

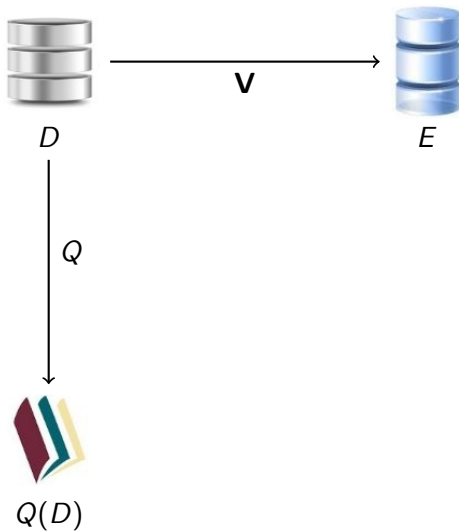
Datalog Rewritings of Regular Path Queries using Views

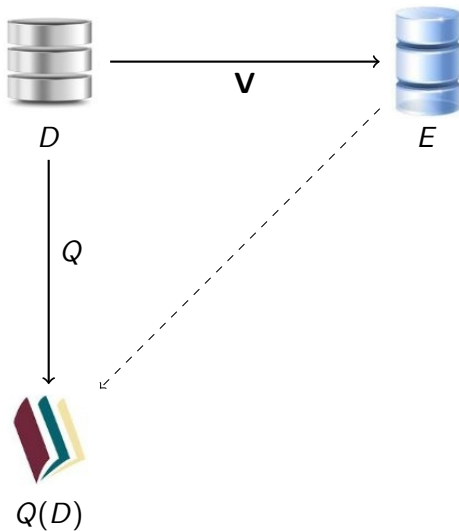
Nadime Francis
with Luc Segoufin and Cristina Sirangelo

ENS-Cachan, Inria

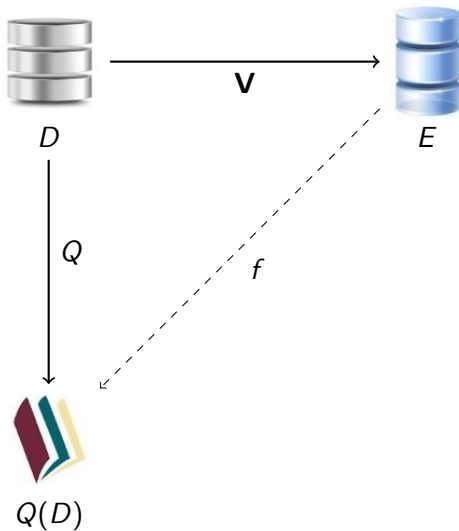
Tuesday, March, 25th
ICDT 2014 - Athens

DETERMINACY AND QUERY REWRITING

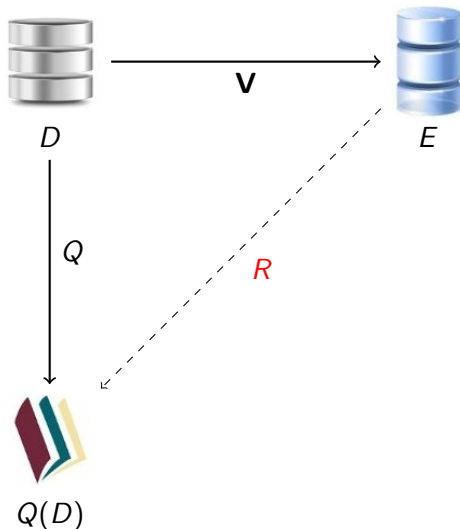




① Can I get $Q(D)$ from E ?



- 1 Can I get $Q(D)$ from E ?
- 2 Does it work for all D ?



- 1 Can I get $Q(D)$ from E ?
- 2 Does it work for all D ?
- 3 If so, can I compute R ?
(in what language ?)

