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)

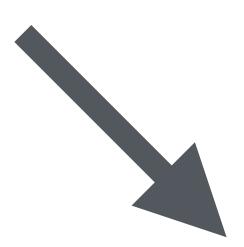
$$\left\{ \begin{array}{c|c} \bar{y} & \exists \bar{z} \ R_1(\bar{x}_1) \land R_2(\bar{x}_2) \land \cdots \land R_k(\bar{x}_k) \end{array} \right\}$$

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

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

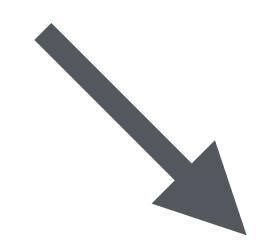
Conjunctive queries correspond so-called join queries in RA. That is, RA queries that only use projection, renaming, selection, and joins.

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 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

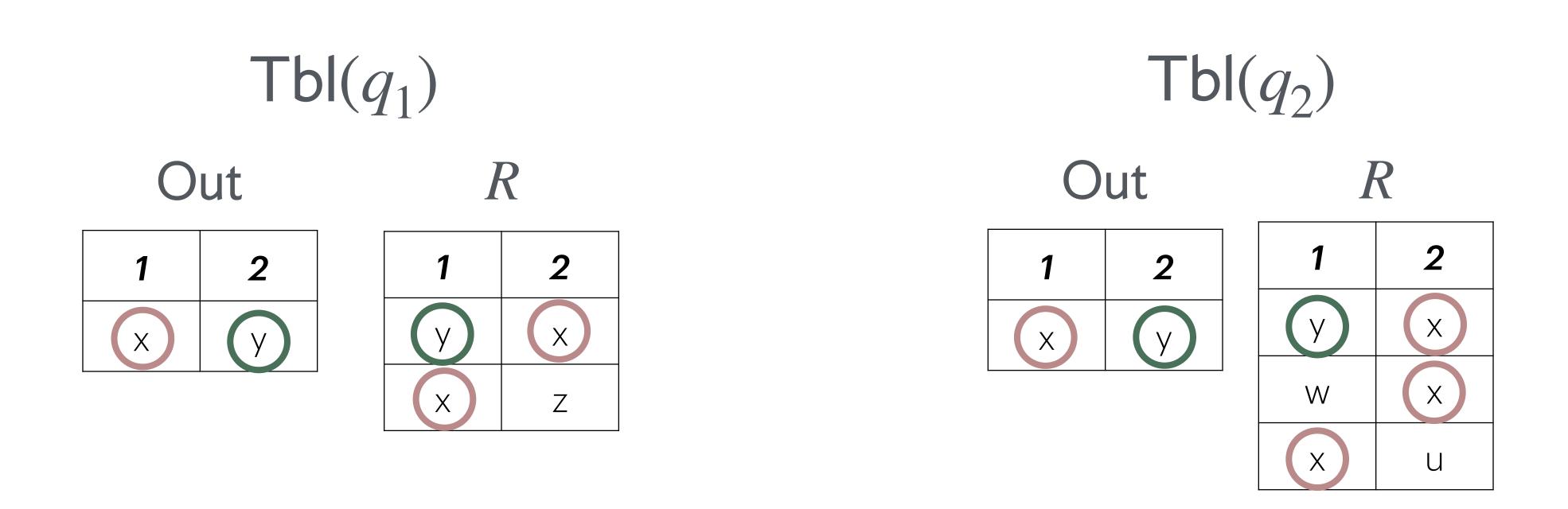
 $\mathsf{Tbl}(q_2)$

Out

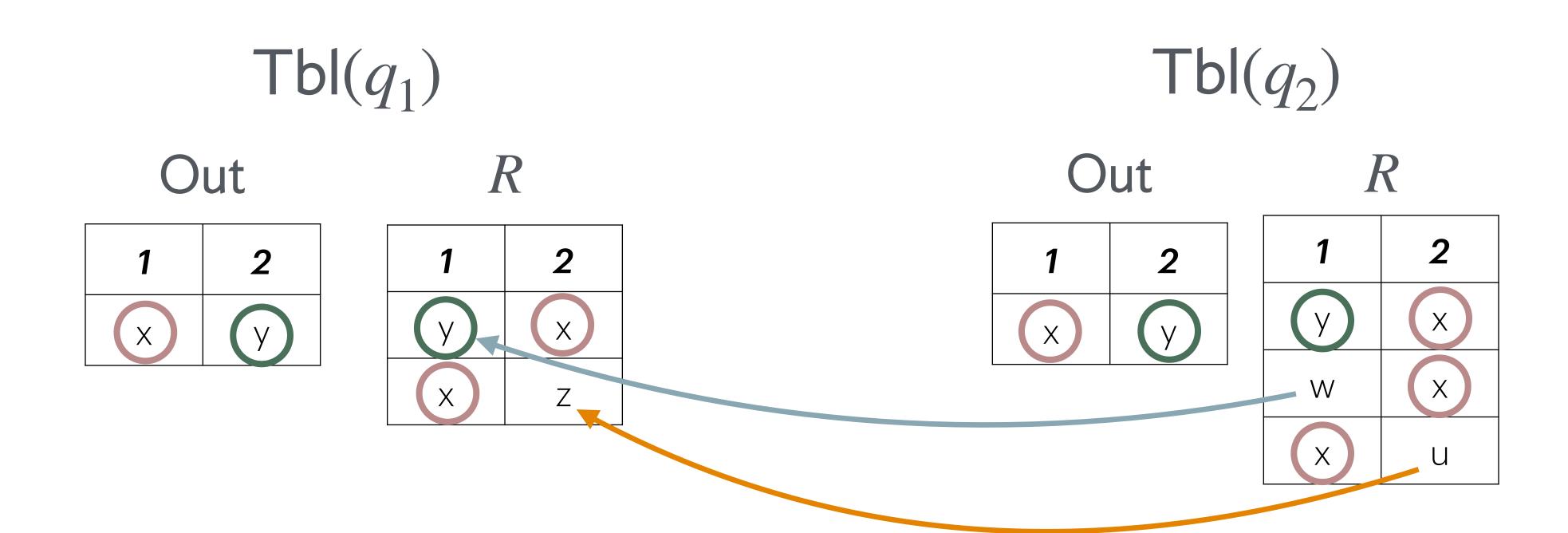
1	2
X	У

R

2
X
X
u

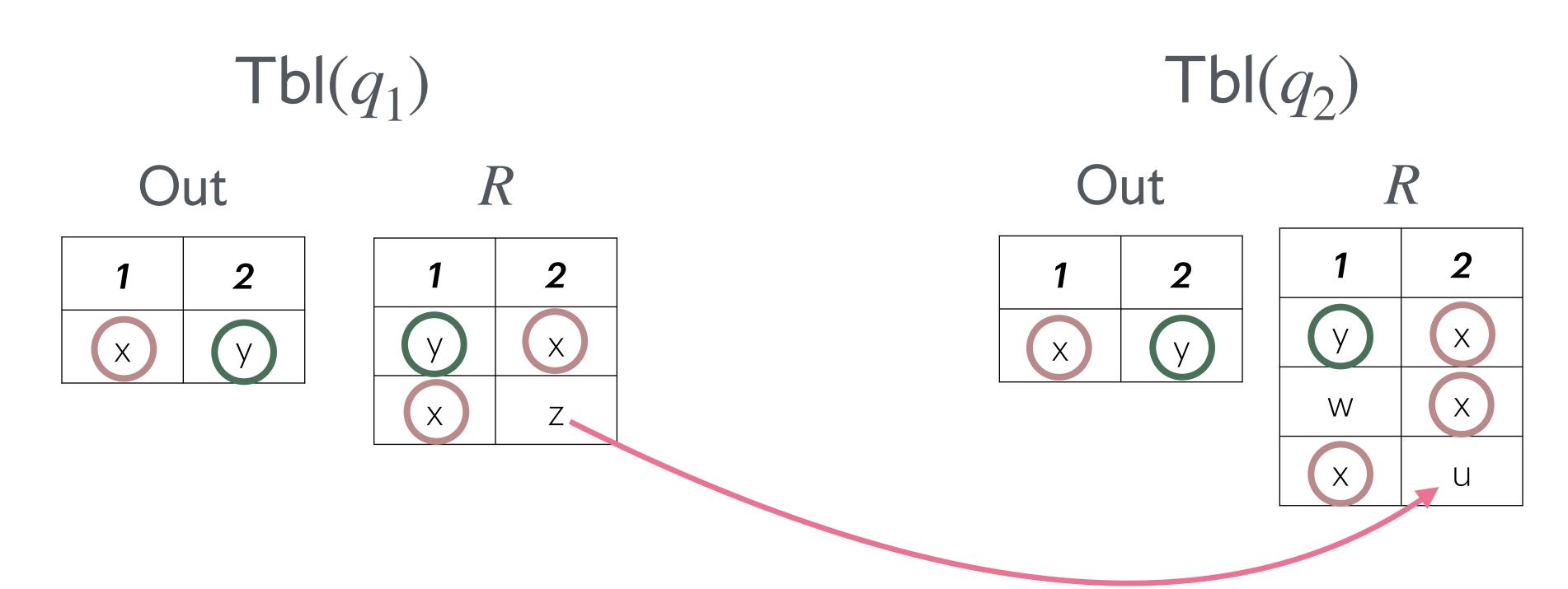


Since Out has only one tuple, a homomorphism h must necessarily have h(x) = x and h(y) = y.



Since 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 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

