

Consistent Query Answering for Primary Keys on Path Queries

Paris Koutris ¹ Xiating Ouyang ¹ Jef Wijsen ²

University of Wisconsin–Madison ¹

University of Mons ²

PODS, Xi'an Shaanxi, China, June 20–25 2021

- **Primary key** constraint as the only integrity constraint
- **Inconsistent** relational databases **violating** the primary constraint
- A **repair** is an inclusion-maximal consistent subinstance

Univ	<u>Acronym</u>	City	Weather	<u>City</u>	Weather
	UW	Madison		Madison	Snow
	UW	Seattle		Seattle	Rain
	UMONS	Mons		Mons	Sunny

Q: Is there a university snowing today?

$Q() : \neg \text{Univ}(\underline{x}, y), \text{Weather}(y, \text{'Snow'})$

The problem statement

Consistent query answering – CERTAINTY(q)

INPUT: an inconsistent database db

QUESTION: Is the **fixed** Boolean query q true in all repairs of db ?



$$q_2() : \neg R(\underline{x}, z), S(\underline{y}, z)$$

$$q_1() : \neg R(\underline{x}, y), S(\underline{y}, z)$$

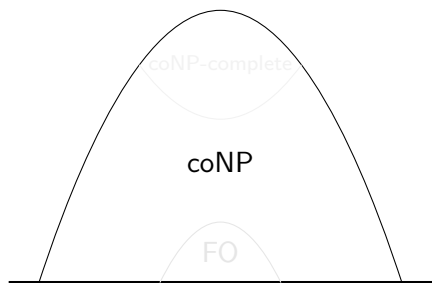
$$\text{CERTAINTY}(q_1) \iff \exists x(\exists y R(\underline{x}, y) \wedge \forall y R(\underline{x}, y) \rightarrow \exists z S(\underline{y}, z))$$

The problem statement

Consistent query answering – CERTAINTY(q)

INPUT: an inconsistent database db

QUESTION: Is the **fixed** Boolean query q true in all repairs of db?



$$q_2() : \neg R(\underline{x}, z), S(\underline{y}, z)$$

$$q_1() : \neg R(\underline{x}, y), S(\underline{y}, z)$$

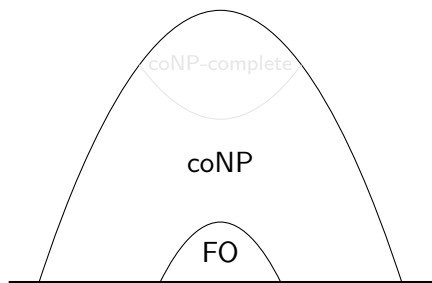
$$\text{CERTAINTY}(q_1) \iff \exists x(\exists y R(\underline{x}, y) \wedge \forall y R(\underline{x}, y) \rightarrow \exists z S(\underline{y}, z))$$

The problem statement

Consistent query answering – CERTAINTY(q)

INPUT: an inconsistent database db

QUESTION: Is the **fixed** Boolean query q true in all repairs of db ?



$$q_2() : \neg R(\underline{x}, z), S(\underline{y}, z)$$

$$q_1() : \neg R(\underline{x}, y), S(\underline{y}, z)$$

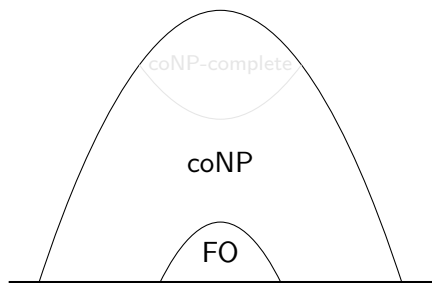
$$\text{CERTAINTY}(q_1) \iff \exists x(\exists y R(\underline{x}, y) \wedge \forall y R(\underline{x}, y) \rightarrow \exists z S(\underline{y}, z))$$

The problem statement

Consistent query answering – CERTAINTY(q)

INPUT: an inconsistent database db

QUESTION: Is the **fixed** Boolean query q true in all repairs of db ?



$$q_2() : \neg R(\underline{x}, z), S(\underline{y}, z)$$

$$q_1() : \neg R(\underline{x}, y), S(\underline{y}, z)$$

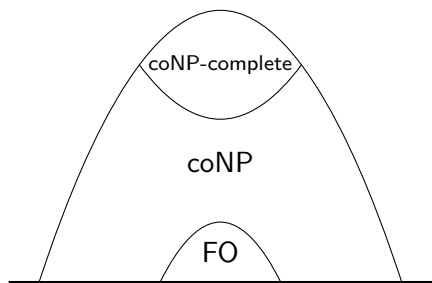
$$\text{CERTAINTY}(q_1) \iff \exists x(\exists y R(\underline{x}, y) \wedge \forall y R(\underline{x}, y) \rightarrow \exists z S(\underline{y}, z))$$