Setting

- Graph databases : D
- Regular Path Queries : Q
- Regular Path Views : $\mathbf{V} = \{V_1, \dots, V_n\}$
- View extensions : E

Determinacy

- $\forall D, D' \quad \mathbf{V}(D) = \mathbf{V}(D') \Rightarrow Q(D) = Q(D')$
- $Q, \mathbf{V} \rightsquigarrow f$: function induced by Q using \mathbf{V}

Rewriting

- Query R such that
 $\forall D \quad R(\mathbf{V}(D)) = f(\mathbf{V}(D))$

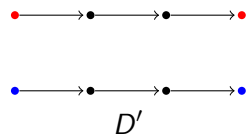
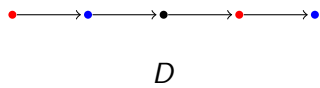
Example 1

$$V_3 = a^3 \quad Q = a^4$$

Example 1

$$V_3 = a^3 \quad Q = a^4$$

The view does not determine the query.



Example 2 (Afrati)

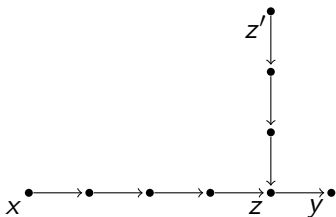
$$V_3 = a^3 \quad V_4 = a^4 \quad Q = a^5$$

Example 2 (Afrati)

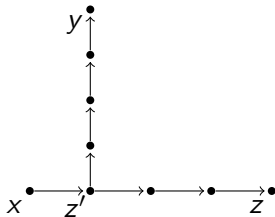
$$V_3 = a^3 \quad V_4 = a^4 \quad Q = a^5$$

$$R(x, y) = \exists z, V_4(x, z) \wedge \forall z' V_3(z', z) \Rightarrow V_4(z', y)$$

$Q \Rightarrow R$



$R \Rightarrow Q$



Example 2 (Afrati)

$$V_3 = a^3 \quad V_4 = a^4 \quad Q = a^5$$

$$R(x, y) = \exists z, V_4(x, z) \wedge \forall z' V_3(z', z) \Rightarrow V_4(z', y)$$

$$R'(x, y) = \exists z, t, V_4(x, z) \wedge V_3(x, t) \wedge \forall z' V_3(z', z) \Rightarrow V_4(z', y)$$

R' is also a rewriting :

$$\forall D, R(\mathbf{V}(D)) = R'(\mathbf{V}(D))$$

R and R' behave differently outside of view images :

$$\exists E, R(E) \neq R'(E)$$

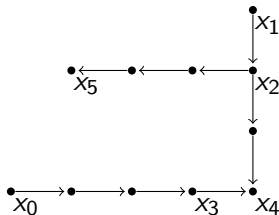
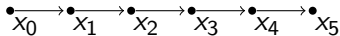
Example 2 (Afrati)

$$V_3 = a^3 \quad V_4 = a^4 \quad Q = a^5$$

$V \rightarrow Q$ but f is not monotone.

$$Q(D) = \{(x_0, x_5)\}$$

$$V(D') \subseteq V(D) \text{ but } Q(D) = \emptyset$$



Monotone determinacy

- $\mathbf{V} \rightarrow Q$ and f is monotone.
- Equivalently :

$$\forall D, D' \quad \mathbf{V}(D) \subseteq \mathbf{V}(D') \Rightarrow Q(D) \subseteq Q(D')$$

Motivations

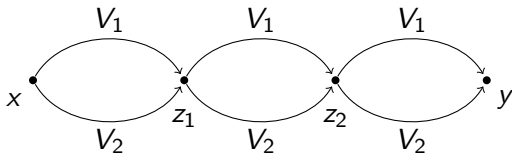
- Many interesting languages are monotone : conjunctive queries, regular path queries, Datalog, etc.
- [Nash et al.] : for CQ Views and Queries, monotone determinacy corresponds to the existence of a CQ rewriting.
- [Calvanese et al.] : decidability for RPQ Views and Queries.

Example 3

$$V_1 = a \mid aa \quad V_2 = aa \mid aaa \quad Q = a(a^6)^* \mid aa(a^6)^*$$

$$R(x, y) = \exists z V_1(x, z) \wedge T^*(z, y)$$

where $T(x, y)$ is defined as :



Our goal is to prove the following theorem :

Theorem

If \mathbf{V} and Q are RPQs and \mathbf{V} determines Q in a monotone way, then there exists a Datalog rewriting of Q using \mathbf{V} .

