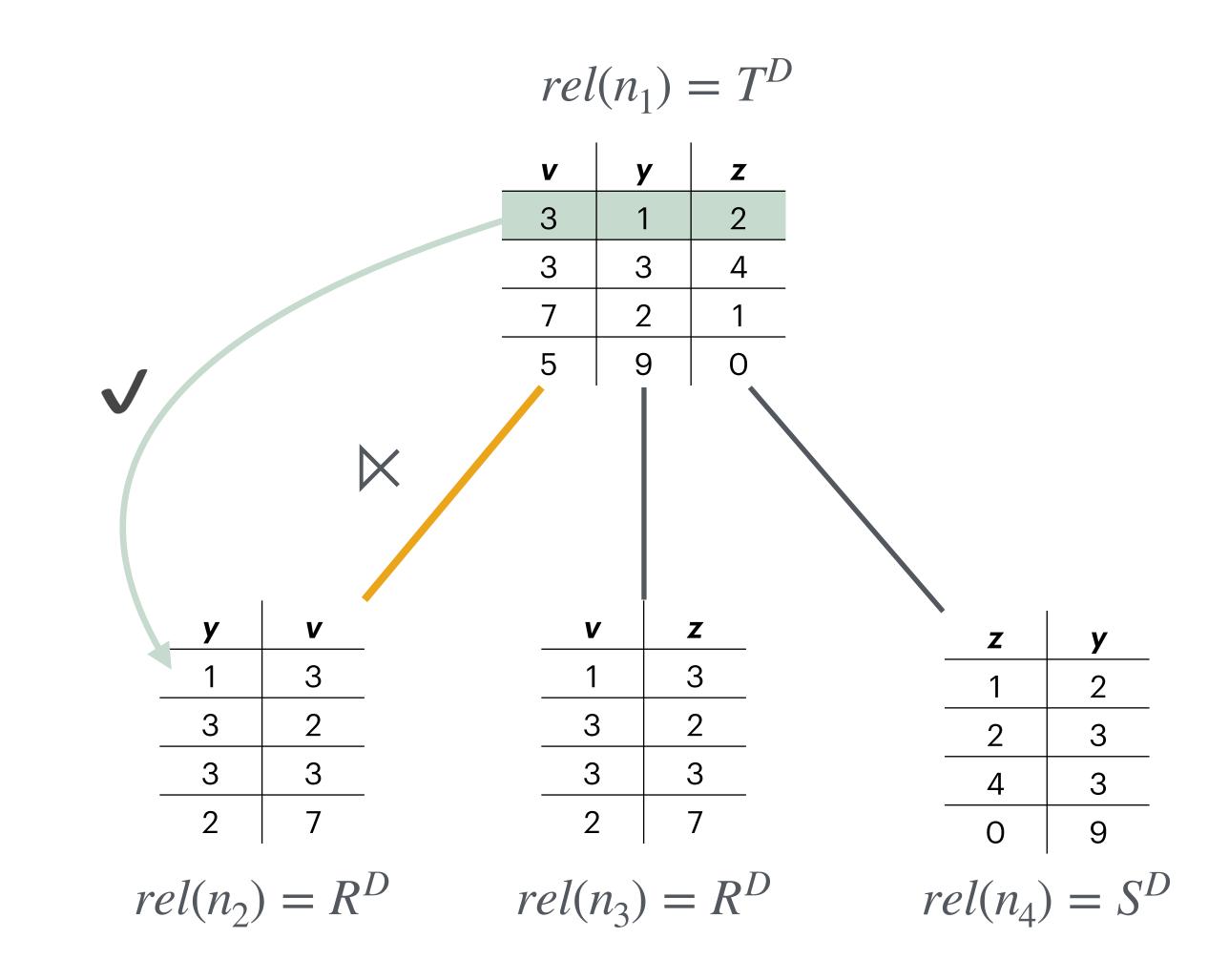
We write rel(N) for the relation associated to node n.



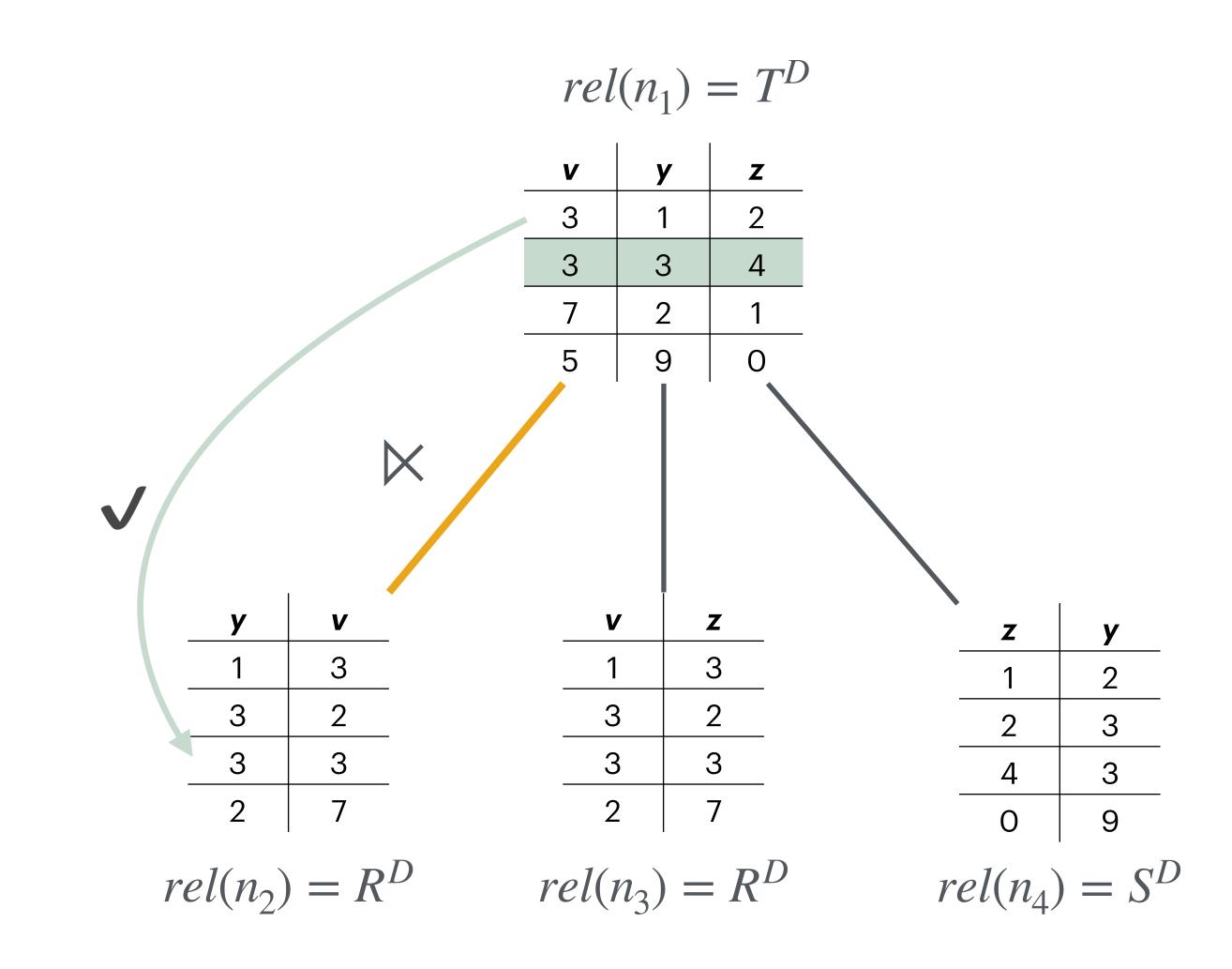
$$R(v,z) \wedge S(z,y) \wedge R(y,v) \wedge T(v,y,z)$$



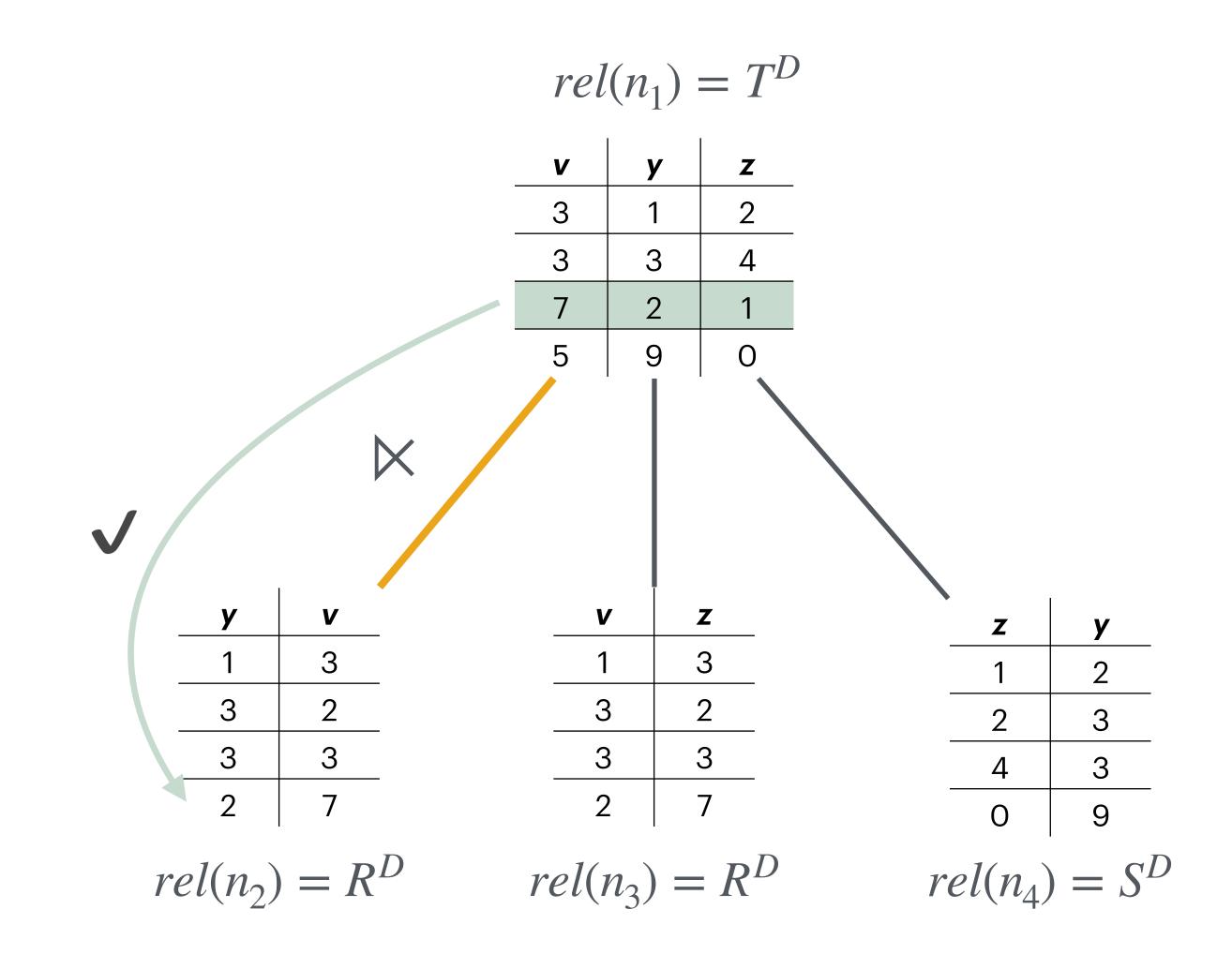
$$R(v,z) \wedge S(z,y) \wedge R(y,v) \wedge T(v,y,z)$$