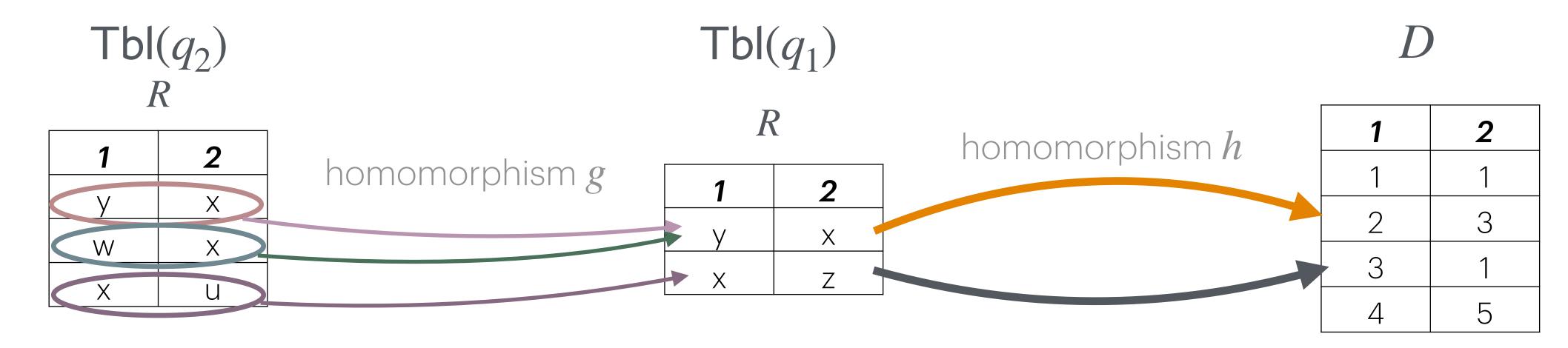
$$q_1 := \{(x, y) \mid \exists z \, R(y, x) \land R(x, z)\}$$

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

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} .

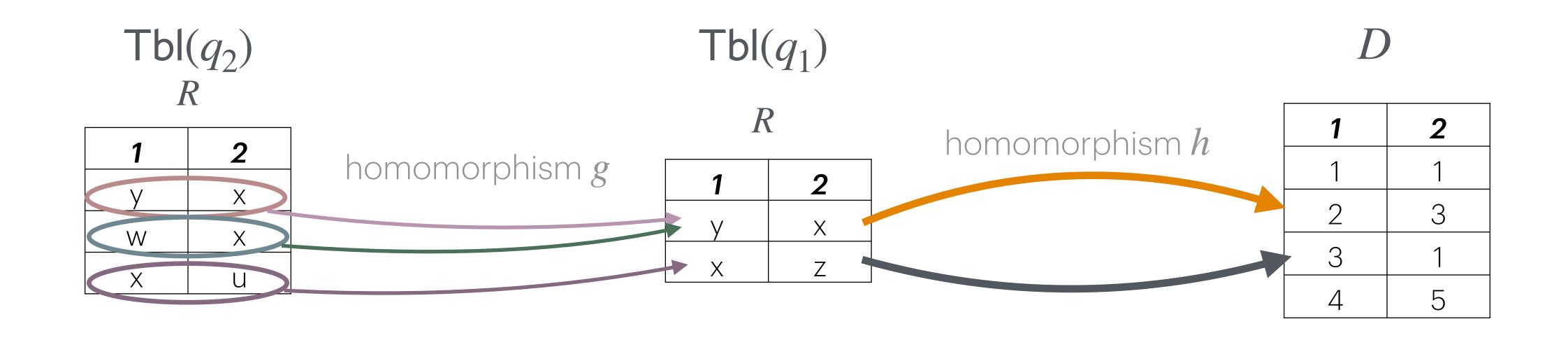
Let g be the homomorphism from $\mathsf{Tbl}(q_2)$ to $\mathsf{Tbl}(q_1)$. It is not hard to see that $h \circ g$ is homomorphism from $\mathsf{Tbl}^*(q_2)$ to D that also maps the output of variables of q_2 to \bar{c} .

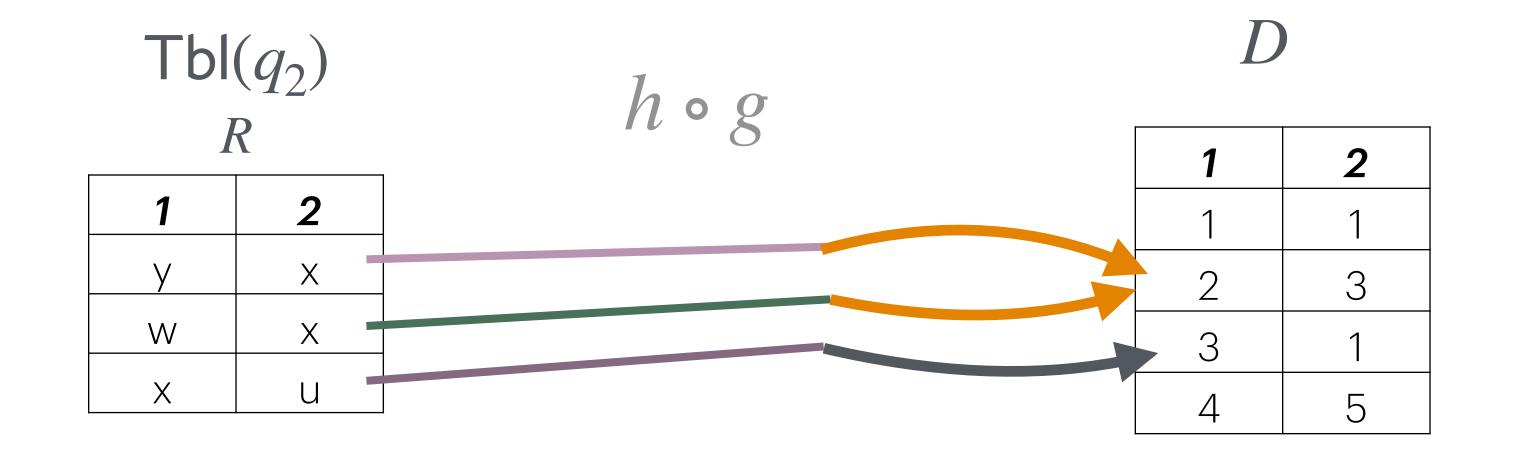


Example queries from before

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

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





Still a homomorphism!

