

### 3 Graph Query Language Semantics

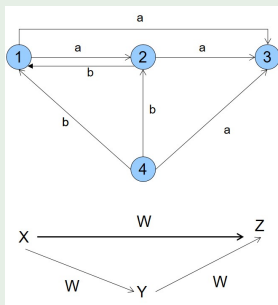
#### Graph Query Languages: SPARQL, Cypher, and Gremlin.

Most widely used query languages in practice but offer significant differences:

- SPARQL operates over RDF graphs, i.e. edge-labelled graphs;
- Cypher is designed to operate over property graphs;
- Gremlin is more imperative in nature than the other two, geared more towards graph traversal than graph pattern matching.

Our interest: semantic issues of graph pattern matching and graph traversal

Given an ELG  $G$  and query  $Q$ :



Two matches:

- $h_1 = \{X \rightarrow 1, Y \rightarrow 2, Z \rightarrow 3, W \rightarrow A\}$
- $h_2 = \{X \rightarrow 4, Y \rightarrow 2, Z \rightarrow 1, W \rightarrow B\}$

## ELG representation by relation $edge(from, label, to)$

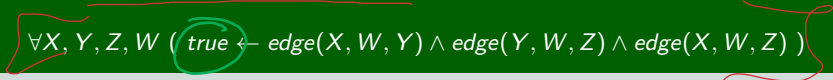
Graph  $G$  as instance of ternary relation  $edge$ :



<i>from</i>	<i>label</i>	<i>to</i>
1	<i>a</i>	2
2	<i>a</i>	3
1	<i>a</i>	3
4	<i>a</i>	3
4	<i>b</i>	2
2	<i>b</i>	1
4	<i>b</i>	1

Query  $Q$  as boolean Conjunctive Query (CQ):

$\forall X, Y, Z, W ( \text{true} \leftarrow edge(X, W, Y) \wedge edge(Y, W, Z) \wedge edge(X, W, Z) )$



# 3.1 Conjunctive Queries

## Definition

A *conjunctive Query*  $Q$  over a database schema  $\mathcal{R}$  is given as

$$ans(\vec{U}) \leftarrow R_1(\vec{U}_1), \dots, R_n(\vec{U}_n),$$

such that for  $1 \leq i \leq n$

- $R_i$  a relation name in  $\mathcal{R}$  and
- $\vec{U}$  and  $\vec{U}_i$  vectors of variables and constants;
- any variable appearing in  $\vec{U}$  appears also in some  $\vec{U}_i$ .
- Left to  $\leftarrow$  is the *head* of the query, and to the right there is the *body*. The atoms in the body are also called *subgoals*.

### Example 1 - quantifiers and connectives are explicitly given

$$\forall X, Y, Z, W \ (ans(W) \leftarrow edge(X, W, Y) \wedge edge(Y, W, Z) \wedge edge(X, W, Z) )$$

### Example 2 - quantifiers and connectives are implicitly given

*Sales(Part, Supplier, Customer),*  
*Part(PName, Type),*  
*Cust(CName, CAddr),*  
*Supp(SName, SAddr).*

$Q : \quad ans(T) \leftarrow [Sales(P, S, C), Part(P, T), Cust(C, A), Supp(S, A)]$

$$ans(\vec{U}) \leftarrow R_1(\vec{U}_1), \dots, R_n(\vec{U}_n).$$

## Answer

- The set of answers  $Q$  w.r.t. an instance  $\mathcal{I}$  of the given relations is denoted  $Q(\mathcal{I})$ .
- If there is a substitution (match, mapping)  $\sigma$  from the variables in  $\vec{U}_1, \dots, \vec{U}_n$  to the constants in  $dom$ , such that  $\sigma(R_1(\vec{U}_1)), \dots, \sigma(R_n(\vec{U}_n)) \in \mathcal{I}$  then by applying the same substitution  $\sigma$  to  $\vec{U}$ , we say that  $\sigma(ans(\vec{U}))$  is an answer in  $Q(\mathcal{I})$ .
- Substitutions are functions - a constant is mapped into itself.

## Problemes

Let  $Q, Q_1, Q_2$  be conjunctive queries.

Containment:  $Q_1 \sqsubseteq Q_2$ , i.e.,  $Q_1(I) \subseteq Q_2(I)$  for any instance  $I$ .