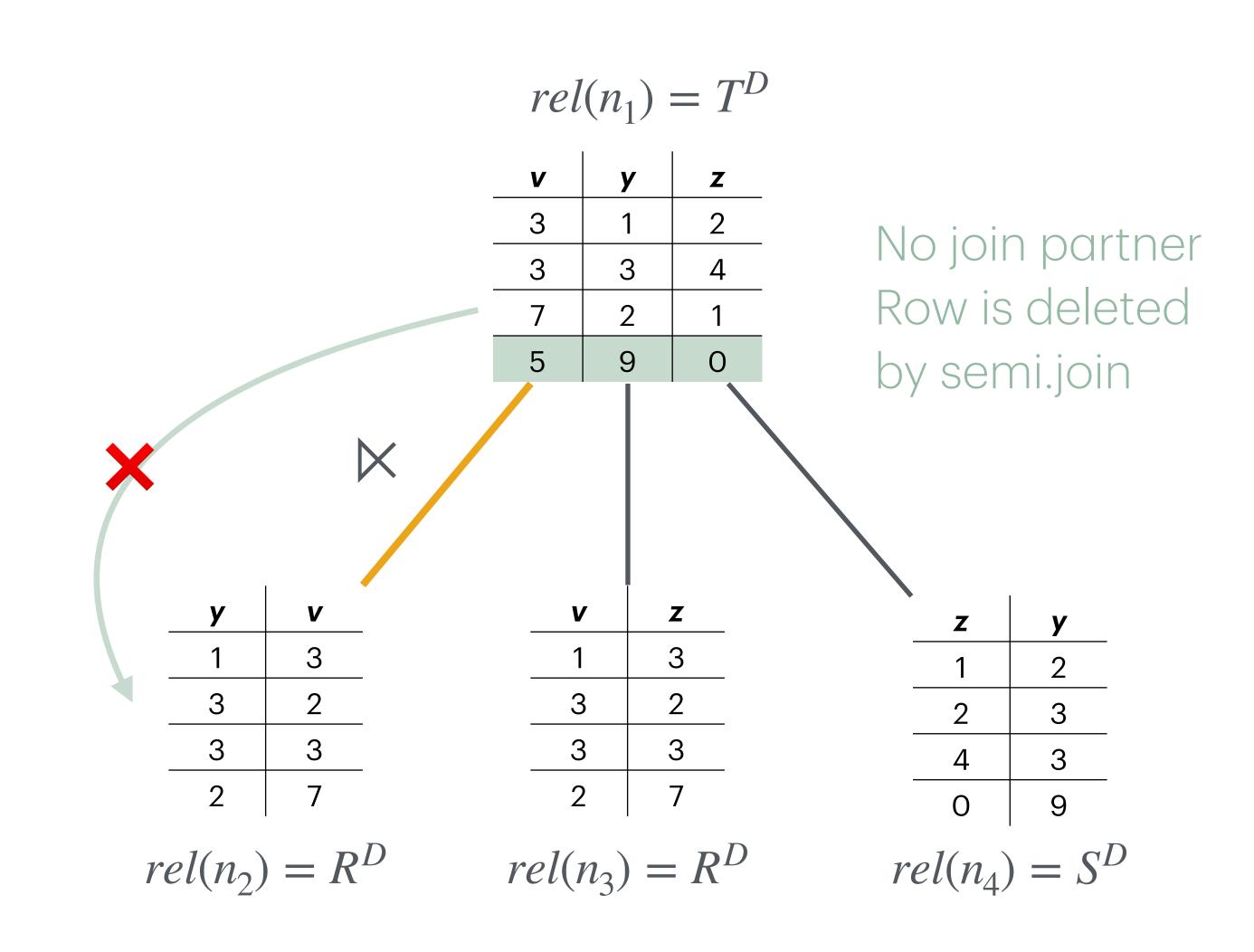
$$R(v,z) \wedge S(z,y) \wedge R(y,v) \wedge T(v,y,z)$$



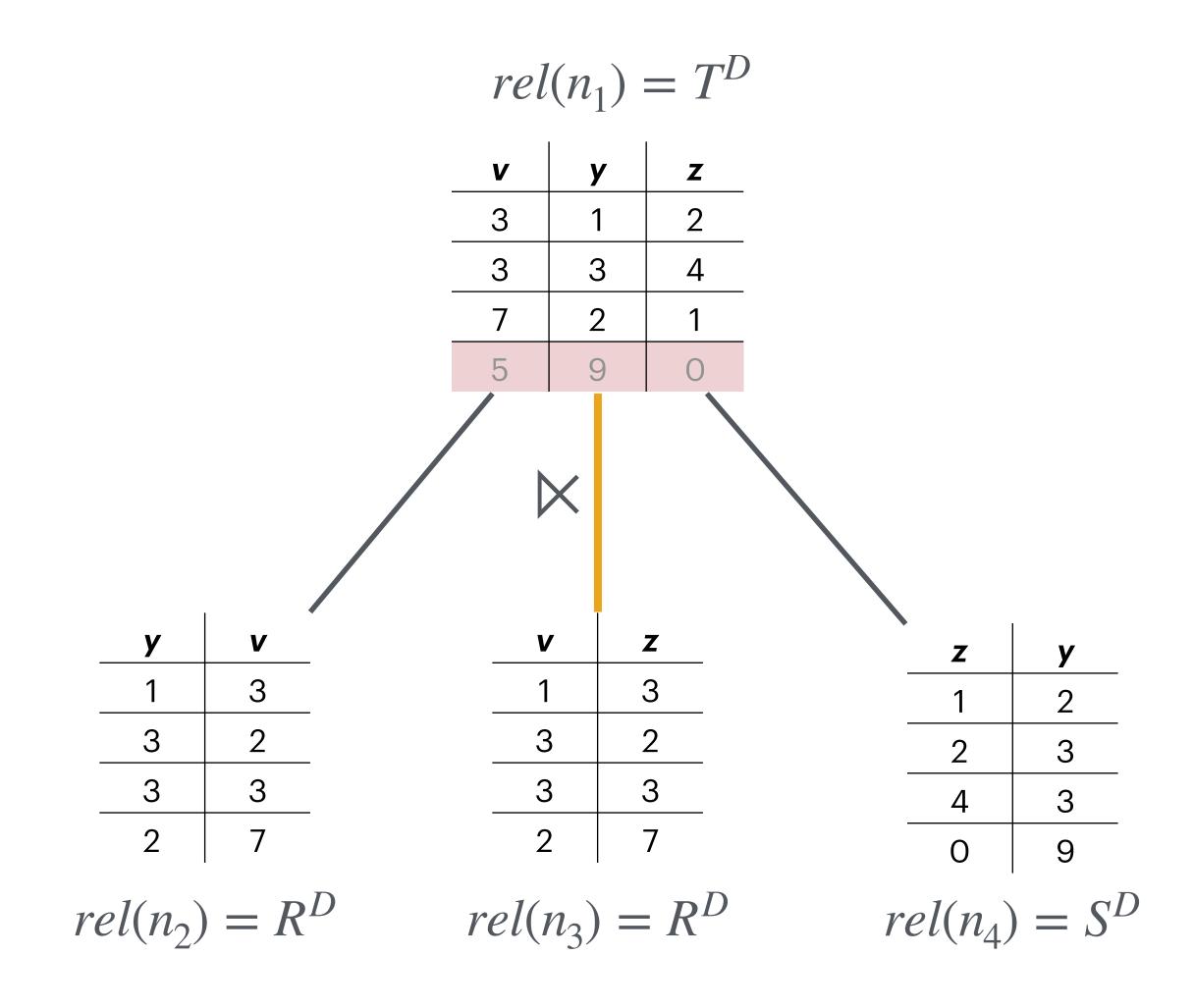
$$R(v,z) \wedge S(z,y) \wedge R(y,v) \wedge T(v,y,z)$$



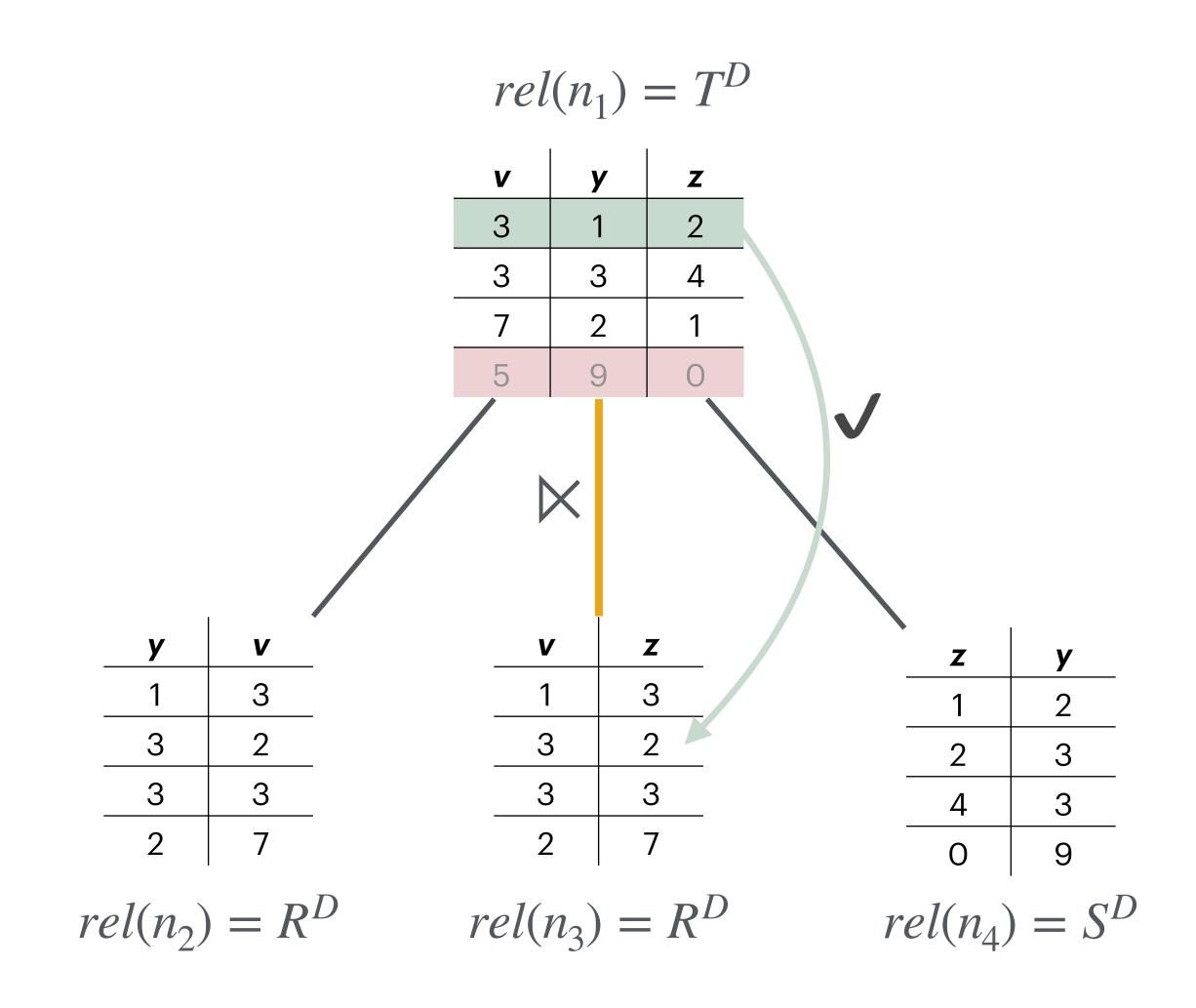
$$R(v,z) \wedge S(z,y) \wedge R(y,v) \wedge T(v,y,z)$$



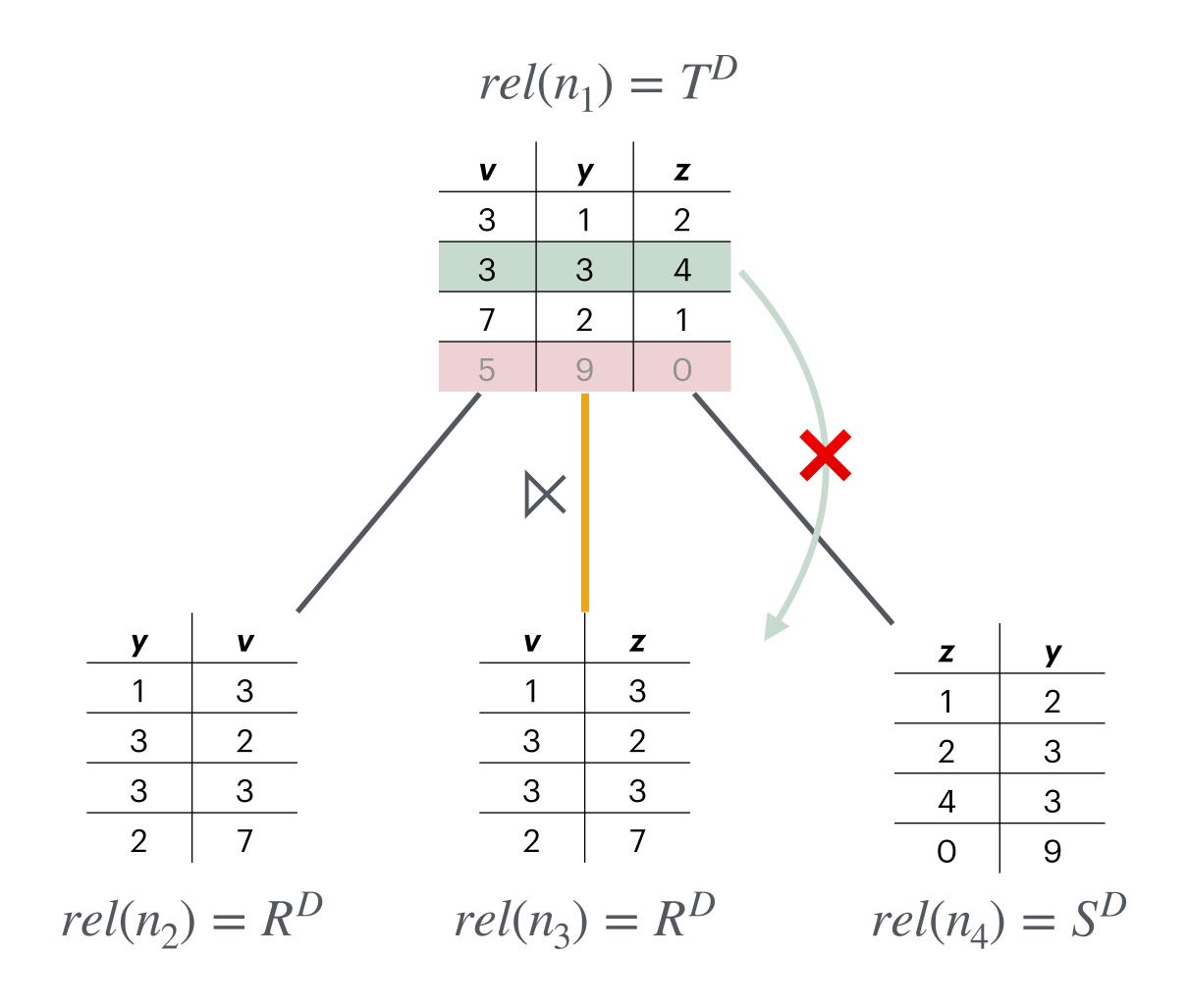
$$R(v,z) \wedge S(z,y) \wedge R(y,v) \wedge T(v,y,z)$$