Output variables map to the same values in D!

Example queries from before

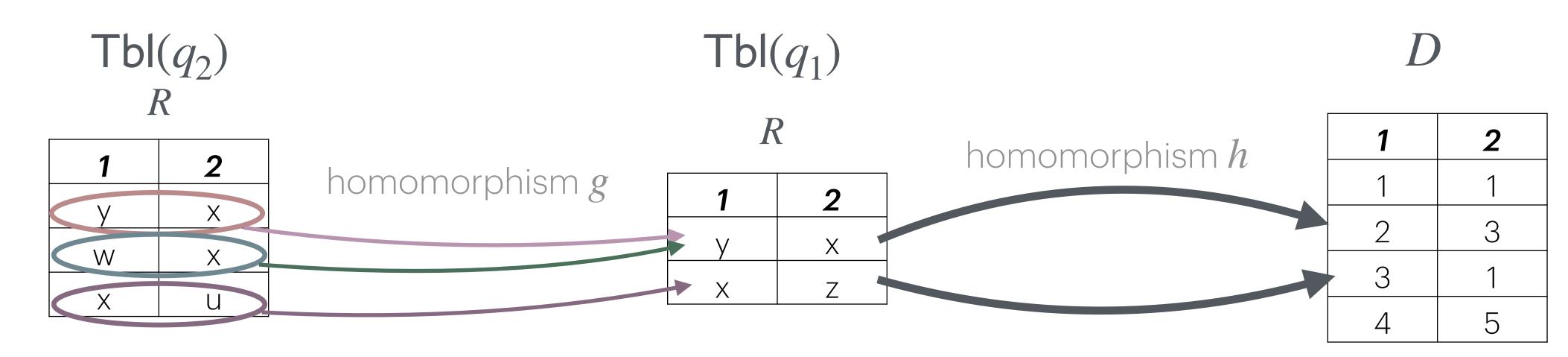
$$q_1 := \{(x, y) \mid \exists z \, R(y, x) \land R(x, z)\}$$

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

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} .



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?

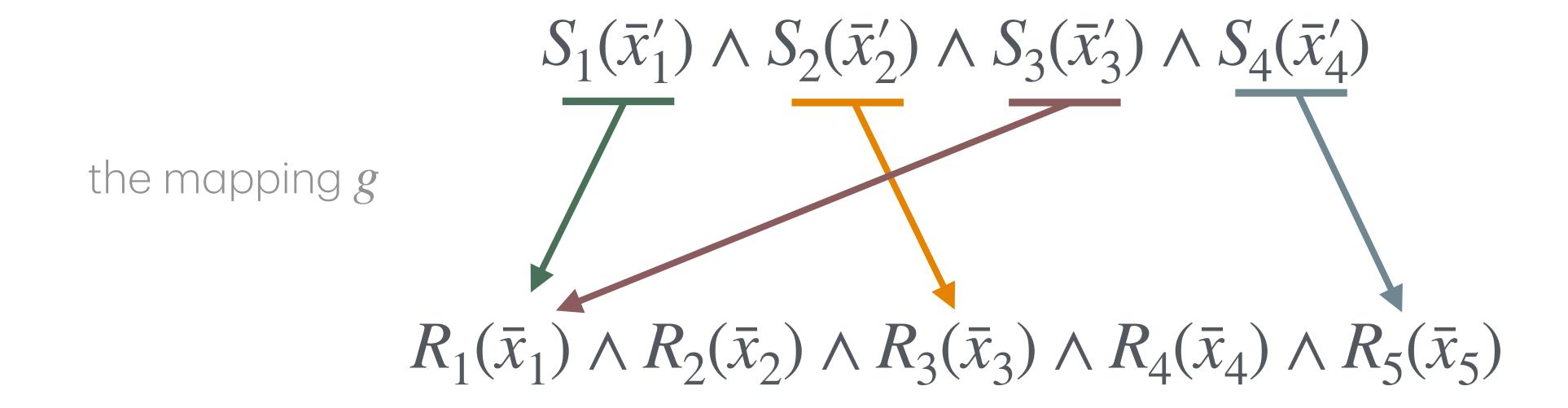
To minimise CQ q, it is enough to check only those queries obtained by deleting atoms from q!