Corollary

Given an RPQ query Q and RPQ views \mathbf{V} , it is decidable whether there exists a Datalog rewriting of Q using \mathbf{V} .

CERTAIN ANSWERS AND CSP

Certain answers (Abiteboul, Duschka)

$$\text{cert}_{Q,\mathbf{V}}(E) = \bigcap_{D \mid E \subseteq \mathbf{V}(D)} Q(D)$$

Remarks

- $\text{cert}_{Q,\mathbf{V}}$ is a rewriting.
- **BUT** $\text{cert}_{Q,\mathbf{V}}$ is coNP-hard to evaluate.
(even if $\mathbf{V} \rightarrow Q$ in a monotone way)
- This does not prevent the existence of PTime rewritings.

Example 4

$$V_1 = \{rg, gr, rb, br, gb, bg\}$$
$$Q = \{\Sigma^* \cdot \alpha_1 \beta_1 \cdot \alpha_2 \beta_2 \cdot \Sigma^* \mid \beta_1 \neq \alpha_2\}$$

Let D be a database such that, for all $(x, y) \in D$,
 $(x, y) \in V_1(D) \Leftrightarrow (y, x) \in V_1(D)$.

Then $(x, y) \in \text{cert}_{Q, \mathbf{v}}(\mathbf{V}(D))$ iff D is not 3-colorable.

Example 4

$$V_1 = \{rg, gr, rb, br, gb, bg\} \quad V_2 = Q$$
$$Q = \{\Sigma^* \cdot \alpha_1 \beta_1 \cdot \alpha_2 \beta_2 \cdot \Sigma^* \mid \beta_1 \neq \alpha_2\}$$

Let D be a database such that, for all $(x, y) \in D$,
 $(x, y) \in V_1(D) \Leftrightarrow (y, x) \in V_1(D)$.

Then $(x, y) \in \text{cert}_{Q, \mathbf{V}}(\mathbf{V}(D))$ iff D is not 3-colorable.

A trivial rewriting is $R = V_2$.

But $\text{cert}_{Q, \mathbf{V}}$ is still as hard to evaluate.

Constraint Satisfaction Problem

- Fix a template structure T .
- $E \in \text{CSP}(T)$ iff $E \xrightarrow{h} T$.

From CSP to queries

- Extend T with source and target nodes
- $(E, u, v) \in \text{CSP}(T)$ iff there exists h such that
 - $E \xrightarrow{h} T$
 - $h(u)$ is a source node
 - $h(v)$ is a target node
- $\text{CSP}(T)$ selects all u, v such that $(E, u, v) \in \text{CSP}(T)$

Theorem (Calvanese et al.)

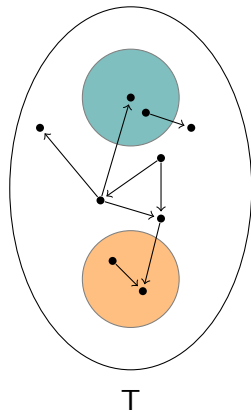
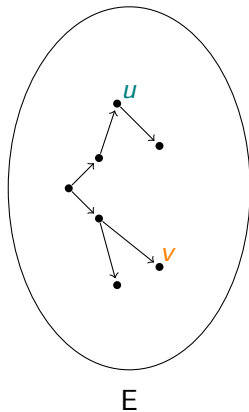
There exists a structure $T_{Q,\mathbf{v}}$ such that, for all D ,

$$(u, v) \in \text{cert}_{Q,\mathbf{v}}(\mathbf{V}(D)) \Leftrightarrow (\mathbf{V}(D), u, v) \notin \text{CSP}(T_{Q,\mathbf{v}})$$

FROM CSP TO DATALOG

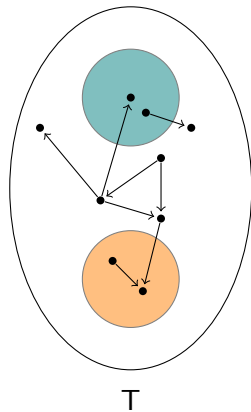
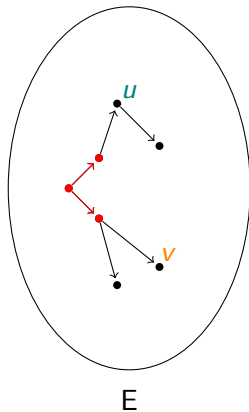
(l, k) -two-player-game

Arena : E, u, v and T with source and target nodes.