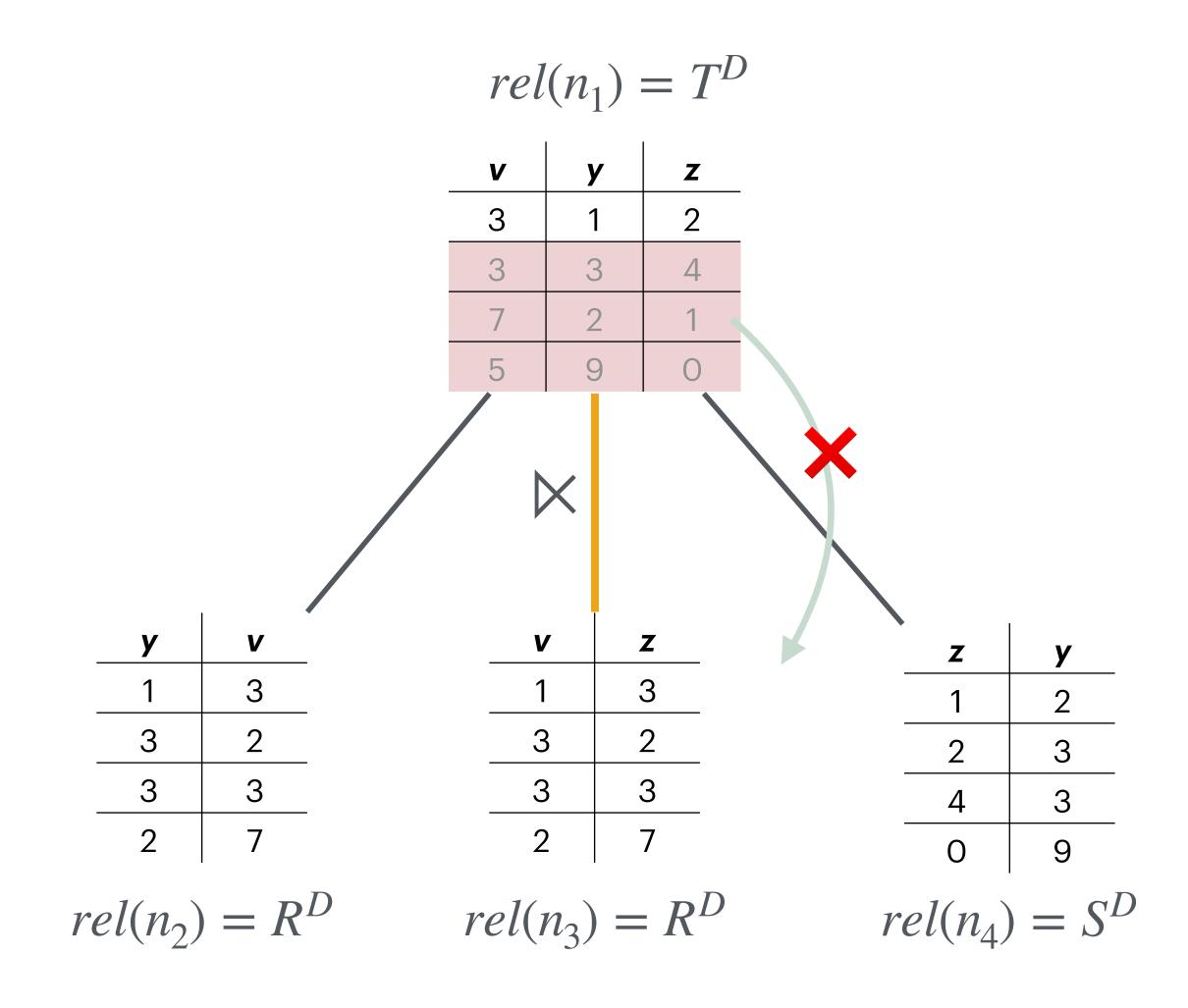
$$R(v,z) \wedge S(z,y) \wedge R(y,v) \wedge T(v,y,z)$$



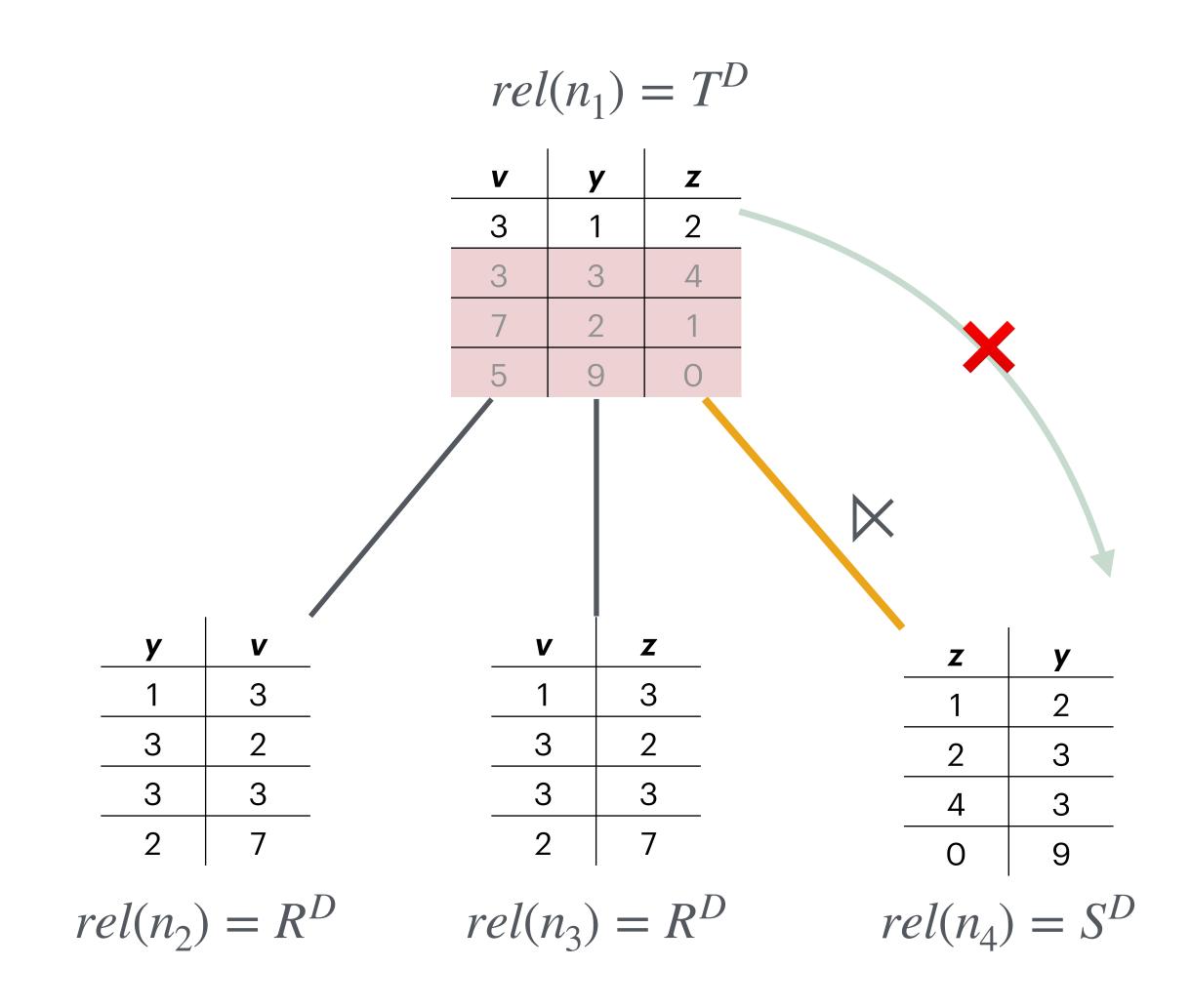
$$R(v,z) \wedge S(z,y) \wedge R(y,v) \wedge T(v,y,z)$$



$$R(v,z) \wedge S(z,y) \wedge R(y,v) \wedge T(v,y,z)$$

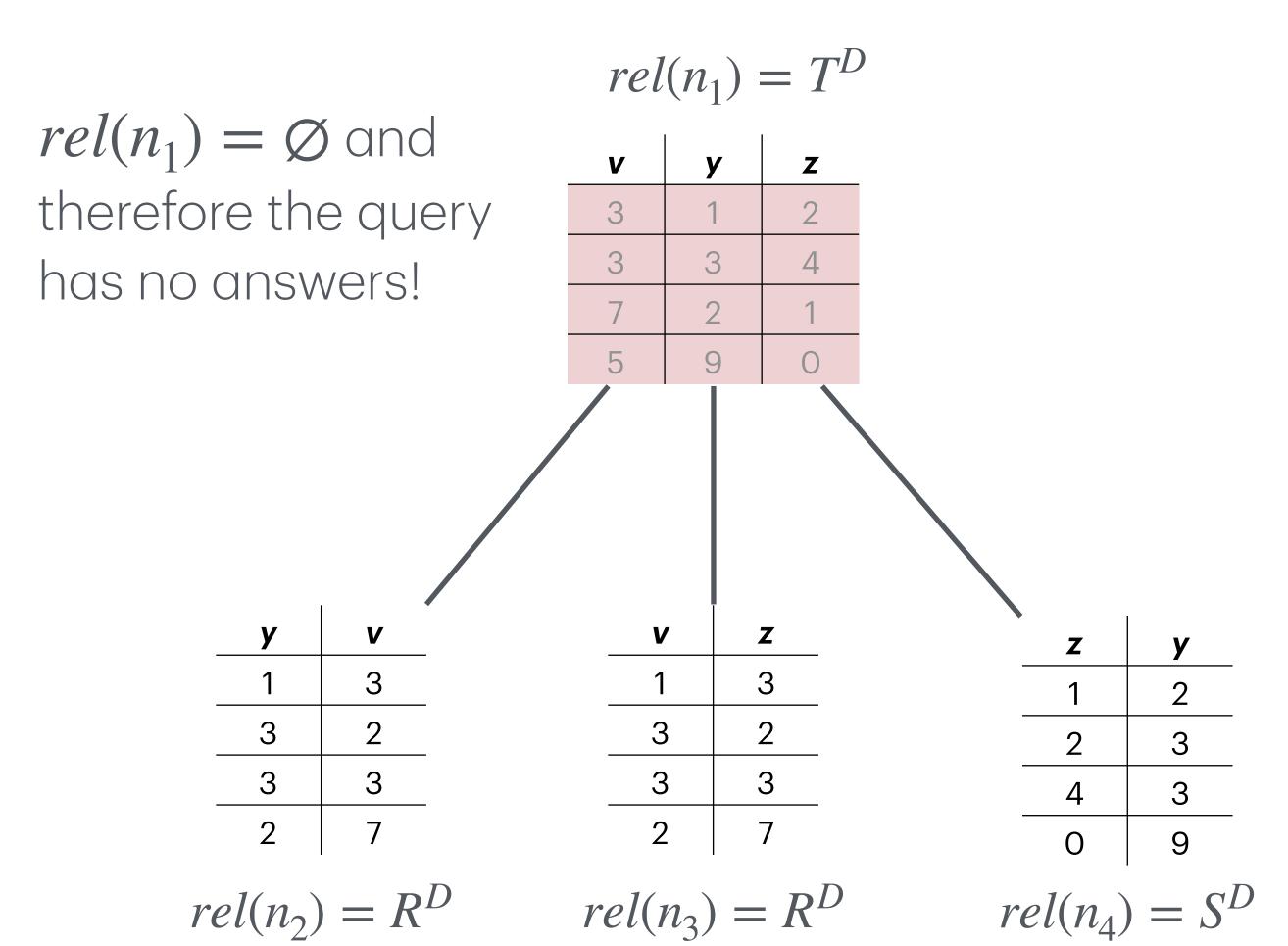


 $R(v,z) \wedge S(z,y) \wedge R(y,v) \wedge T(v,y,z)$ 



 $R(v,z) \wedge S(z,y) \wedge R(y,v) \wedge T(v,y,z)$ 

3. At the end, for the root r we have  $rel(r) \neq \emptyset$  if and only if  $q(D) \neq \emptyset$ .



 $R(v,z) \wedge S(z,y) \wedge R(y,v) \wedge T(v,y,z)$ 

 $rel(n_1) = \emptyset$  and therefore the query has no answers!