Equivalence:  $Q_1 \equiv Q_2$ , i.e.,  $Q_1 \sqsubseteq Q_2$  and  $Q_2 \sqsubseteq Q_1$ .

Minimization: Given  $Q_1$ , construct an equivalent query  $Q_2$ , which has as most as many subgoals in its body as  $Q_1$  and is minimal in the sense, that any query  $Q_3$  being equivalent to  $Q_2$  has at least as many subgoals in the body as  $Q_2$ .

$Q_2$  is called *minimal*.

## Example 1

Relation *edge*:

<i>from</i>	<i>label</i>	<i>to</i>
1	<i>a</i>	2
2	<i>a</i>	3
1	<i>a</i>	3
4	<i>a</i>	3
4	<i>b</i>	2
2	<i>b</i>	1
4	<i>b</i>	1

Containment relationship?

$Q :$        $ans(X, Z) \leftarrow edge(X, W, Y), edge(Y, W, Z), edge(X, W, Z)$

$Q' :$        $ans(X, Z) \leftarrow edge(X, W, Y), edge(Y, W, Z), edge(X, W, Z),$   
                   $edge(X', W', Y), edge(Y, W', Z'), edge(X', W', Z')$

$Q'' :$        $ans(X, X') \leftarrow edge(X, W, Y), edge(Y, W, Z), edge(X, W, Z),$   
                   $edge(X', W', Y), edge(Y, W', Z'), edge(X', W', Z')$



## Example 2

Relation *edge*:

<i>from</i>	<i>label</i>	<i>to</i>
1	<i>a</i>	2
2	<i>a</i>	3
1	<i>a</i>	3
4	<i>a</i>	3
4	<i>b</i>	2
2	<i>b</i>	1
4	<i>b</i>	1

Containment relationship?

$Q :$        $ans(X, X) \leftarrow edge(X, W, Y), edge(Y, W, Z), edge(X, W, Z),$   
                   $edge(X', W', Y), edge(Y, W', Z'), edge(X', W', Z')$

$Q' :$        $ans(X, X') \leftarrow edge(X, W, Y), edge(Y, W, Z), edge(X, W, Z),$   
                   $edge(X', W', Y), edge(Y, W', Z'), edge(X', W', Z')$

### Example 3

*Sales(Part, Supplier, Customer),  
Part(PName, Type),  
Cust(CName, CAddr),  
Supp(SName, SAddr).*

Containment relationship?

$Q : \quad \text{ans}(T) \leftarrow \text{Sales}(P, S, C), \text{Part}(P, T), \text{Cust}(C, A), \text{Supp}(S, A)$

$Q' : \quad \text{ans}(T) \leftarrow \text{Sales}(P, S, C), \text{Part}(P, T), \text{Cust}(C, A), \text{Supp}(S, A),$   
 $\text{Sales}(P', S', C'), \text{Part}(P', T)$

$\Rightarrow$  Minimization!

## Lemma

Let

$$Q_1 : \quad \text{ans}(\vec{U}) \leftarrow R_1(\vec{U}_1), \dots, R_n(\vec{U}_n)$$

$$Q_2 : \quad \text{ans}(\vec{U}) \leftarrow S_1(\vec{V}_1), \dots, S_m(\vec{V}_m)$$

be conjunctive queries, where

$$\{R_1(\vec{U}_1), \dots, R_n(\vec{U}_n)\} \supseteq \{S_1(\vec{V}_1), \dots, S_m(\vec{V}_m)\}$$

Then  $Q_1 \sqsubseteq Q_2$ .

?

↑

Having more constraint, you  
will have less answers

## Substitution

- A *substitution*  $\theta$  over a set of variables  $\mathcal{V}$  is a mapping from  $\mathcal{V}$  to  $\mathcal{V} \cup \text{dom}$ , where  $\text{dom}$  a corresponding domain.
- We extend  $\theta$  to constants  $a \in \text{dom}$  and relation names  $R \in \mathcal{R}$ , where  $\theta(a) = a$ , resp.  $\theta(R) = R$ .

Note, differently to a *match*, variables may be renamed, i.e. mapped to variables.