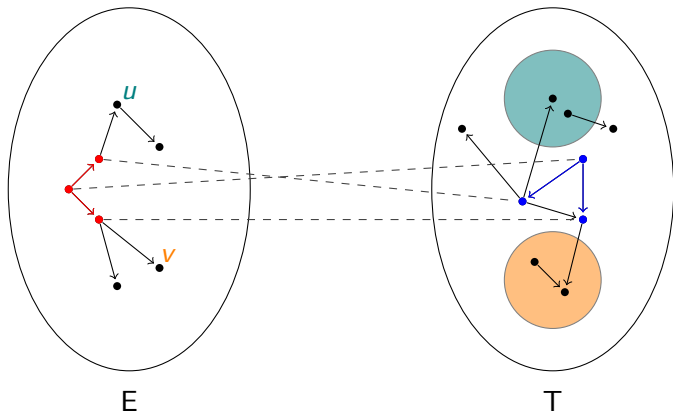
(l, k) -two-player-game

Player 1 starts by selecting k nodes on E .



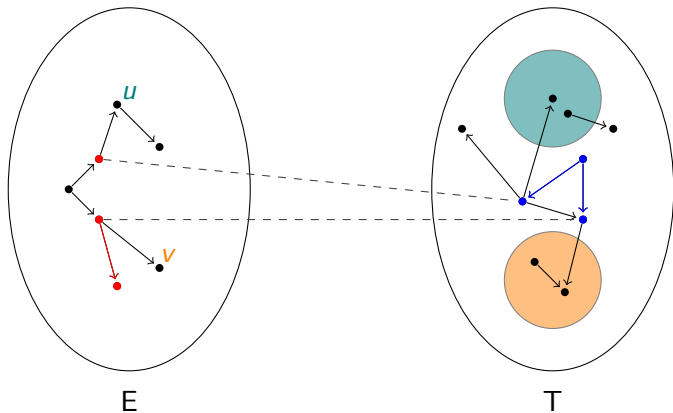
(l, k) -two-player-game

Player 2 provides a homomorphism from those nodes to T.



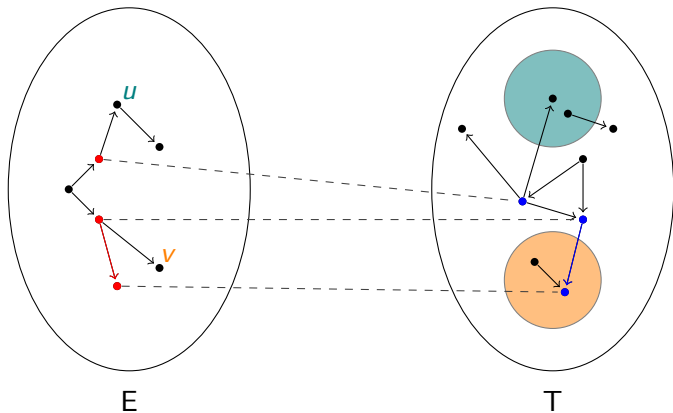
(l, k) -two-player-game

Player 1 selects k nodes, keeping at most l previous nodes.