The problem statement

Consistent query answering – CERTAINTY(q)

INPUT: an inconsistent database db

QUESTION: Is the **fixed** Boolean query q true in all repairs of db ?



$$q_2() : \neg R(\underline{x}, z), S(\underline{y}, z)$$

$$q_1() : \neg R(\underline{x}, y), S(\underline{y}, z)$$

$$\text{CERTAINTY}(q_1) \iff \exists x(\exists y R(\underline{x}, y) \wedge \forall y R(\underline{x}, y) \rightarrow \exists z S(\underline{y}, z))$$

Complexity classification

Conjecture

For any BCQ q , CERTAINTY(q) is either in PTIME or coNP-complete.

Complexity classification

Conjecture

For any BCQ q , CERTAINTY(q) is either in PTIME or coNP-complete.

Classification	BCQ Class	Result
FO , non- FO	\mathcal{C}_{forest} (SJF _(self-join-free))	[Fuxman and Miller, ICDT'05]
FO , non- FO	SJF α -acyclic queries	[Wijsen, PODS'10]
P , coNP-comp.	SJF two atoms	[Kolaitis and Pema, 12]
P , coNP-comp.	SJF simple keys	[Koutris and Suciu, ICDT'14]
FO , P \ FO , coNP-comp.	SJF	[Koutris and Wijsen, PODS'15]
FO , L-comp., coNP-comp. \star	SJF	[Koutris and Wijsen, ICDT'19]
FO	SJF path	(implied by [KW, PODS'15])

Complexity classification

Conjecture

For any BCQ q , CERTAINTY(q) is either in PTIME or coNP-complete.

Classification	BCQ Class	Result
FO , non- FO	C_{forest} (SJF _(self-join-free))	[Fuxman and Miller, ICDT'05]
FO , non- FO	SJF α -acyclic queries	[Wijsen, PODS'10]
P , coNP-comp.	SJF two atoms	[Kolaitis and Pema, 12]
P , coNP-comp.	SJF simple keys	[Koutris and Suciu, ICDT'14]
FO , P \ FO , coNP-comp.	SJF	[Koutris and Wijsen, PODS'15]
FO , L-comp., coNP-comp. \star	SJF	[Koutris and Wijsen, ICDT'19]
FO \leftarrow	SJF path	(implied by [KW, PODS'15])

$$q() : -R_1(\underline{x}_1, x_2), R_2(\underline{x}_2, x_3), \dots, R_n(\underline{x}_n, x_{n+1}) \implies R_1 R_2 \dots R_n$$

distinct variables x_i , relation names R_i

Complexity classification

Conjecture

For any BCQ q , CERTAINTY(q) is either in PTIME or coNP-complete.

Classification	BCQ Class	Result
FO , non- FO	C_{forest} (SJF _(self-join-free))	[Fuxman and Miller, ICDT'05]
FO , non- FO	SJF α -acyclic queries	[Wijsen, PODS'10]
P , coNP-comp.	SJF two atoms	[Kolaitis and Pema, 12]
P , coNP-comp.	SJF simple keys	[Koutris and Suciu, ICDT'14]
FO , $\mathbf{P} \setminus \mathbf{FO}$, coNP-comp.	SJF	[Koutris and Wijsen, PODS'15]
FO , L-comp., coNP-comp. \star	SJF	[Koutris and Wijsen, ICDT'19]
FO \leftarrow	SJF path	(implied by [KW, PODS'15])
FO , NL-comp., P -comp., coNP-comp.	SJF path	(Our result)

$$q() : -R_1(\underline{x}_1, x_2), R_2(\underline{x}_2, x_3), \dots, R_n(\underline{x}_n, x_{n+1}) \implies R_1 R_2 \dots R_n$$

distinct variables x_i , relation names R_i

$$q() : -R(\underline{x}, y), R(\underline{y}, z), X(\underline{z}, w) \implies R R X$$

The curse of self-joins

Theorem (Deletion Propagation, TODS 2012)

For every CQ without self-joins, deletion propagation is either APX-hard or solvable (in polynomial time) by the unidimensional algorithm.