## Example

Consider

$Q :$   $\text{ans}(T) \leftarrow \text{Sales}(P, S, C), \text{Part}(P, T), \text{Cust}(C, A), \text{Supp}(S, A)$

$Q' :$   $\text{ans}(T) \leftarrow \text{Sales}(P, S, C), \text{Part}(P, T), \text{Cust}(C, A), \text{Supp}(S, A),$   
 $\text{Sales}(P', S', C'), \text{Part}(P', T)$

and  $\theta$ :

$X$	$P$	$P'$	$S$	$S'$	$C$	$C'$	$T$	$A$
$\theta(X)$	$P$	$P$	$S$	$S$	$C$	$C$	$T$	$A$

## Containment Mapping (Homomorphism)

Given conjunctive queries

$$\begin{array}{ll} Q_1 : & \text{ans}(\vec{U}) \leftarrow R_1(\vec{U}_1), \dots, R_n(\vec{U}_n) \\ Q_2 : & \text{ans}(\vec{V}) \leftarrow S_1(\vec{V}_1), \dots, S_m(\vec{V}_m) \end{array}$$

Substitution  $\theta$  is called containment mapping from  $Q_2$  to  $Q_1$ , if  $Q_2$  can be transformed by means of  $\theta$  to become part of  $Q_1$ :

- $\theta(\text{ans}(\vec{V})) = \text{ans}(\vec{U})$ ,
- for  $i = 1, \dots, m$  there exists a  $j \in \{1, \dots, n\}$ , such that  $\theta(S_i(\vec{V}_i)) = R_j(\vec{U}_j)$ .

## Example

$Q :$        $ans(T) \leftarrow Sales(P, S, C), Part(P, T), Cust(C, A), Supp(S, A)$

$Q' :$        $ans(T) \leftarrow Sales(P, S, C), Part(P, T), Cust(C, A), Supp(S, A),$   
                   $Sales(P', S', C'), Part(P', T)$

$\theta:$

$X$	$P$	$P'$	$S$	$S'$	$C$	$C'$	$T$	$A$
$\theta(X)$	$P$	$P$	$S$	$S$	$C$	$C$	$T$	$A$

$\theta$  is a containment mapping.


## Theorem

Let

$$\begin{array}{ll} Q_1 : & \text{ans}(\vec{U}) \leftarrow R_1(\vec{U}_1), \dots, R_n(\vec{U}_n) \\ Q_2 : & \text{ans}(\vec{V}) \leftarrow S_1(\vec{V}_1), \dots, S_m(\vec{V}_m) \end{array}$$

be conjunctive queries.

$Q_1 \sqsubseteq Q_2$  iff there exists a containment mapping  $\theta$  from  $Q_2$  to  $Q_1$ .



$$\begin{array}{ll}
 Q_1 : & \text{ans}(\vec{U}) \leftarrow R_1(\vec{U}_1), \dots, R_n(\vec{U}_n) \\
 Q_2 : & \text{ans}(\vec{V}) \leftarrow S_1(\vec{V}_1), \dots, S_m(\vec{V}_m) \\
 & Q_1 \sqsubseteq Q_2?
 \end{array}$$

Proof " $\Leftarrow$ ", i.e. there exists a containment mapping  $\theta$  from  $Q_2$  to  $Q_1$ :

Let  $\mathcal{I}$  be a database instance and let  $\mu \in Q_1(\mathcal{I})$ :

There exists a substitution  $\tau$ , such that  $\tau(\vec{U}_j) \in \mathcal{I}(R_j)$ ,  $j \in \{1, \dots, n\}$  and  $\mu = \tau(\vec{U})$ .

Consider a substitution  $\tau' = \tau \circ \theta^1$  and further  $\tau'(S_i(\vec{V}_i))$ ,  $i \in \{1, \dots, m\}$ .

There holds  $\tau'(\vec{V}_i) \in \mathcal{I}(S_i)$ ,  $i \in \{1, \dots, m\}$  and therefore also  $\mu = \tau'(\vec{V})$ .

i.e.,  $\mu \in Q_2(\mathcal{I})$ .

□

---

${}^1\tau'(\cdot) = \tau(\theta(\cdot))$



Proof " $\Rightarrow$ " is based on a *canonical instance* of a query  $Q$

Let  $Q$  be a conjunctive query  $ans(\vec{U}) \leftarrow R_1(\vec{U}_1), \dots, R_n(\vec{U}_n)$  over a database schema  $\mathcal{R}$ .

The *canonical instance*  $\mathcal{I}_Q$  to  $Q$  is an instance of  $\mathcal{R} = \{R_1, \dots, R_n\}$  constructed as follows.