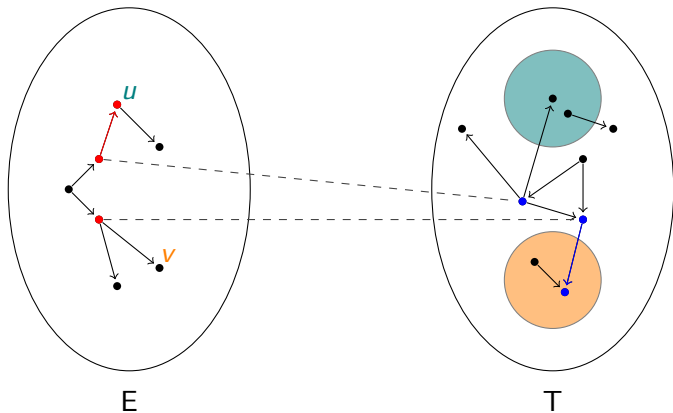
(l, k) -two-player-game

Player 2 extends previous homomorphism to new nodes.



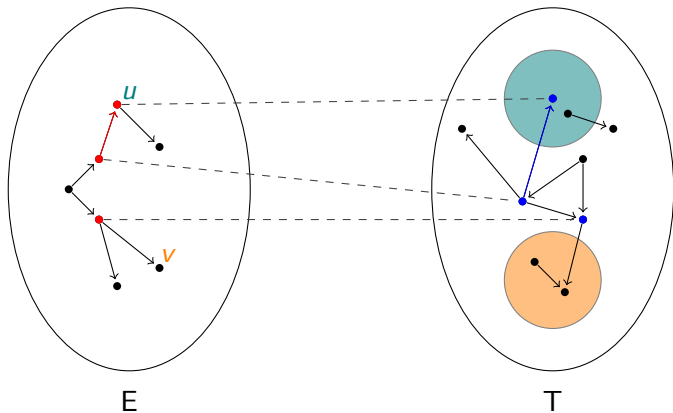
(l, k) -two-player-game

If **Player 1** selects u ...



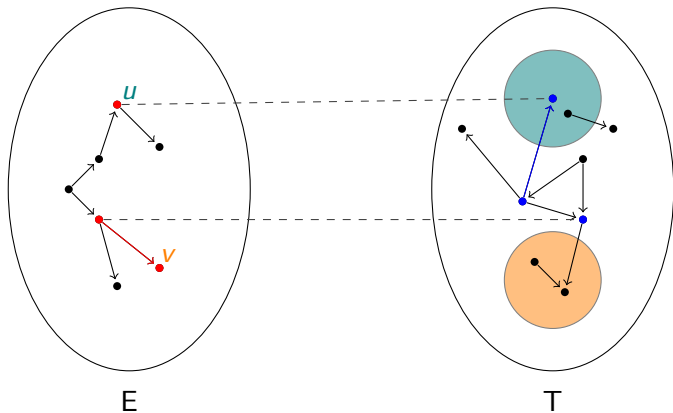
(l, k) -two-player-game

Then **Player 2** must send it to a **source** node.



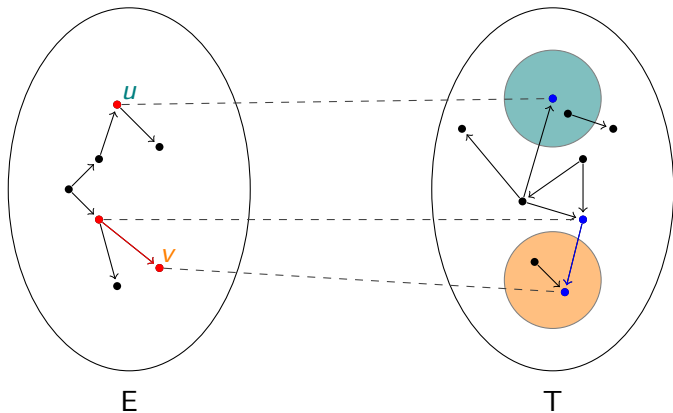
(l, k) -two-player-game

Similarly, if **Player 1** selects v ...