Theorem (Pricing, JACM 2015)

Let Q be a CQ without self-joins. The data complexity for $\text{PRICE}(Q)$ is either in PTIME or NP-complete.

Theorem (Query resilience, PODS 2020)

Let q be a single-self-join-CQ with at most two occurrences of the self-join relation. The problem $\text{RES}(q)$ is either in PTIME or NP-complete.

Theorem (Our result)

Let q be a path query. The problem $\text{CERTAINTY}(q)$ is either in FO, NL-complete, PTIME-complete or coNP-complete.

The curse of self-joins

Theorem (Deletion Propagation, TODS 2012)

For every CQ without self-joins, deletion propagation is either APX-hard or solvable (in polynomial time) by the unidimensional algorithm.

Theorem (Pricing, JACM 2015)

Let Q be a CQ without self-joins. The data complexity for $\text{PRICE}(Q)$ is either in PTIME or NP-complete.

Theorem (Query resilience, PODS 2020)

Let q be a single-self-join-CQ with at most two occurrences of the self-join relation. The problem $\text{RES}(q)$ is either in PTIME or NP-complete.

Theorem (Our result)

Let q be a path query. The problem $\text{CERTAINTY}(q)$ is either in FO, NL-complete, PTIME-complete or coNP-complete.

- 1 Handling self-joins
- 2 Classification result
- 3 Proof sketch

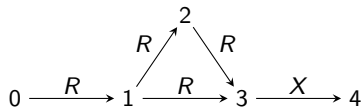
- 1 Handling self-joins
- 2 Classification result
- 3 Proof sketch

The notion of “rewinding”

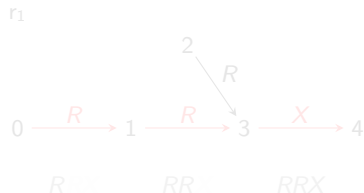
INPUT to CERTAINTY(*RRX*):

<i>R</i>	<u><i>A</i>₁</u>	<i>A</i> ₂	
	0	1	
	1	2	-
	1	3	
	2	3	-

<i>X</i>	<u><i>B</i>₁</u>	<i>B</i> ₂
	3	4



QUESTION: Do all repairs contain an *RRX* path?

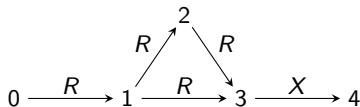


The notion of “rewinding”

INPUT to CERTAINTY(*RRX*):

<i>R</i>	<u><i>A</i>₁</u>	<i>A</i> ₂	
	0	1	
	1	2	-
	1	3	
	2	3	-

<i>X</i>	<u><i>B</i>₁</u>	<i>B</i> ₂
	3	4



QUESTION: Do all repairs contain an *RRX* path?

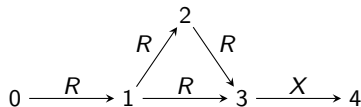


The notion of “rewinding”

INPUT to CERTAINTY(*RRX*):

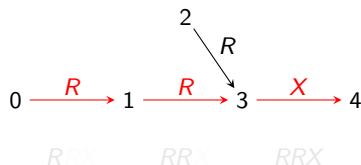
<i>R</i>	<i>A</i> ₁	<i>A</i> ₂	-
	0	1	
	1	2	
	1	3	
	2	3	

<i>X</i>	<i>B</i> ₁	<i>B</i> ₂
	3	4



QUESTION: Do all repairs contain an *RRX* path?

*r*₁



*r*₂

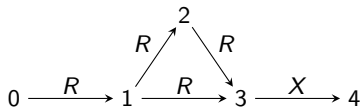


The notion of “rewinding”

INPUT to CERTAINTY(*RRX*):

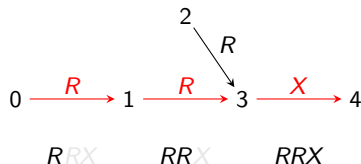
<i>R</i>	<u><i>A</i>₁</u>	<i>A</i> ₂	
	0	1	
	1	2	-
	1	3	
	2	3	-

<i>X</i>	<u><i>B</i>₁</u>	<i>B</i> ₂
	3	4



QUESTION: Do all repairs contain an *RRX* path?