Let  $\tau$  be a substitution, which assigns to any  $X$  in  $Q$  an unique constant  $a_X$ .

- For any subgoal  $R(t_1, \dots, t_n)$  in the body of  $Q$  insert a tuple of the form  $(\tau(t_1), \dots, \tau(t_n))$  into  $\mathcal{I}_Q(R)$ ; thus  $\tau(R(t_1, \dots, t_n)) \in \mathcal{I}_Q(R)$ .

No other tuples are inserted into  $\mathcal{I}_Q(R)$ .

$\tau$  is called *canonical substitution*.

## Example

$Q :$   $ans(T) \leftarrow Sales(P, S, C), Part(P, T), Cust(C, A), Supp(S, A)$

$Q' :$   $ans(T) \leftarrow Sales(P, S, C), Part(P, T), Cust(C, A), Supp(S, A),$   
 $Sales(P', S', C'), Part(P', T)$

$I_Q :$

*Subst. A. →*

<i>Sales</i>	<i>Part</i>	<i>Cust</i>	<i>Supp</i>
$a_P \quad a_S \quad a_C$	$a_P \quad a_T$	$a_C \quad a_A$	$a_S \quad a_A$

$I_{Q'} :$

<i>Sales</i>	<i>Part</i>	<i>Cust</i>	<i>Supp</i>
$a_P \quad a_S \quad a_C$	$a_P \quad a_T$	$a_C \quad a_A$	$a_S \quad a_A$
$a_{P'} \quad a_{S'} \quad a_{C'}$	$a_{P'} \quad a_T$		

$$Q_1 : \quad \text{ans}(\vec{U}) \leftarrow R_1(\vec{U}_1), \dots, R_n(\vec{U}_n)$$

$$Q_2 : \quad \text{ans}(\vec{V}) \leftarrow S_1(\vec{V}_1), \dots, S_m(\vec{V}_m)$$

$$Q_1 \sqsubseteq Q_2?$$

Proof " $\Rightarrow$ ", i.e. we assume  $Q_1 \sqsubseteq Q_2$ :

Consider  $\mathcal{I}_{Q_1}$  and the corresponding canonical substitution  $\tau$ .

Then  $\tau(\text{ans}(\vec{U})) \in Q_1(\mathcal{I}_{Q_1})$ .

Because of  $Q_1 \sqsubseteq Q_2$  further  $\tau(\text{ans}(\vec{U})) \in Q_2(\mathcal{I}_{Q_1})$ .

Thus, there exists a substitution  $\rho$ , such that  $\rho(S_i(\vec{V}_i)) = \tau(R_j(\vec{U}_j))$ ,  $1 \leq i \leq m$ ,  $j \in \{1, \dots, n\}$  und  $\rho(\text{ans}(\vec{V})) = \tau(\text{ans}(\vec{U}))$ .

$\tau^{-1} \circ \rho$  is a containment mapping from  $Q_2$  to  $Q_1$

□

## Corollary

Let

$$Q_1 : \quad \text{ans}(\vec{U}) \leftarrow R_1(\vec{U}_1), \dots, R_n(\vec{U}_n)$$

$$Q_2 : \quad \text{ans}(\vec{V}) \leftarrow S_1(\vec{V}_1), \dots, S_m(\vec{V}_m)$$

be conjunctive queries,  $\mathcal{I}_{Q_1}$  the canonical instance to  $Q_1$  with canonical substitution  $\tau$ .

$$Q_1 \sqsubseteq Q_2, \text{ iff } \tau(\text{ans}(\vec{U})) \in Q_2(\mathcal{I}_{Q_1}).$$

**Proof:** It remains to show, whenever  $\tau(\text{ans}(\vec{U})) \in Q_2(\mathcal{I}_{Q_1})$ , then  $Q_1 \sqsubseteq Q_2$ .

For any  $S_j$  there exists a nonempty  $R_i$ , such that  $S_j = R_i$ .

Further, there exists a substitution  $\rho$ , such that for  $S_j(\vec{V}_j)$  we have  $\rho(\vec{V}_j) \in \mathcal{I}_{Q_1}(R_i)$ .

$\tau^{-1} \circ \rho$  is a containment mapping from  $Q_2$  to  $Q_1$ .