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:

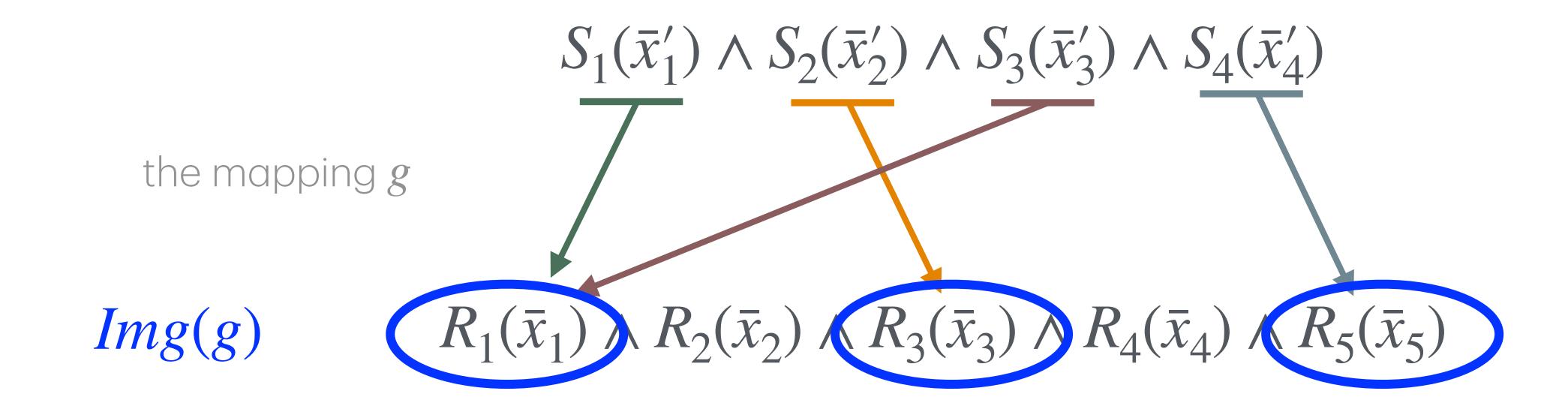
$$f: \mathsf{Tbl}(q) \to \mathsf{Tbl}(q')$$
 and $g: \mathsf{Tbl}(q') \to \mathsf{Tbl}(q)$

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 $\bar{x}_{i_{\alpha}} = h(\bar{x}'_{\alpha})$.

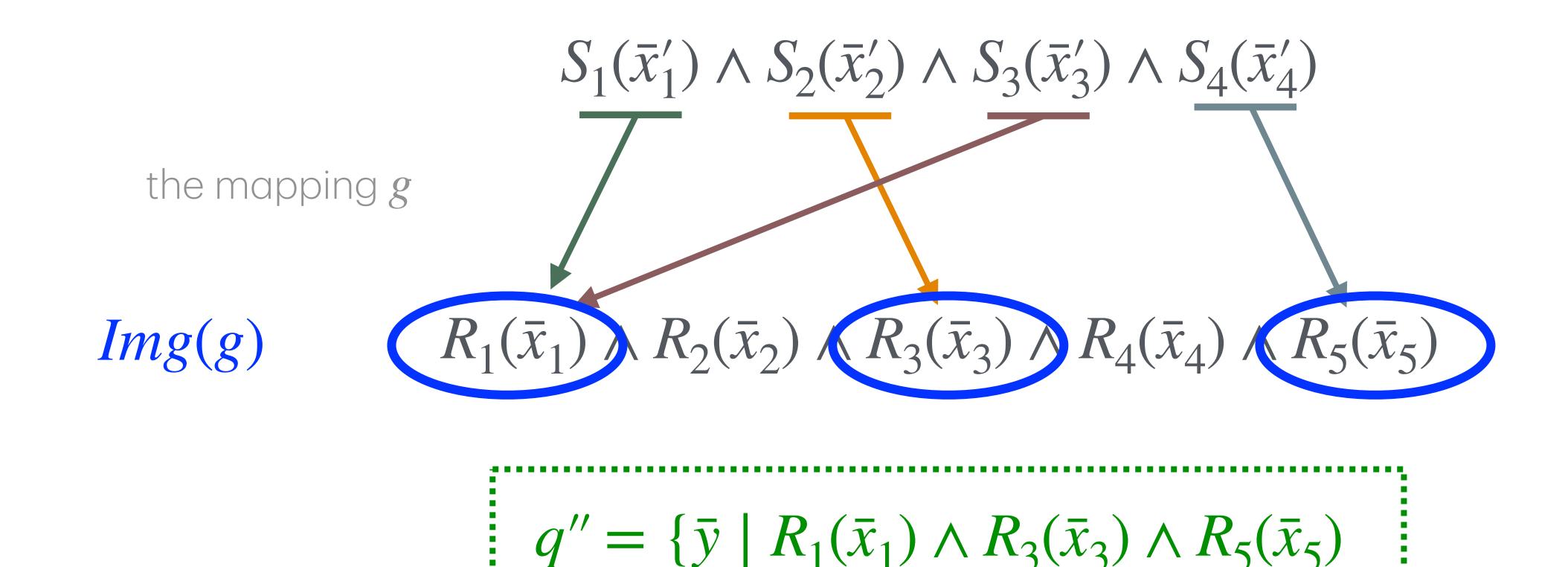


We have that $g: \mathrm{Tbl}(q') \to \mathrm{Tbl}(q)$ maps every $S_{\alpha}(\bar{x}'_{\alpha}) \in \mathrm{Tbl}(q')$ into some $R_{i_{\alpha}}(\bar{x}_{i_{\alpha}}) \in \mathrm{Tbl}(q)$ with $S_{\alpha} = R_{i_{\alpha}}$ and $\bar{x}_{i_{\alpha}} = h(\bar{x}'_{\alpha})$.

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.