*r*₁



*r*₂

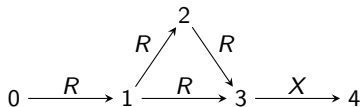


The notion of “rewinding”

INPUT to CERTAINTY(*RRX*):

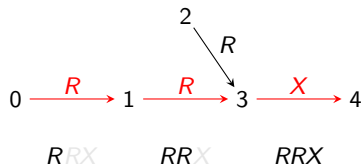
<i>R</i>	A_1	A_2	
	$\frac{0}{1}$	$\frac{1}{2}$	
	$\frac{1}{2}$	$\frac{3}{3}$	

<i>X</i>	B_1	B_2
	$\frac{3}{3}$	$\frac{4}{4}$

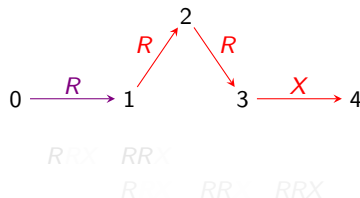


QUESTION: Do all repairs contain an *RRX* path?

r_1



r_2

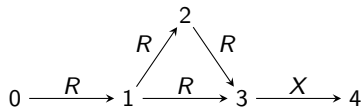


The notion of “rewinding”

INPUT to CERTAINTY(*RRX*):

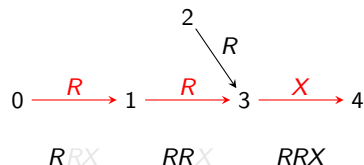
<i>R</i>	A_1	A_2	-
	$\frac{0}{1}$	$\frac{1}{2}$	-
	$\frac{1}{2}$	$\frac{3}{3}$	-

<i>X</i>	B_1	B_2
	$\frac{3}{3}$	$\frac{4}{4}$

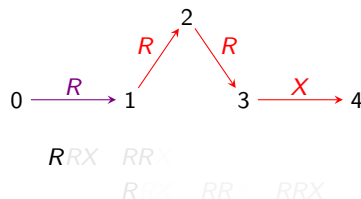


QUESTION: Do all repairs contain an *RRX* path?

r_1



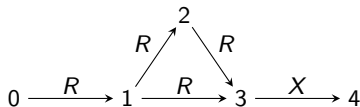
r_2



The notion of “rewinding”

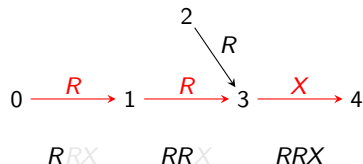
INPUT to CERTAINTY(*RRX*):

<i>R</i>	<i>A</i> ₁	<i>A</i> ₂	-	<i>X</i>	<i>B</i> ₁	<i>B</i> ₂
	0	1	-		3	4
	1	2	-			
	1	3	-			
	2	3	-			

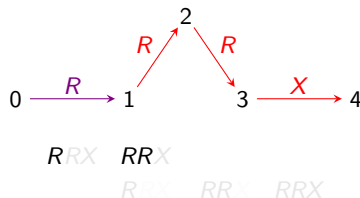


QUESTION: Do all repairs contain an *RRX* path?

*r*₁



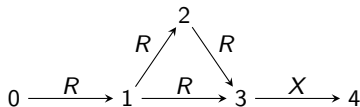
*r*₂



The notion of “rewinding”

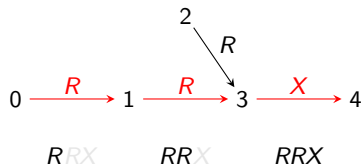
INPUT to CERTAINTY(*RRX*):

<i>R</i>	A_1	A_2	-	<i>X</i>	B_1	B_2
	$\frac{0}{1}$	$\frac{1}{2}$	-		$\frac{3}{3}$	$\frac{4}{4}$
	$\frac{1}{2}$	$\frac{3}{3}$	-			

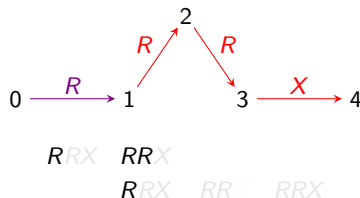


QUESTION: Do all repairs contain an *RRX* path?

r_1



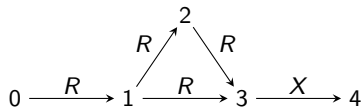
r_2



The notion of “rewinding”

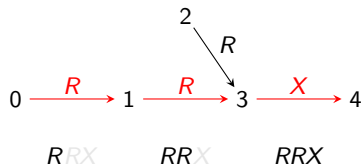
INPUT to CERTAINTY(*RRX*):

<i>R</i>	<i>A</i> ₁	<i>A</i> ₂	-	<i>X</i>	<i>B</i> ₁	<i>B</i> ₂
	0	1	-		3	4
	1	2	-			
	1	3	-			
	2	3	-			

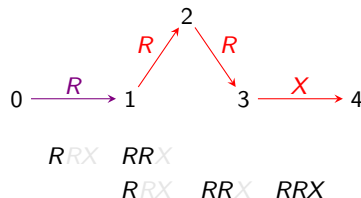


QUESTION: Do all repairs contain an *RRX* path?