## Example

$$ans(a_T) \in Q(\mathcal{I}_{Q'})$$

and

$$ans(a_T) \in Q'(\mathcal{I}_Q).$$

Theorem:

Query containment for conjunctive queries is NP-complete.

Query answering? Possible in polynomial time w.r.t. size of the database (ignoring size of the query).

## 3.2 Minimization of Conjunctive Queries

### Problem

A query  $Q'$  is a subquery of a query  $Q$ , if the body of  $Q'$  is a subset of the body of  $Q$ .

Given  $Q_1$ , construct an equivalent query  $Q_2$ , which has as most as many subgoals in its body as  $Q_1$  and is minimal in the sense, that any query  $Q_3$  being equivalent to  $Q_2$  has at least as many subgoals in the body as  $Q_2$ .

Can minimization be done by deleting subgoals from  $Q_1$ , i.e. the result  $Q_2$  is a subquery of  $Q_1$ ?

### Example:

$Q :$              $ans(T) \leftarrow Sales(P, S, C), Part(P, T), Cust(C, A), Supp(S, A)$   
 $Q' :$              $ans(T) \leftarrow Sales(P, S, C), Part(P, T), Cust(C, A), Supp(S, A),$   
                        $Sales(P', S', C'), Part(P', T)$

$Q$  is minimal and equivalent to  $Q'$ .

## Theorem

Let  $Q_1 : \text{ans}(\vec{U}) \leftarrow R_1(\vec{U}_1), \dots, R_n(\vec{U}_n)$  be a conjunctive query.  
Then there exists a minimal conjunctive query  $Q_2$  equivalent to  $Q_1$ ,

$$Q_2 : \text{ans}(\vec{V}) \leftarrow S_1(\vec{V}_1), \dots, S_m(\vec{V}_m),$$

such that  $\{S_1(\vec{V}_1), \dots, S_m(\vec{V}_m)\} \subseteq \{R_1(\vec{U}_1), \dots, R_n(\vec{U}_n)\}$ .

## Proof

We can assume the existence of some conjunctive query  $Q_3$  which is minimal and equivalent to  $Q_1$ .

Because of equivalence, there exists containment mappings  $\theta$  from  $Q_1$  to  $Q_3$ , and also  $\lambda$  from  $Q_3$  to  $Q_1$ .

Let w.o.l.g.  $\{S_1(\vec{V}_1), \dots, S_m(\vec{V}_m)\}$  be that subset of subgoals from  $Q_1$ , which are images with respect to  $\lambda$  and let  $Q_2$  be a conjunctive query built out of these subgoals and no others.

- (i) We have  $Q_1 \sqsubseteq Q_2$  as  $Q_1$  may have additional subgoals to the subgoals also being subgoals of  $Q_2$ .
- (ii)  $Q_2 \sqsubseteq Q_1$  as  $\lambda \circ \theta$  is a containment mapping, i.e. each subgoal of  $Q_1$  is guaranteed to be mapped on one subgoal of  $Q_2$ .
- (iii) Minimality follows as, because of  $\lambda$ ,  $Q_2$  cannot have more subgoals than  $Q_3$ .



Query minimization is NP-hard.

We can compute all possible containment mappings over query  $Q$  and select one, whose image produces a minimal set of subgoals.

### Algorithm *Conjunctive Query Minimization*

- Chose a subgoal from  $Q$  and remove it to obtain a new query  $Q'$ . We have  $Q \sqsubseteq Q'$ .
- Check if  $Q' \sqsubseteq Q$ ; if so, then  $Q'$  is equivalent and we can continue the process of removing another subgoal.
- If not, try to remove another atom from  $Q$ .



## Example:

$$Q : \quad \text{ans}(X, Z) \leftarrow R(X, 5, Z_1), R(X_1, 5, Z_2), R(X_1, 5, Z)$$

$Q$  can be minimized to  $Q'$

$$Q' : \quad \text{ans}(X, Z) \leftarrow R(X, 5, Z_1), R(X_1, 5, Z)$$

However, not to  $Q'' : \text{ans}(X, Z) \leftarrow R(X, 5, Z)$ , as  $Q''$  and  $Q$ , respectively  $Q''$  and  $Q'$  are not equivalent.

$X \rightarrow X$   
 $Z_1 \rightarrow Z$   
 $X_1 \rightarrow X$   
 $Z_2 \rightarrow Z$   
 $Z \rightarrow Z$