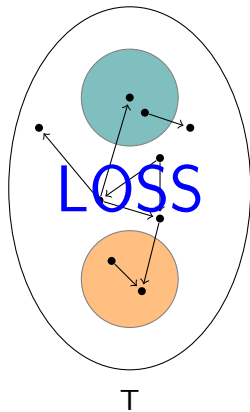
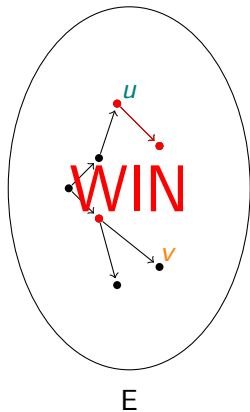
(l, k) -two-player-game

Then **Player 2** must send it to a **target** node.



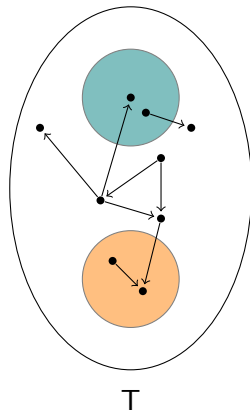
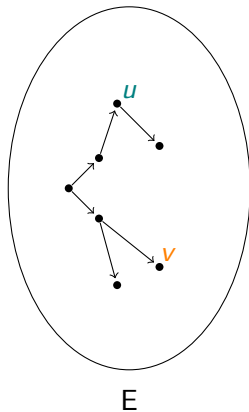
(l, k) -two-player-game

Then **Player 1** wins.



(l, k) -two-player-game

Player 2 wins if he can play forever.



Remark

- The (l, k) -two-player game on $T_{Q, \mathbf{v}}$ is an over-approximation of $\text{CSP}(T_{Q, \mathbf{v}})$.
- If $(E, u, v) \in \text{CSP}(T_{Q, \mathbf{v}})$, then Player 2 has a winning strategy.
- It is an under-approximation of $\text{cert}_{Q, \mathbf{v}}$.
- If Player 1 has a winning strategy, then $(u, v) \in \text{cert}_{Q, \mathbf{v}}(E)$.

Theorem (Feder, Vardi)

There exists a Datalog $_{l,k}$ program $Q_{l,k}$ such that $(u, v) \in Q_{l,k}(E)$ if and only if Player 1 has a winning strategy for the game played on (E, u, v) and $T_{Q, \mathbf{v}}$.

The case of simple paths

- $\exists l, k$ such that $Q_{l,k}$ is exact on views of simple paths.
- $(u, v) \in Q_{l,k}(\mathbf{V}(\pi))$ if and only if $(u, v) \in \text{cert}_{Q,\mathbf{V}}(\mathbf{V}(\pi))$.
- This is the technical part of the proof !

Lifting paths to arbitrary databases

Assume that some candidate rewriting R is :

- Exact on views of simple paths.
- Sound on views of arbitrary databases.
- Closed under homomorphism.

Then it is exact on views of arbitrary databases.

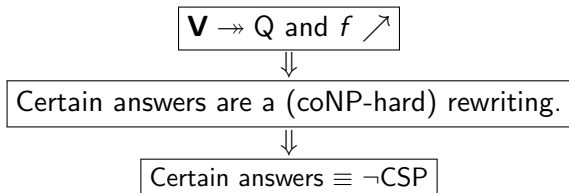
CONCLUSION AND PERSPECTIVES

$$\mathbf{V} \rightarrow \mathbf{Q} \text{ and } f \nearrow$$

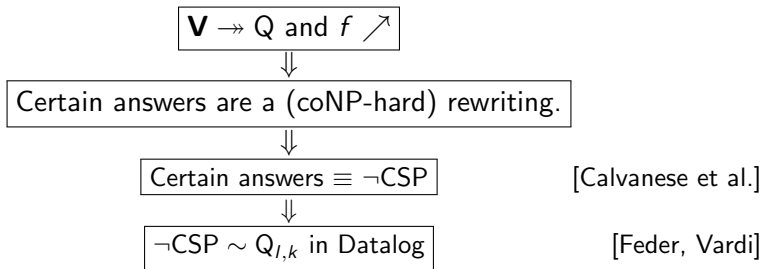
$\mathbf{V} \rightarrow \mathbf{Q}$ and $f \nearrow$