*r*₁



*r*₂

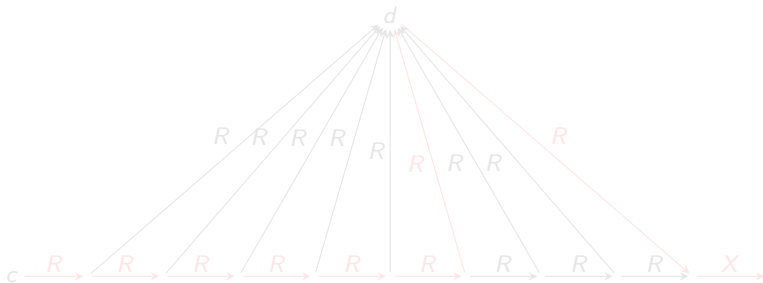


The notion of “rewinding” (cont.)

Proposition

The following statements are equivalent:

1. db is a “yes”-instance for $\text{CERTAINTY}(RRX)$; and
2. $\exists c$ such that in all repairs, there exists a path of $\underline{RR \cdot R^* \cdot X}$ starting at c .



“Reachability”, “NL-complete”

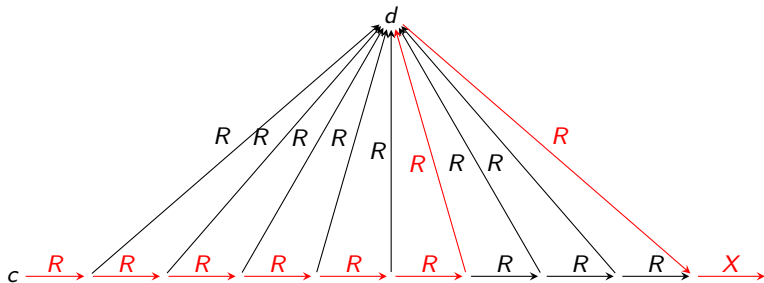
How to find the regular expression?

The notion of “rewinding” (cont.)

Proposition

The following statements are equivalent:

1. db is a “yes”-instance for $\text{CERTAINTY}(RRX)$; and
2. $\exists c$ such that in all repairs, there exists a path of $\underline{RR \cdot R^* \cdot X}$ starting at c .



“Reachability”, “NL-complete”

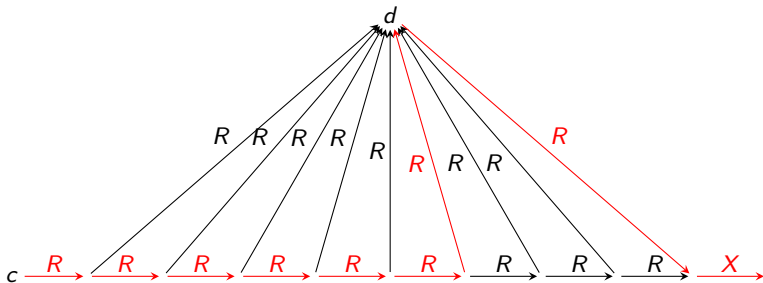
How to find the regular expression?

The notion of “rewinding” (cont.)

Proposition

The following statements are equivalent:

1. db is a “yes”-instance for $\text{CERTAINTY}(RRX)$; and
2. $\exists c$ such that in all repairs, there exists a path of $\underline{RR \cdot R^* \cdot X}$ starting at c .



“Reachability”, “NL-complete”

How to find the regular expression?

From path query to NFA



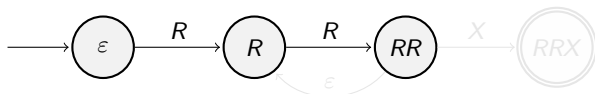
NFA(RRX)

From path query to NFA



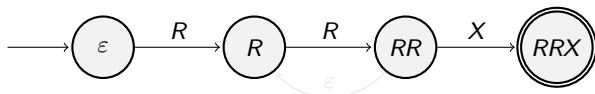
NFA(RRX)

From path query to NFA



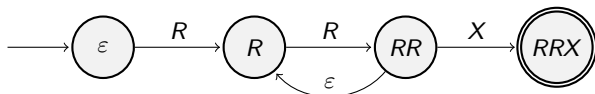
NFA(RRX)