 $rel(n_1) = T^D$ 

V	y	Z
3	1	2
3	3	4
7	2	1
5	9	0

How much did it cost to run the algorithm?

Atoms(q) - 1 many semi-joins.

Semi-joins each require  $O(n \log n)$  time.

In general:  $\tilde{O}(|Atoms(q)|\cdot|D|)$ . A big improvement over NP.

У	V
1	3
3	2
3	3
2	7

$$rel(n_2) = R^D$$

V	Z
1	3
3	2
3	3
2	7

$$rel(n_3) = R^D$$

<b>Z</b>	У
1	2
2	3
4	3
0	9

$$rel(n_4) = S^D$$

#### Beyond $\alpha$ -Acyclicity

In most real-world settings the vast majority of queries are  $\alpha$ -acyclic. A large study [1] of over 56mio unique SPARQL queries (a graph query language) found that over 99.9% of queries were acyclic.

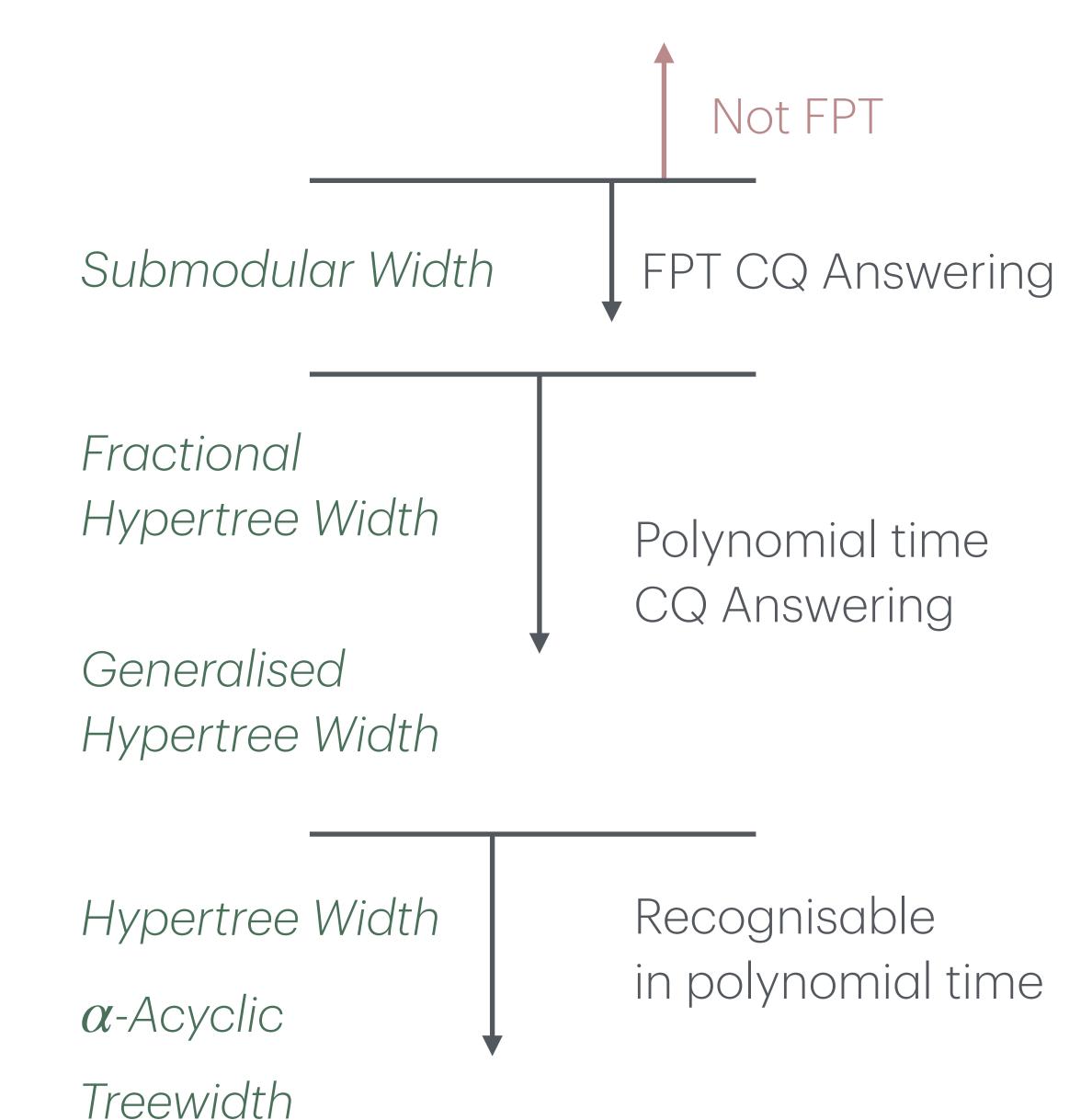
[1] Bonifati, Angela, Wim Martens, and Thomas Timm. "An analytical study of large SPARQL query logs." https://www.vldb.org/pvldb/vol11/p149-bonifati.pdf

However, cyclic queries do come up, and we would like to know if we can do something like Yannakakis' algorithm for those queries. This has lead to the development of *hypertree decompositions* and related notions. They generalize what we've seen here beyond acyclicity!

Explore more in a literature exercise!

#### A Glimpse Beyond

- The connection between structure and complexity has been an active research topic over the last decades.
- lacktriangle There is a rich theory generalising the concept of lpha-acyclicity to not only identify linear time queries.



# Enumerating Answers

# Efficient decision is nice, but what if we want to get the answers?

#### Enumeration

- ◆ For enumeration, quantification of variables will matter again.
- ◆ We will focus on the special case of full CQs, i.e., queries where no variables are existentially quantified.

Example full CQ

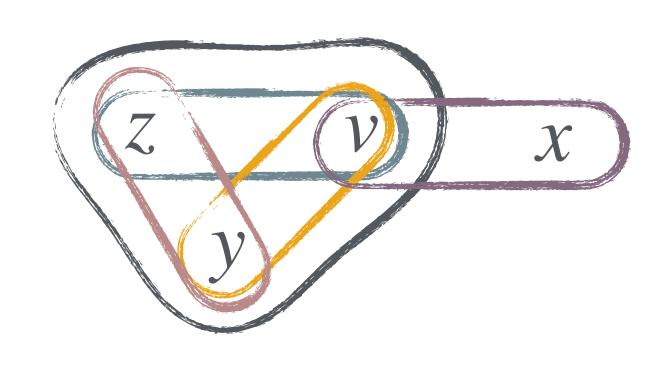
$$q = \{ (v, x, y, z) \mid R(v, z) \land S(z, y) \land R(y, v) \land T(v, y, z) \land P(x, v) \}$$

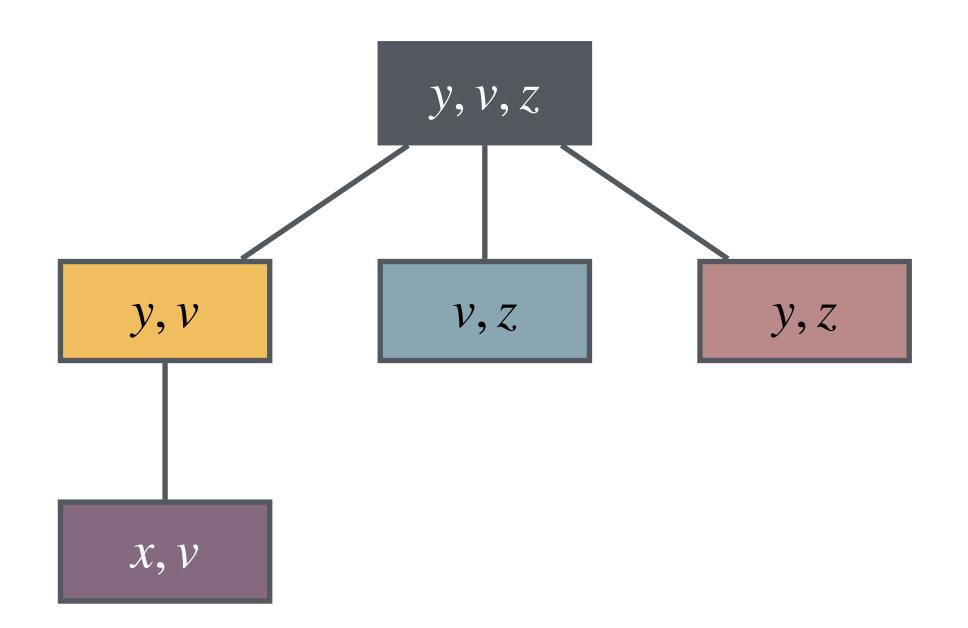
#### Yannakakis' Algorithm — for Enumeration

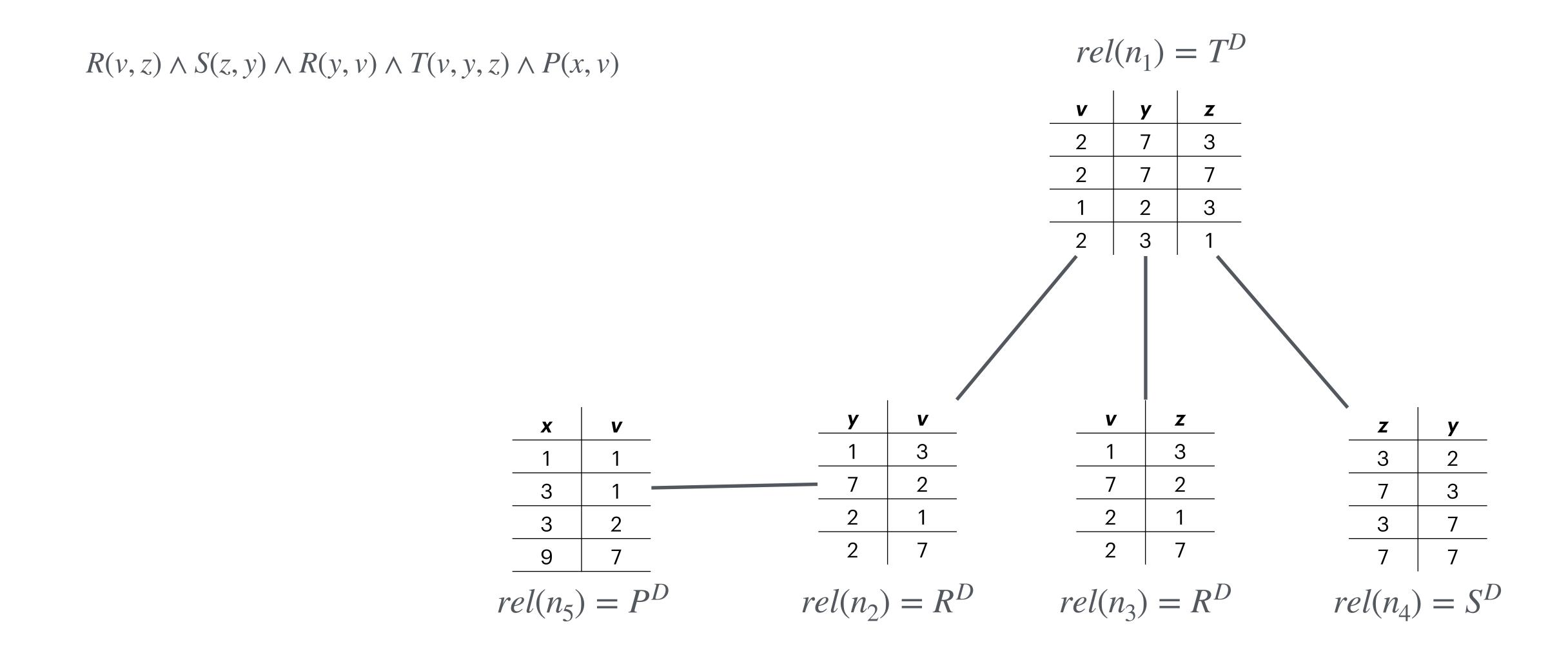
Let  $(T, \lambda)$  be a join tree of a CQ q and root it arbitrarily. Given a database D we can decide  $q(D) \neq \emptyset$  as follows:

- 1. Assign to rel(n) to each node n as before.
- 2. In a bottom-up traversal: update rel(n) to  $\left(rel(n) \ltimes rel(c_1)\right) \ltimes \cdots \ltimes rel(c_\ell)$  where  $c_1, \ldots, c_\ell$  are the children of n.
- 3. In a top-down traversal: update rel(n) to  $(rel(n) \ltimes rel(p))$  where p is the parent of n.
- 4. Every tuple left after this process will be part of an output, we can collect them one by one in a final bottom-up process.

$$q = \{ (v, x, y, z) \mid R(v, z) \land S(z, y) \land R(y, v) \land T(v, y, z) \land P(x, v) \}$$

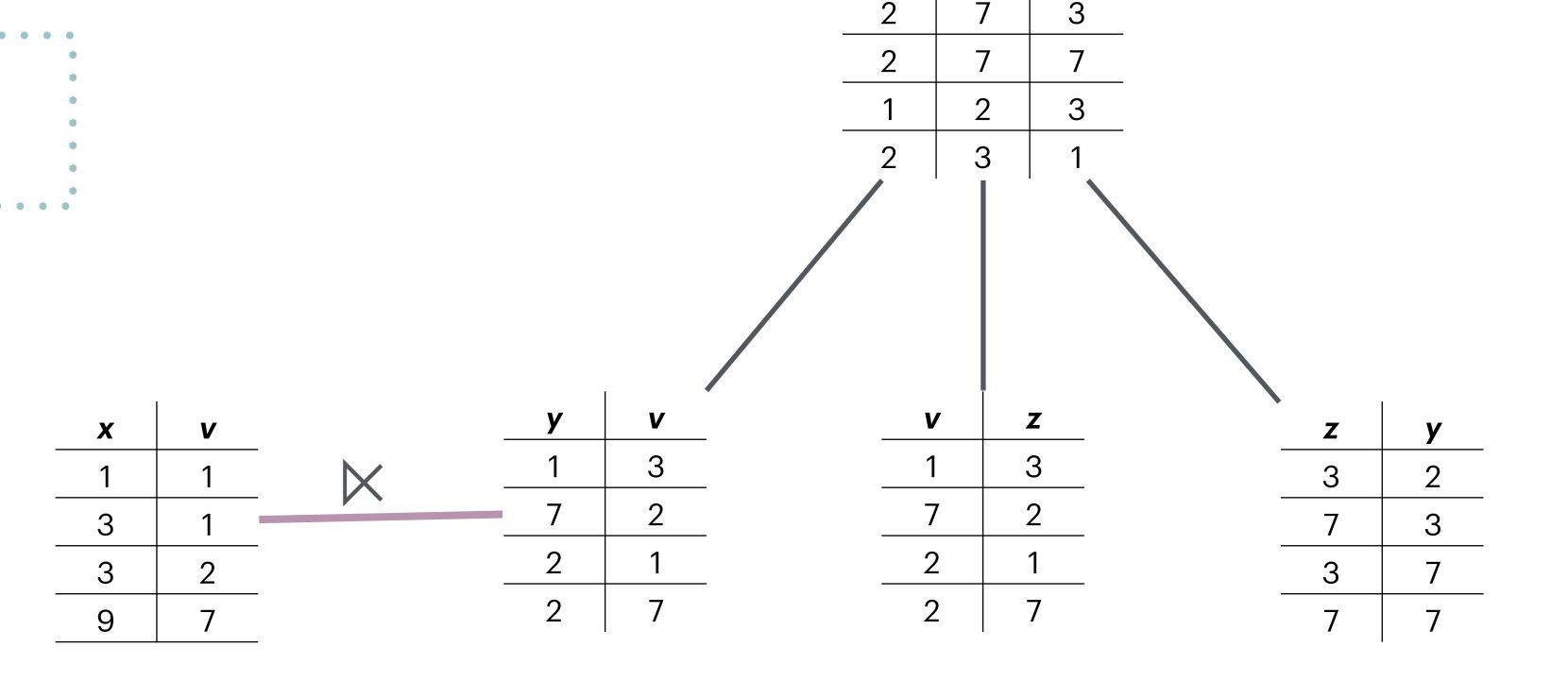






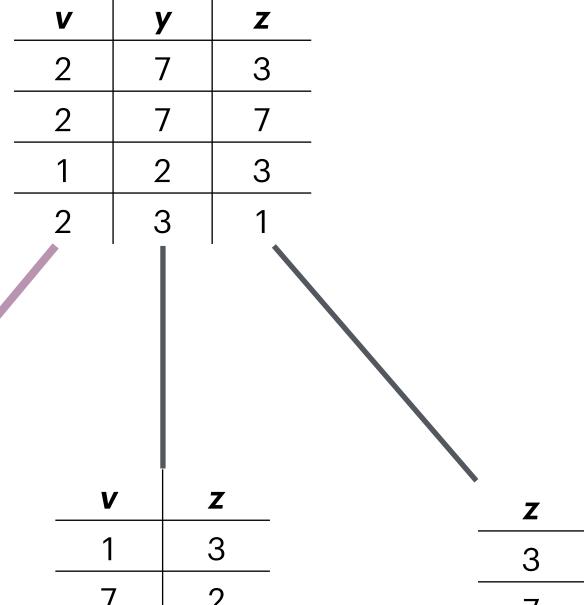
 $R(v,z) \wedge S(z,y) \wedge R(y,v) \wedge T(v,y,z) \wedge P(x,v)$ 