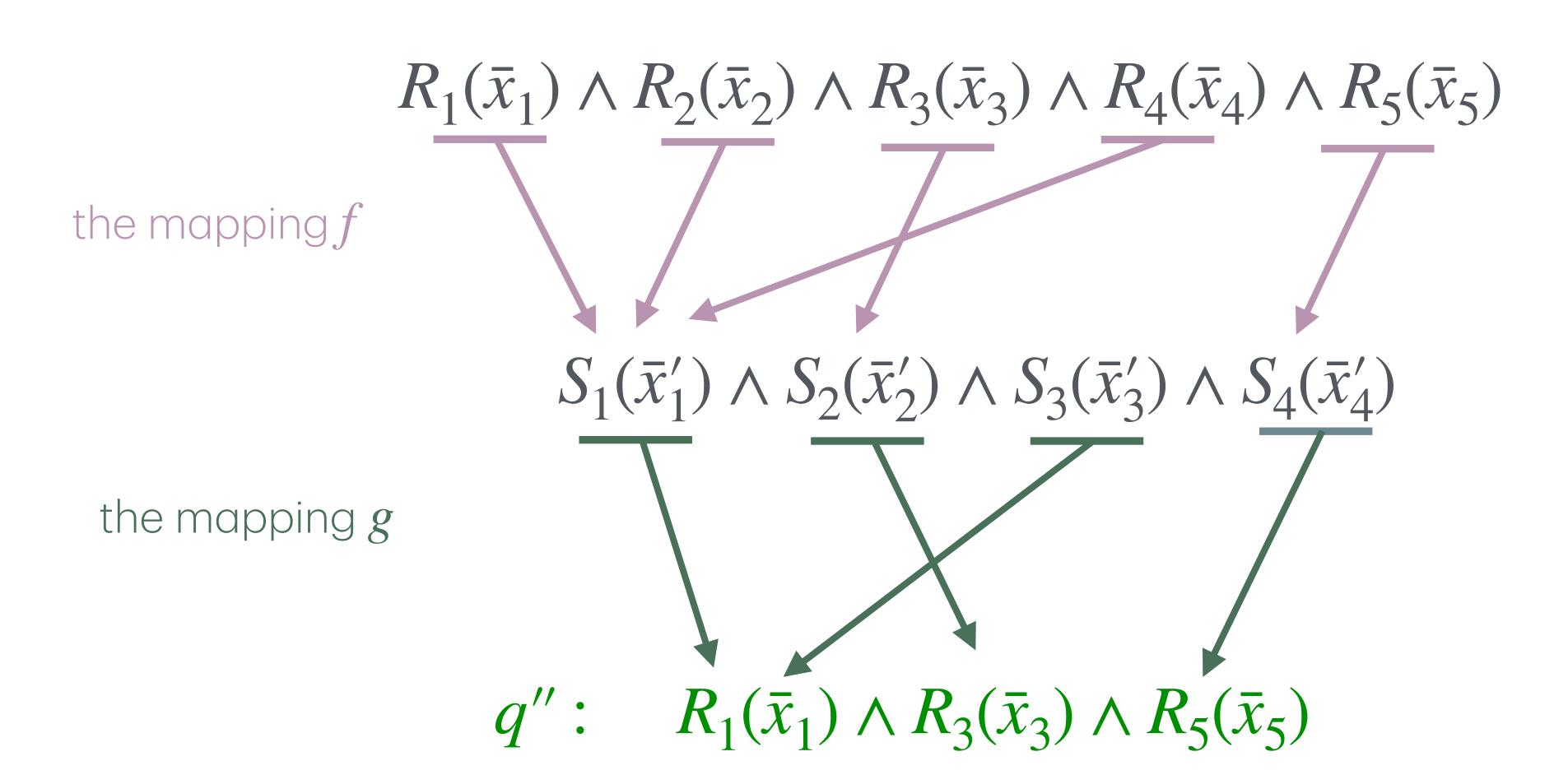
$$q'': R_{1}(\bar{x}_{1}) \wedge R_{3}(\bar{x}_{3}) \wedge R_{5}(\bar{x}_{5})$$

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

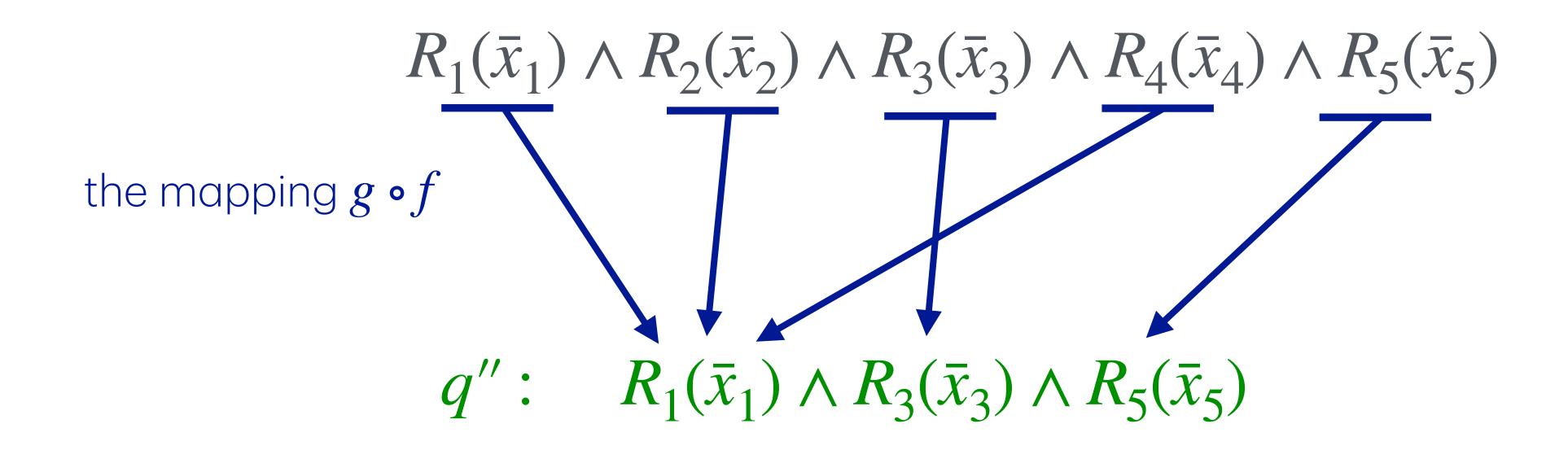
We use the Homomorphism Theorem to show that q'' is also equivalent to q:

- ♦ 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

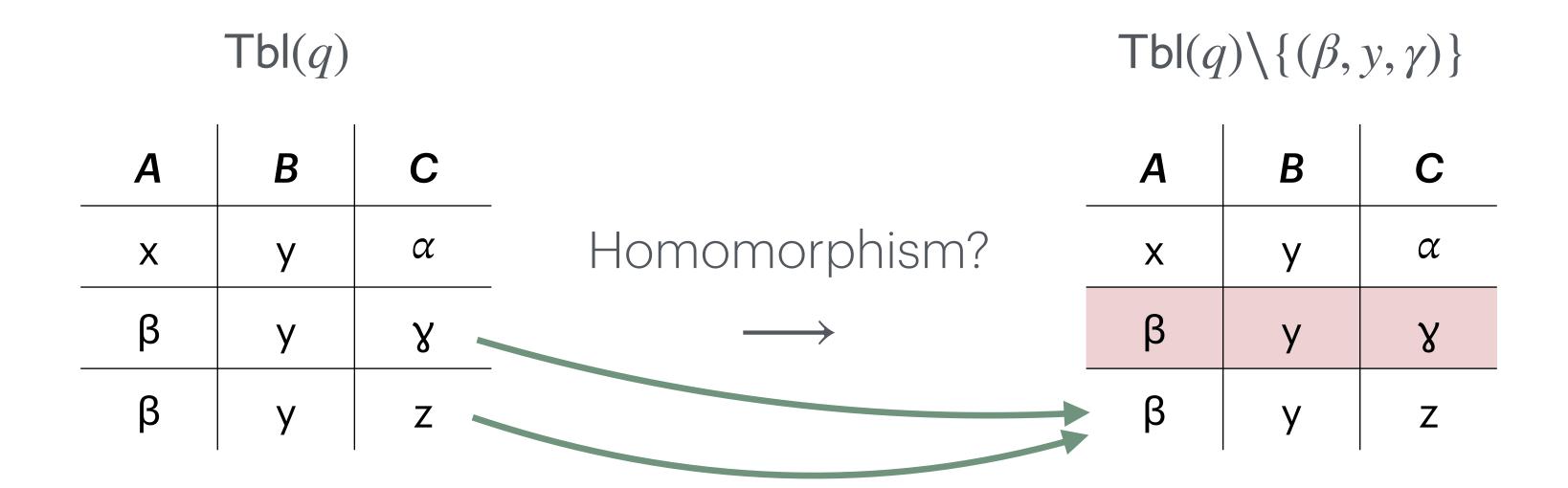
$$q = \{(x, y, z) \mid \exists \alpha \beta \gamma \ R(x, y, \alpha) \land R(\beta, y, \gamma) \land R(\beta, y, z) \}$$

	Tbl(q)			Tbl(q	$(y)\setminus\{(x,y)\}$	$y, \alpha)$
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, β map to themselves and $h(\gamma) = z$

CQ Minimisation Example

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

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 $q(D) \neq \emptyset$?

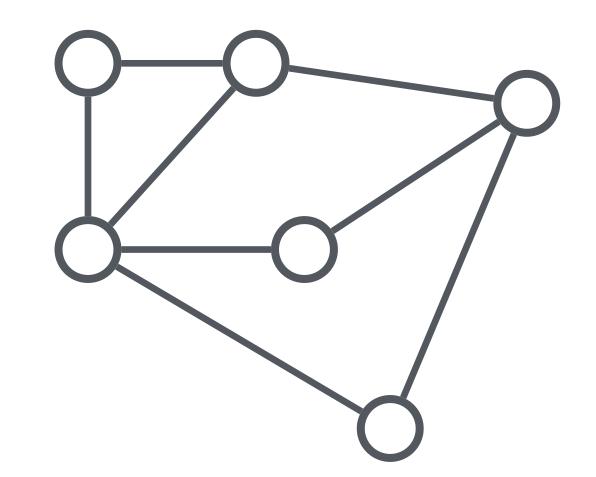
If there is any homomorphism from $\mathsf{Tbl}^*(q)$ to D.