From path query to NFA



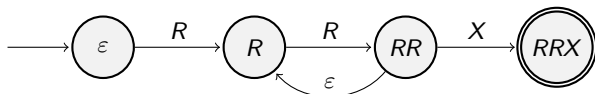
NFA(RRX)

From path query to NFA



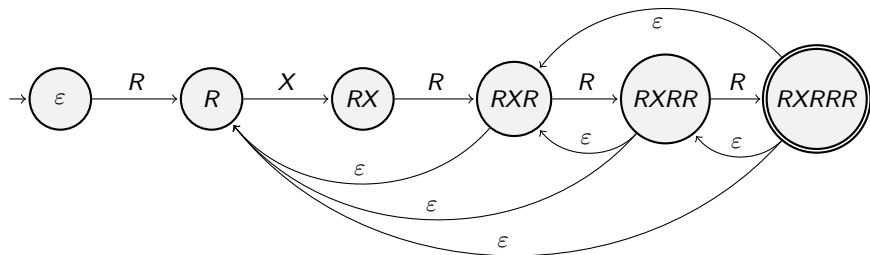
NFA(RRX)

From path query to NFA



NFA(RRX)

From path query to NFA (cont.)



NFA($RXRRR$)

- 1 Handling self-joins
- 2 Classification result**
- 3 Proof sketch

Our result

NL-hard

$$q_2 = RXRY$$

$$RXRXRY \in \text{NFA}(q_2)$$

C_1 : q is a prefix of every word in $\text{NFA}(q)$

FO-rewritable

$$q_1 = RXRX$$

$$RXRXRX \in \text{NFA}(q_1)$$

Our result

coNP-complete

$$q_4 = \text{RXRXRYRY} \quad \text{RXRXRYRXRYRY} \in \text{NFA}(q_4)$$

C_3 : q is a factor of every word in $\text{NFA}(q)$

PTIME

NL-hard

$$q_2 = \text{RXRY} \quad \text{RXRXRY} \in \text{NFA}(q_2)$$

C_1 : q is a prefix of every word in $\text{NFA}(q)$

FO-rewritable

$$q_1 = \text{RXRX} \quad \text{RXRXRX} \in \text{NFA}(q_1)$$

Our result

coNP-complete

$$q_4 = RXXRYRY \quad RXXRYRXRYRY \in \text{NFA}(q_4)$$

C_3 : q is a factor of every word in $\text{NFA}(q)$

PTIME

$$q_3 = RXYRY$$

C_2 : Whenever $q = uRvRw$, q is a factor of $uRvRvRw$; and whenever $q = uRv_1Rv_2Rw$ for consecutive occurrences of R , $v_1 = v_2$ or Rw is a prefix of Rv_1 .

NL-hard

$$q_2 = RXY \quad RXXRY \in \text{NFA}(q_2)$$

C_1 : q is a prefix of every word in $\text{NFA}(q)$

FO-rewritable

$$q_1 = RXX \quad RXXRX \in \text{NFA}(q_1)$$

Our result

coNP-complete

$$q_4 = RXXRYRY \quad RXXRYRXRYRY \in \text{NFA}(q_4)$$

C_3 : q is a factor of every word in $\text{NFA}(q)$

PTIME-complete

$$q_3 = RXYRY$$

C_2 : Whenever $q = uRvRw$, q is a factor of $uRvRvRw$; and whenever $q = uRv_1Rv_2Rw$ for consecutive occurrences of R , $v_1 = v_2$ or Rw is a prefix of Rv_1 .

NL-complete

$$q_2 = RXY \quad RXXRY \in \text{NFA}(q_2)$$

C_1 : q is a prefix of every word in $\text{NFA}(q)$

FO-rewritable

$$q_1 = RXX \quad RXXRX \in \text{NFA}(q_1)$$

C_1 , C_2 and C_3 are decidable

C_1 : q is a prefix of every word in $\text{NFA}(q)$

\iff Whenever $q = u \cdot Rv \cdot Rw$, q is a prefix of $u \cdot Rv \cdot Rv \cdot Rw$.

C_3 : q is a factor of every word in $\text{NFA}(q)$

\iff Whenever $q = u \cdot Rv \cdot Rw$, q is a factor of $u \cdot Rv \cdot Rv \cdot Rw$.