2. The bottom up semi-joins



 $R(v,z) \wedge S(z,y) \wedge R(y,v) \wedge T(v,y,z) \wedge P(x,v)$ 

2. The bottom up semi-joins



X	V
1	1
3	1
3	2
9	7

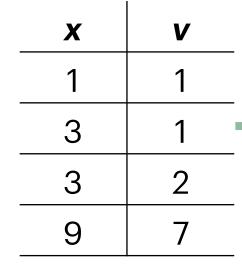
<b>y</b>	V
7	2
2	1
2	7

V	Z
1	3
7	2
2	1
2	7

 $R(v,z) \wedge S(z,y) \wedge R(y,v) \wedge T(v,y,z) \wedge P(x,v)$ 

2. The bottom up semi-joins





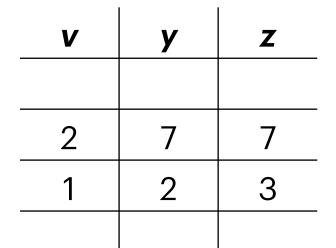
<b>y</b>	V
7	2
2	1
2	7

	Z
1	3
7	2
2	1
2	7

_	Z	У
_	3	2
	7	3
	3	7
_	7	7

 $R(v,z) \wedge S(z,y) \wedge R(y,v) \wedge T(v,y,z) \wedge P(x,v)$ 

2. The bottom up semi-joins



X	V
1	1
3	1
3	2
9	7

V
2
1
7

<b>v</b>	Z
1	3
7	2
2	1
2	7

_	Z	У
_	3	2
	7	3
	3	7
-	7	7

 $R(v,z) \wedge S(z,y) \wedge R(y,v) \wedge T(v,y,z) \wedge P(x,v)$ 

2. The bottom up semi-joins

V	у	Z
2	7	7
1	2	3

X	V
1	1
3	1
3	2
9	7
	·

<b>y</b>	V
7	2
2	1
2	7

V	Z
1	3
7	2
2	1
2	7

Z	У
3	2
7	3
3	7
7	7

 $R(v,z) \wedge S(z,y) \wedge R(y,v) \wedge T(v,y,z) \wedge P(x,v)$ 

3. The top-down semi-joins. Recall, only from parent to child.

V	y	Z
2	7	7
1	2	3

<b>V</b>
1
1
2
7

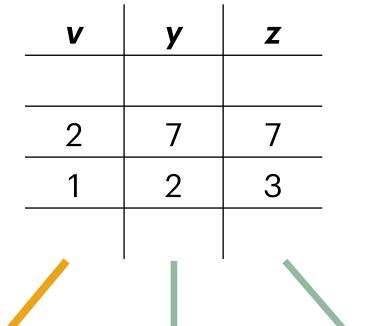
y	V
7	2
2	1
2	7

V	Z
1	3
7	2
2	1
2	7

_	Z	У
_	3	2
	7	3
	3	7
_	7	7

 $R(v,z) \wedge S(z,y) \wedge R(y,v) \wedge T(v,y,z) \wedge P(x,v)$ 

3. The top-down semi-joins.Recall, only from parent to child.



X	V
1	1
3	1
3	2
9	7
9	7

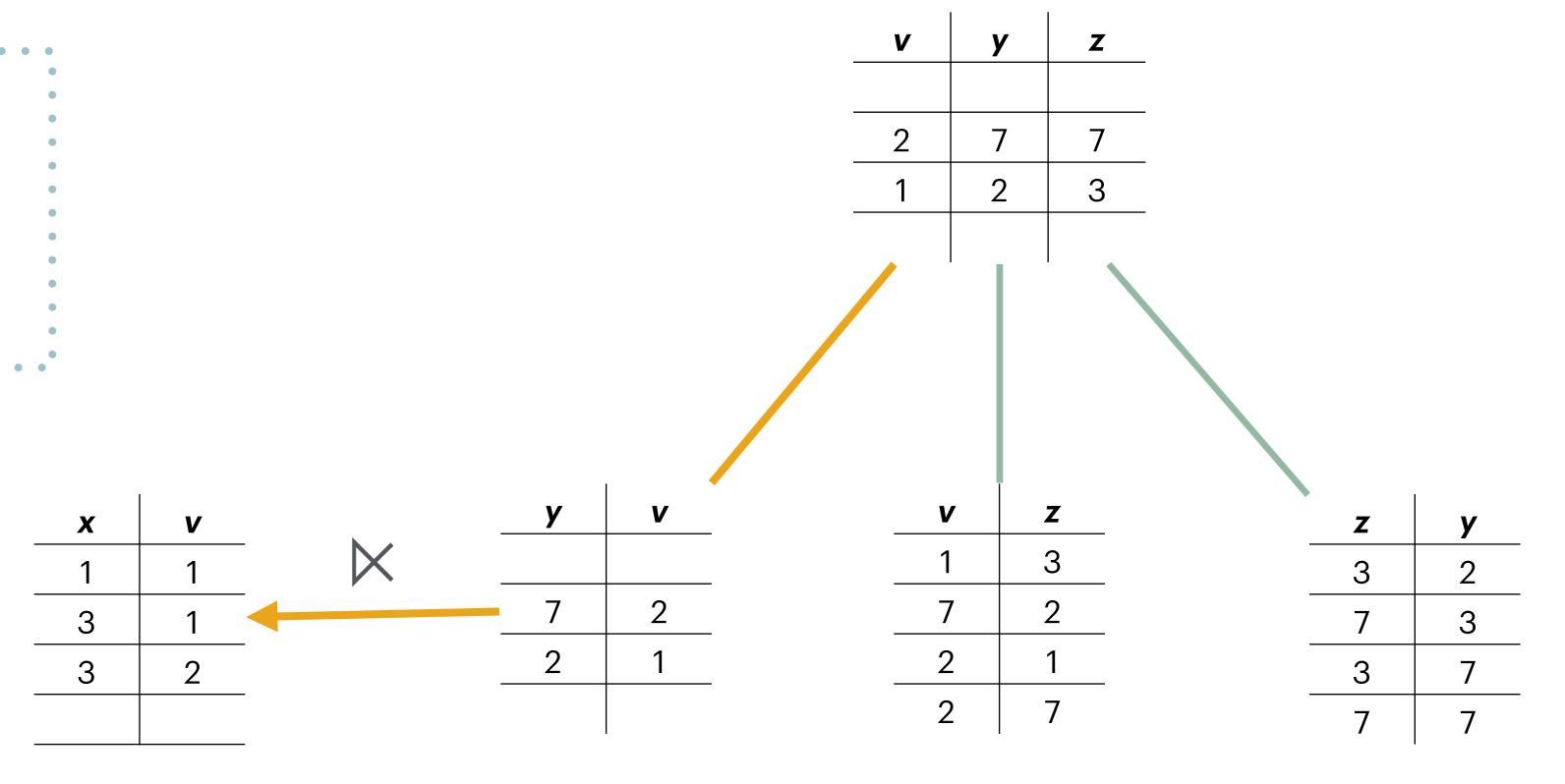
<b>y</b>	V
7	2
2	1

V	Z
1	3
7	2
2	1
2	7

 Z	У
 3	2
 7	3
3	7
7	7

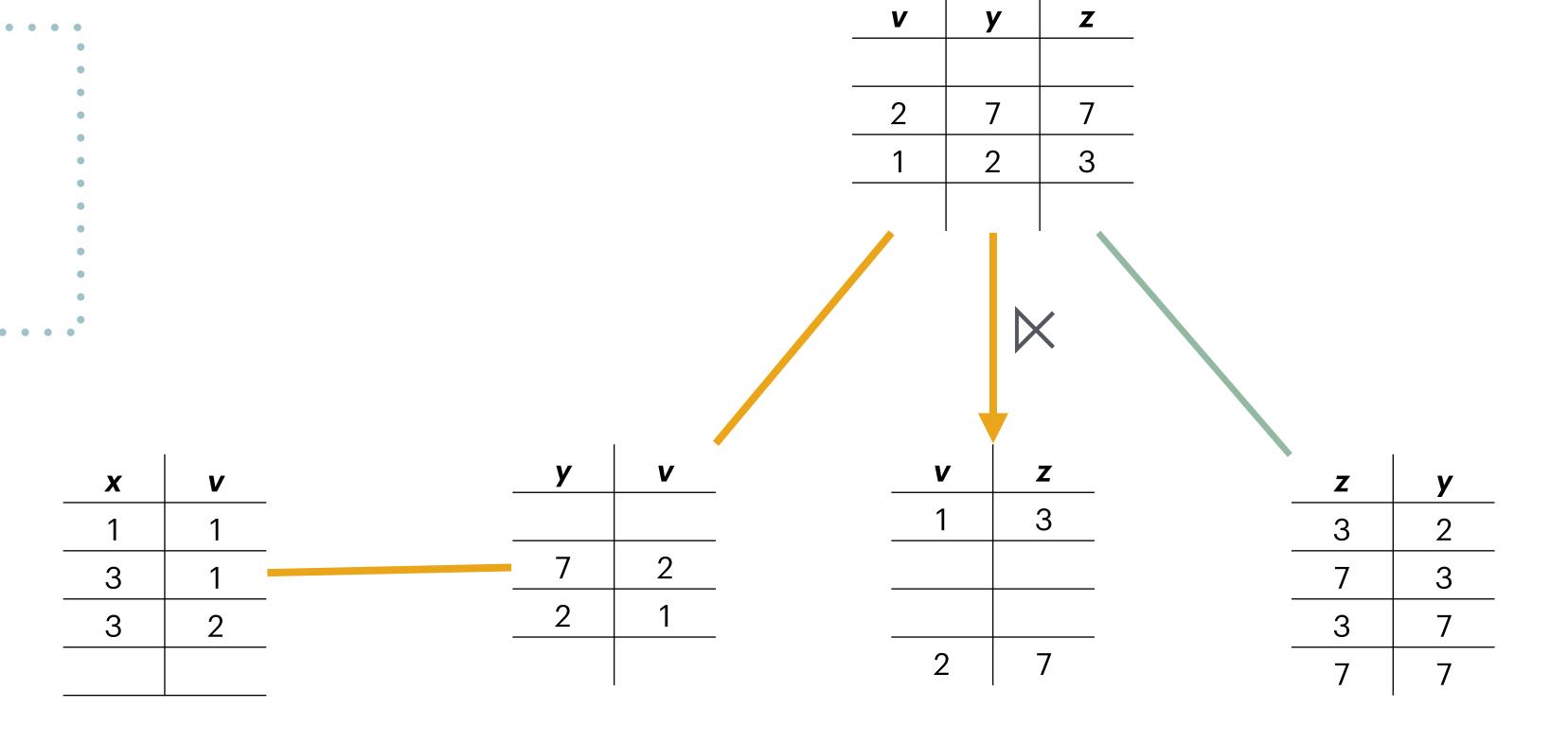
 $R(v,z) \wedge S(z,y) \wedge R(y,v) \wedge T(v,y,z) \wedge P(x,v)$ 

3. The top-down semi-joins.Recall, only from parent to child.



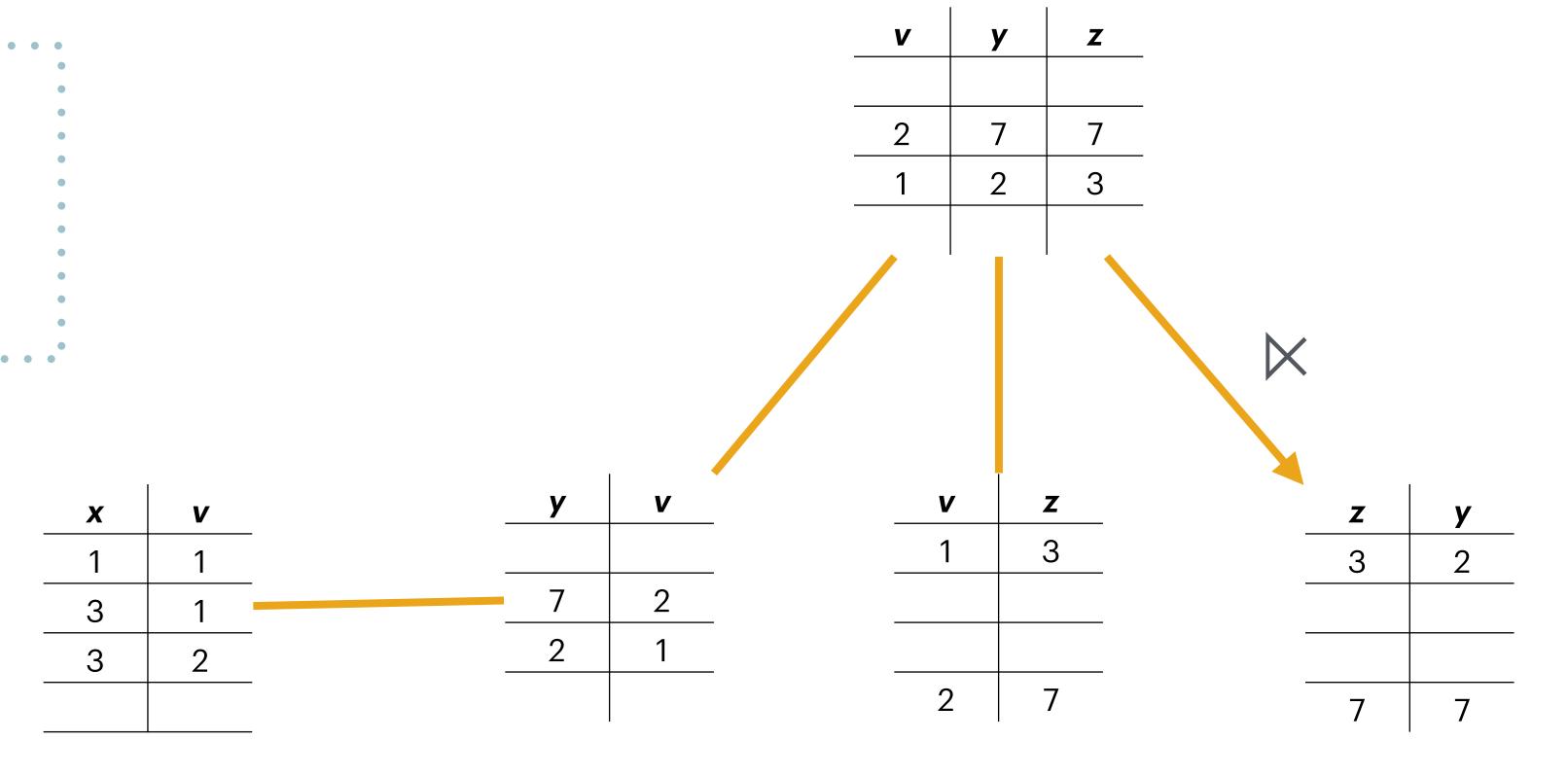
 $R(v,z) \wedge S(z,y) \wedge R(y,v) \wedge T(v,y,z) \wedge P(x,v)$ 

3. The top-down semi-joins. Recall, only from parent to child.



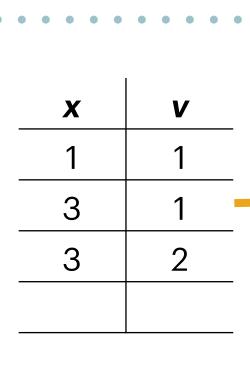
 $R(v,z) \wedge S(z,y) \wedge R(y,v) \wedge T(v,y,z) \wedge P(x,v)$ 

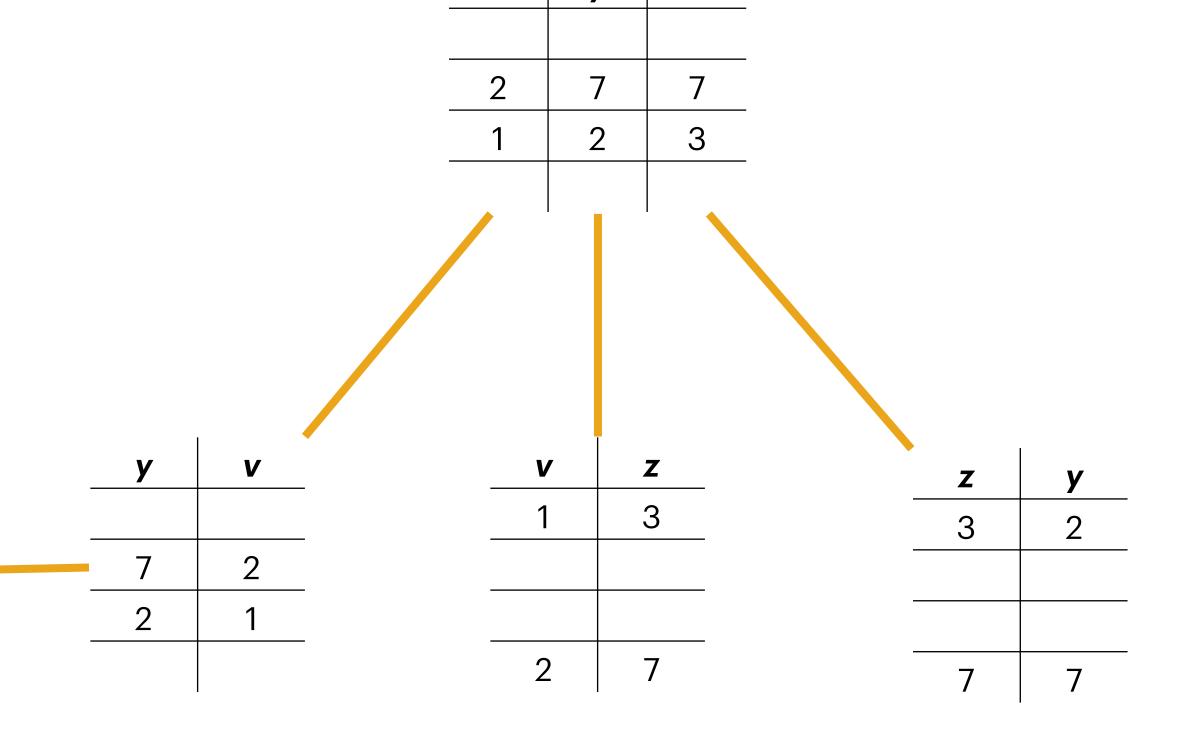
3. The top-down semi-joins.Recall, only from parent to child.