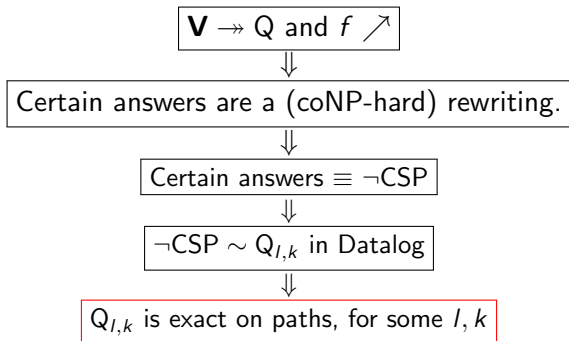


Certain answers are a (coNP-hard) rewriting.



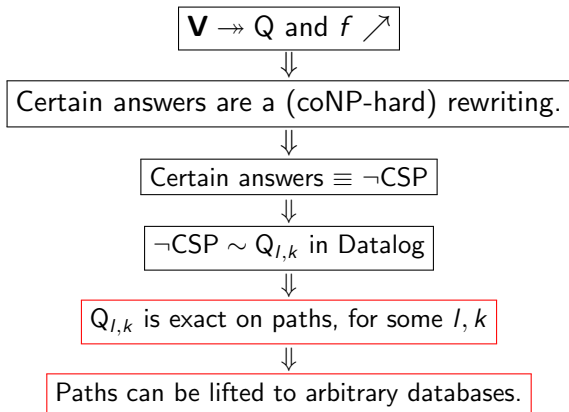
[Calvanese et al.]





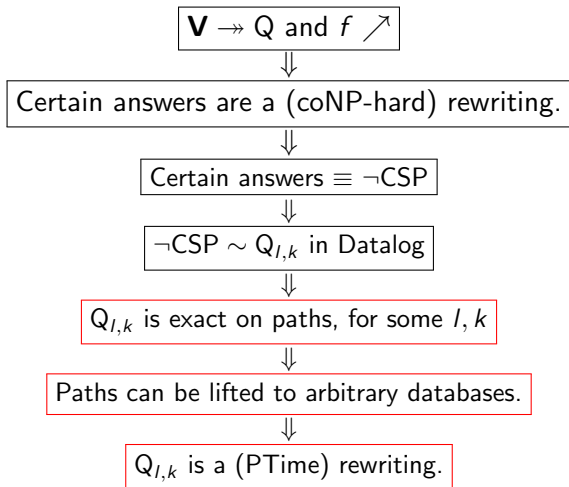
[Calvanese et al.]

[Feder, Vardi]



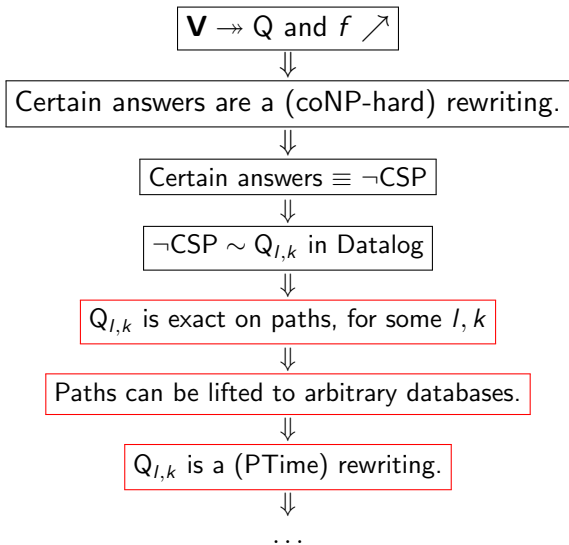
[Calvanese et al.]

[Feder, Vardi]



[Calvanese et al.]

[Feder, Vardi]



[Calvanese et al.]

[Feder, Vardi]

...



Perspectives

...



Perspectives

Rewriting language

- Datalog_{*l,k*} is not user-friendly.
- All examples in Linear Datalog.

View and query languages

- More expressive languages.
- Closed language ?

Non monotone determinacy

- [Afrati] : first-order rewriting for path views and queries.
- Work in progress : towards RPQs → add disjunction.