1 Handling self-joins

2 Classification result

3 Proof sketch

Lemma (PTIME)

Let q be a path query satisfying C_3 . The following statements are equivalent:

1. db is a “yes”-instance for $\text{CERTAINTY}(q)$; and
2. $\exists c$ such that in all repairs, there exists a path accepted by $\text{NFA}(q)$ starting at c .

Moreover, item 2 can be decided in **PTIME** using dynamic programming/least fixedpoint logic.

C_3 : q is a factor of every word in $\text{NFA}(q)$



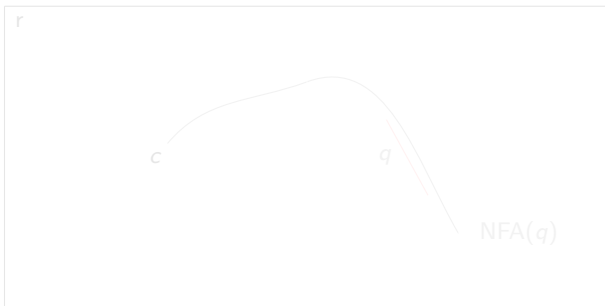
Lemma (PTIME)

Let q be a path query satisfying C_3 . The following statements are equivalent:

1. db is a “yes”-instance for $CERTAINTY(q)$; and
2. $\exists c$ such that in all repairs, there exists a path accepted by $NFA(q)$ starting at c .

Moreover, item 2 can be decided in **PTIME** using dynamic programming/least fixedpoint logic.

C_3 : q is a factor of every word in $NFA(q)$



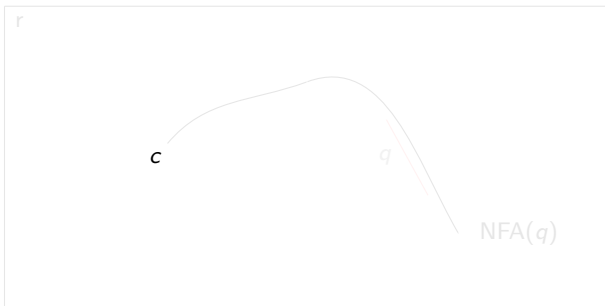
Lemma (PTIME)

Let q be a path query satisfying C_3 . The following statements are equivalent:

1. db is a “yes”-instance for $CERTAINTY(q)$; and
2. $\exists c$ such that in all repairs, there exists a path accepted by $NFA(q)$ starting at c .

Moreover, item 2 can be decided in **PTIME** using dynamic programming/least fixedpoint logic.

C_3 : q is a factor of every word in $NFA(q)$



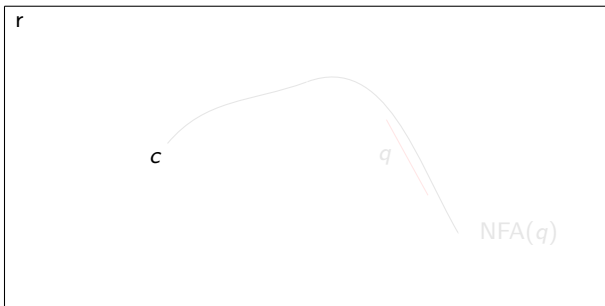
Lemma (PTIME)

Let q be a path query satisfying C_3 . The following statements are equivalent:

1. db is a “yes”-instance for $CERTAINTY(q)$; and
2. $\exists c$ such that in all repairs, there exists a path accepted by $NFA(q)$ starting at c .

Moreover, item 2 can be decided in **PTIME** using dynamic programming/least fixedpoint logic.

C_3 : q is a factor of every word in $NFA(q)$



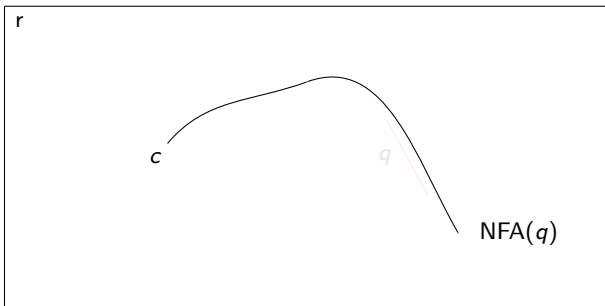
Lemma (PTIME)

Let q be a path query satisfying C_3 . The following statements are equivalent:

1. db is a “yes”-instance for $\text{CERTAINTY}(q)$; and
2. $\exists c$ such that in all repairs, there exists a path accepted by $\text{NFA}(q)$ starting at c .

Moreover, item 2 can be decided in **PTIME** using dynamic programming/least fixedpoint logic.