 $R(v,z) \wedge S(z,y) \wedge R(y,v) \wedge T(v,y,z) \wedge P(x,v)$ 

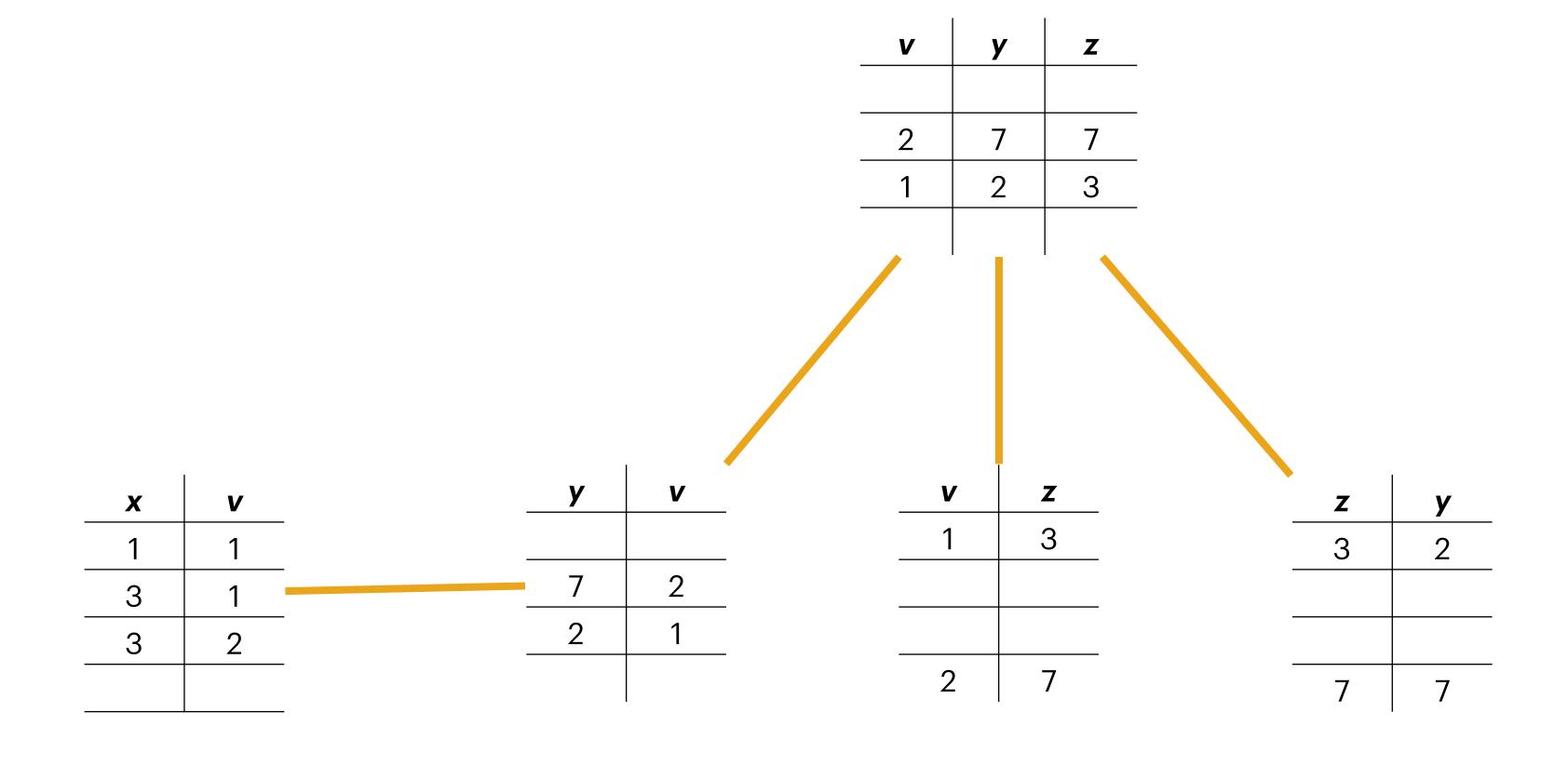
4. All tuples that are left-over are part of a solution. Joining the tables now creates no intermediate results that are not part of the output.





Answers q(D):

V	X	y	Z
2	3	7	7
1	3	2	3
1	1	2	3



### How Efficient is This?



After both semi-join phases of Yannakakis' Algorithm, we are left only with tuples that are part of answers. That means that no unnecessary effort is performed if we now join the left-over tuples to compute the answer in time  $\tilde{O}(|D| + |q(D)|)$  for a fixed  $\alpha$ -acyclic query.

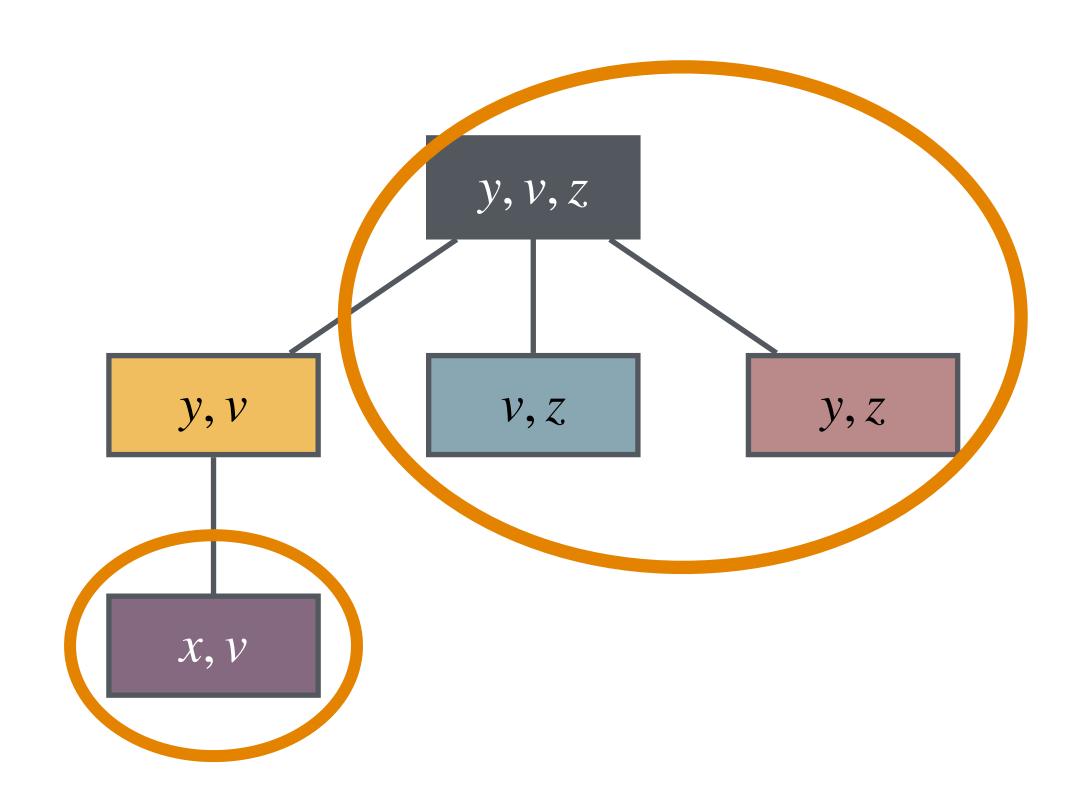
However, with our tree structure we can do something more refined. We can enumerate the answers with **small delay between answers**. That is, the algorithm produces tuples one after another, where the time between outputting each tuple does not depend on the database.

### Quantification Matters

Consider a version of our example query with existential quantification:

$$q = \{ (x, z) \mid \exists vy \ R(v, z) \land S(z, y) \land R(y, v) \land T(v, y, z) \land P(x, v) \}$$

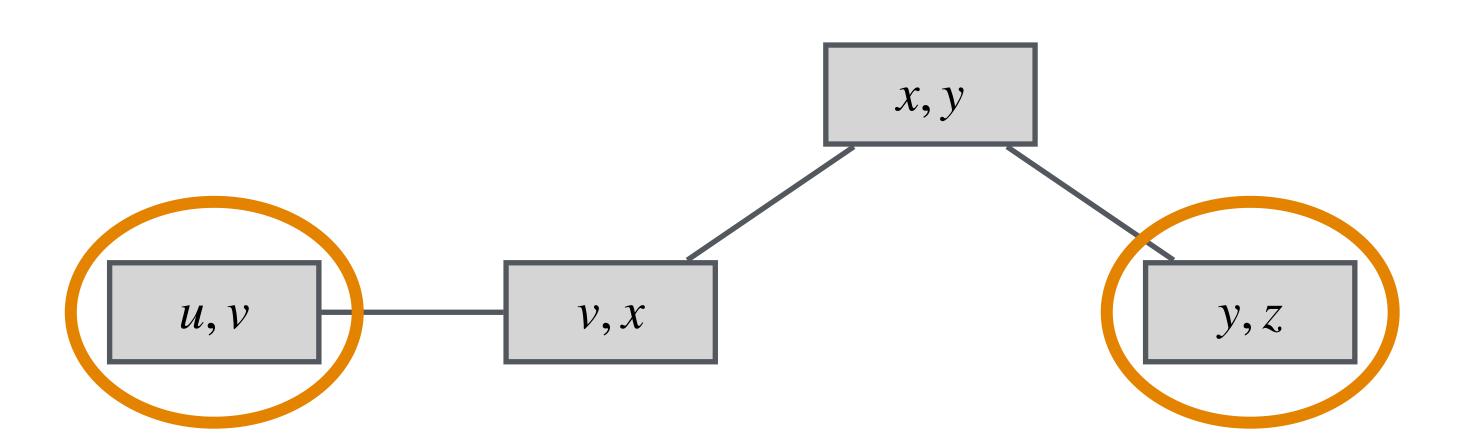
When collecting the final answers, we need to create intermediate results that extend to y, v to connect the two parts of the join tree that contain the output variables.



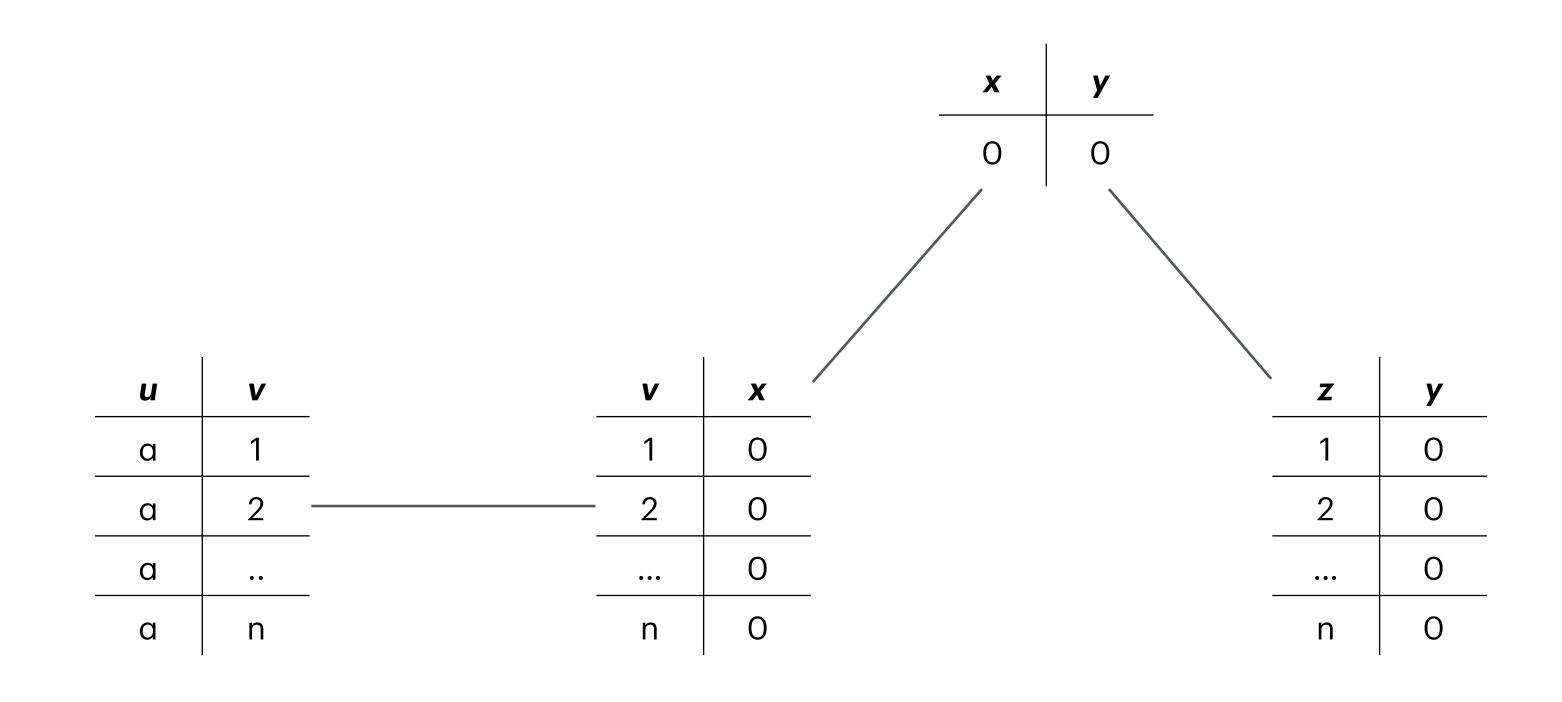
### Quantification Matters

Let's simplify our example to see how this can break the linear behaviour of Yannakakis.