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

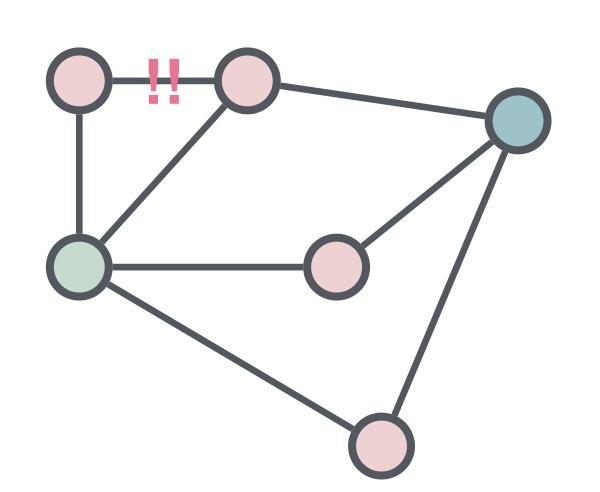
NP-Hardness

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

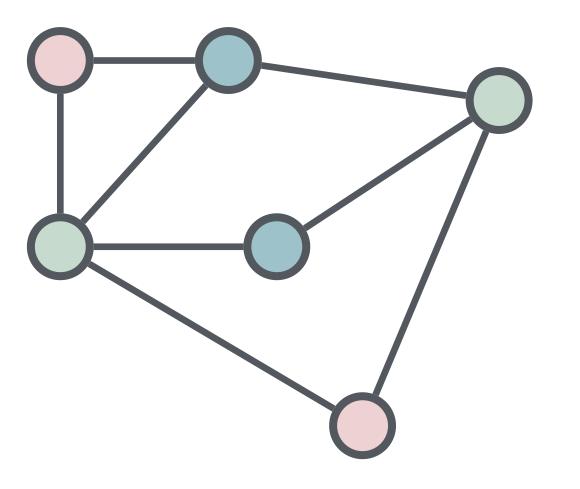
That is, can we color the vertices of G with red, green, and blue such that no edge is between two vertices of the same colour?



Not a 3-colouring

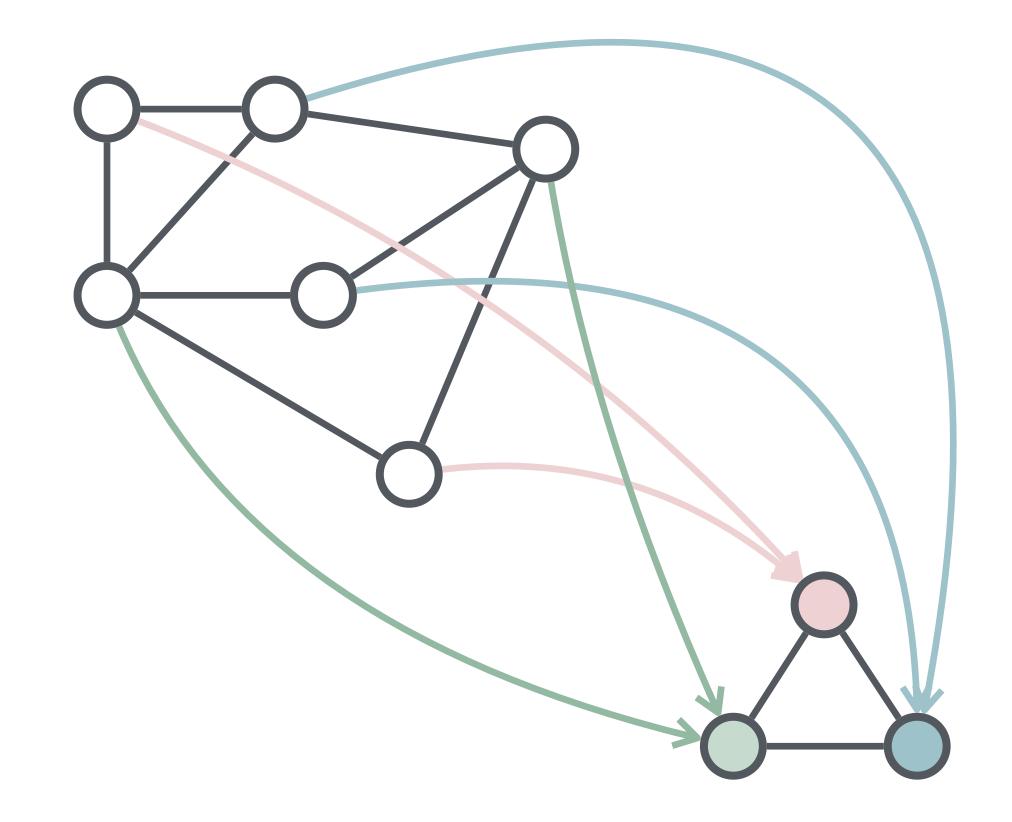


Valid 3-colouring



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.



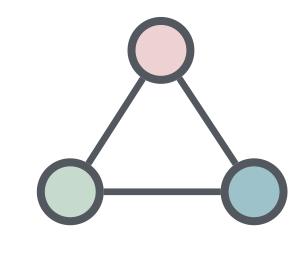
NP-Hardness

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.