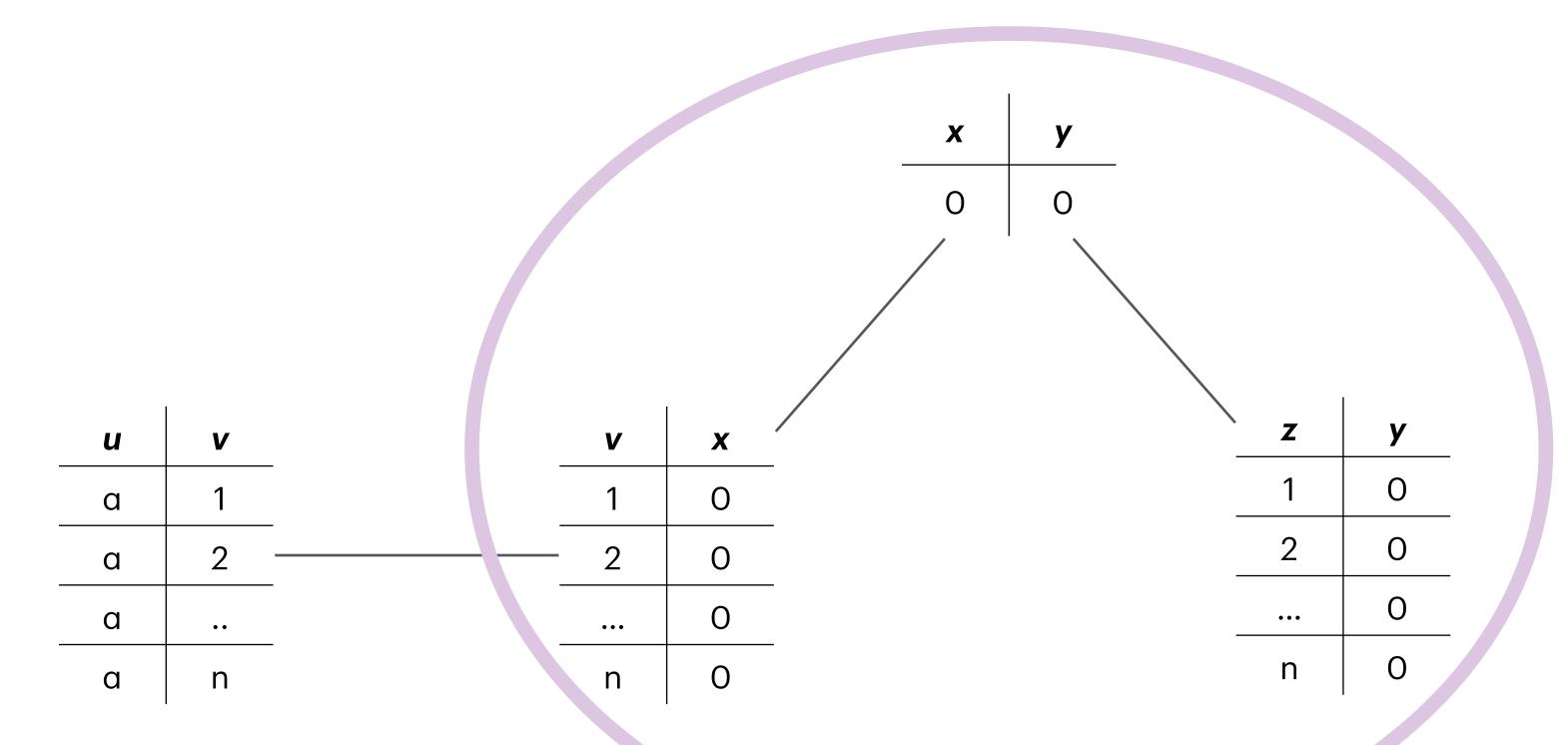
$$q = \{ (u, z) \mid \exists vxy \ T(u, v) \land R(v, x) \land S(x, y) \land R(z, y) \}$$



Assume these relations after the semi-join phase

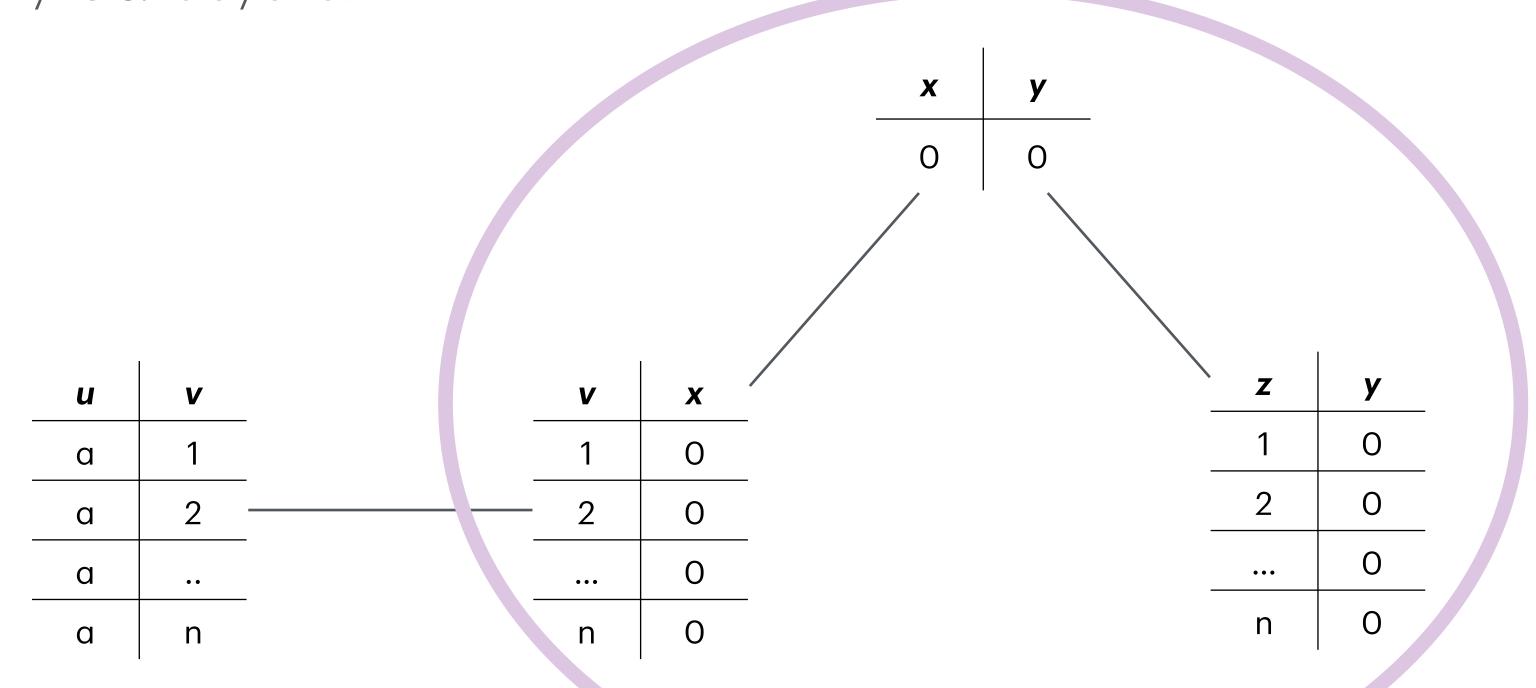


To see which assignments to  $\boldsymbol{u}$  and  $\boldsymbol{z}$  work together, we need to compute joins over variables that are unrelated to the output.



Joining these three tables creates a relation of size  $n^2$ !

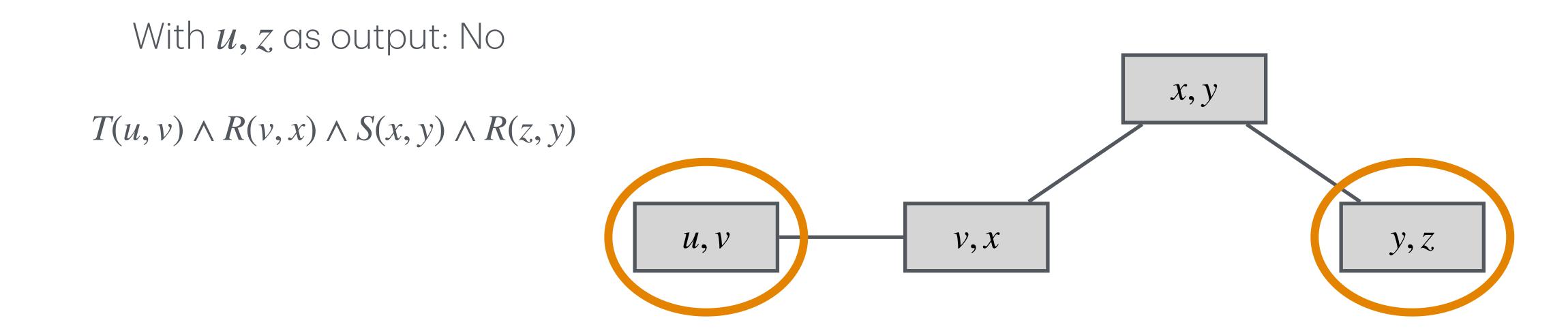
The actual result tuples are:  $(a,1),(a,2),\dots,(a,n)$ . Meaning |q(D)|=n, but creating the joins requires  $\Theta(n^2)$  time! The time bound  $\tilde{O}(|D|+|q(D)|)$  doesn't hold even though the query is  $\alpha$ -acyclic!



Joining these three tables creates a relation of size  $n^2$ !

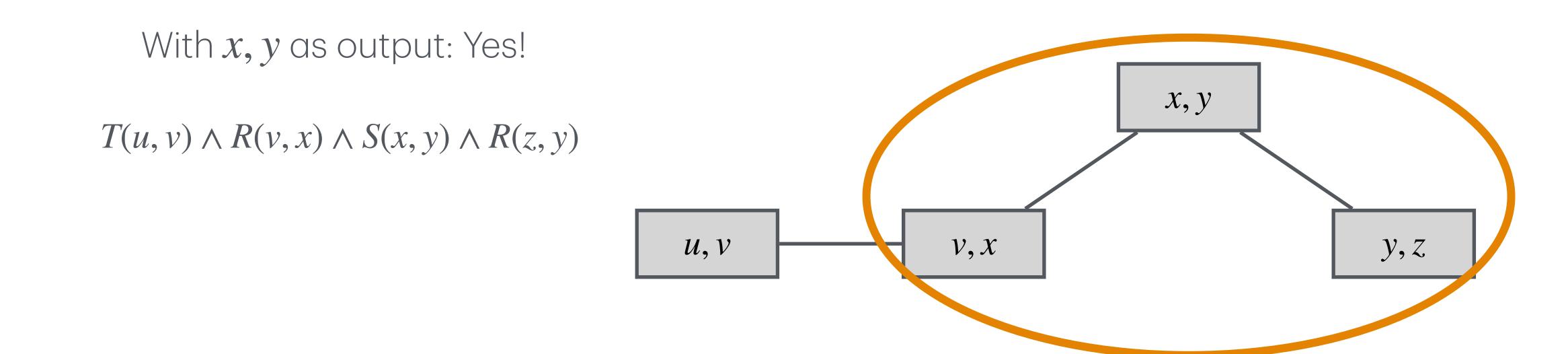
#### Quantification Details

- lacktriangle We can obtain the same guarantees (and in particular constant-delay enumeration) only on a limited number of  $\alpha$ -acyclic CQs when there is existential quantification!
- ◆ Intuitively, this requires the part of the join tree that contains the output variables to be connected:



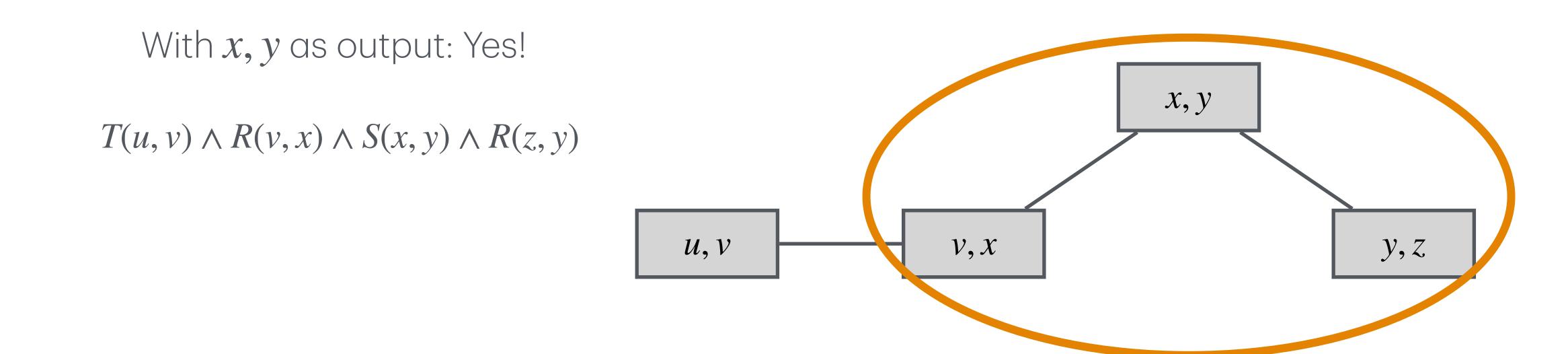
### Quantification Details

- lacktriangle We can obtain the same guarantees (and in particular constant-delay enumeration) only on a limited number of lpha-acyclic CQs when there is existential quantification!
- ◆ Intuitively, this requires the part of the join tree that contains the output variables to be connected:



### Quantification Details

- lacktriangle We can obtain the same guarantees (and in particular constant-delay enumeration) only on a limited number of lpha-acyclic CQs when there is existential quantification!
- ◆ Intuitively, this requires the part of the join tree that contains the output variables to be connected:



### Summary

- ◆ CQs are the core element of many data retrieval tasks.
- ◆ Unlike FO/RA queries, we can decide containment, equivalence and even minimality of CQs.
- ◆ Answering CQs is NP-complete, but acyclic queries can always be answered efficiently —> Structure of CQs is a key factor in their complexity.
- ◆ Yannakakis' Algorithm!
- ◆ Enumeration or counting versions bring their own challenges with them like quantification and more fine-grained notions of complexity.