C_3 : q is a factor of every word in $\text{NFA}(q)$



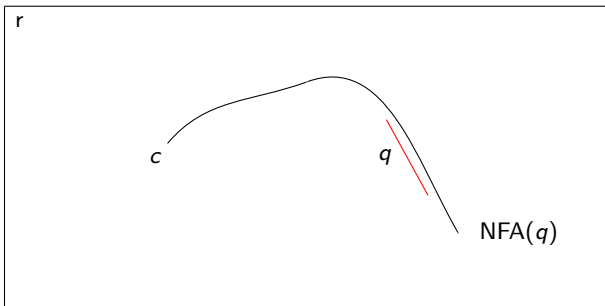
Lemma (PTIME)

Let q be a path query satisfying C_3 . The following statements are equivalent:

1. db is a “yes”-instance for $CERTAINTY(q)$; and
2. $\exists c$ such that in all repairs, there exists a path accepted by $NFA(q)$ starting at c .

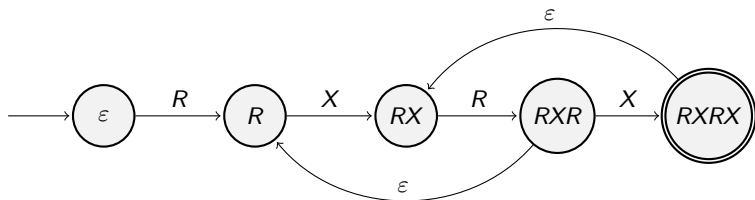
Moreover, item 2 can be decided in **PTIME** using dynamic programming/least fixedpoint logic.

C_3 : q is a factor of every word in $NFA(q)$

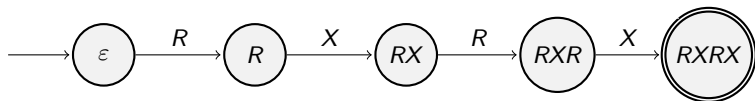


FO-rewritability

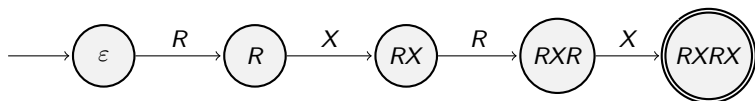
C_1 : q is a prefix of every word in $NFA(q)$



C_1 : q is a prefix of every word in $NFA(q)$



C_1 : q is a prefix of every word in $NFA(q)$



Lemma (FO)

Let q be a path query satisfying C_1 . The following statements are equivalent:

1. db is a “yes”-instance for $CERTAINTY(q)$; and
2. $\exists c$ such that in all repairs, there exists a path of q starting at c .

Moreover, item 2 can be decided in **FO**.

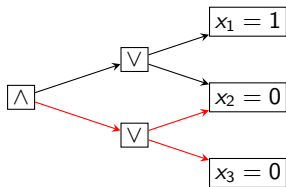
Hardness results

Lemma

Let q be a path query. Then we have

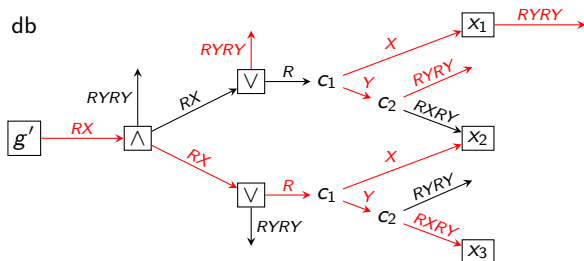
- if q does not satisfy C_1 , then $\text{CERTAINTY}(q)$ is **NL**-hard (via Reachability);
- if q does not satisfy C_2 , then $\text{CERTAINTY}(q)$ is **PTIME**-hard (via Monotone Circuit Value); and
- if q does not satisfy C_3 , then $\text{CERTAINTY}(q)$ is **coNP**-hard (via Satisfiability).

c



$$q = \underbrace{\quad}_u \underbrace{RX}_{Rv_1} \underbrace{RY}_{Rv_2} \underbrace{RY}_{Rw}$$

db



Conjecture

Let q be a CQ. Then we have

- *CERTAINTY(q) is either in PTIME or coNP-complete, and it is decidable which of the two cases applies; and*
 - *it is decidable whether or not CERTAINTY(q) is in FO.*
-
- Acyclic queries with self-joins
 - Multiple key constraints, negated atom, aggregation etc.

Conclusion

coNP-complete

$$q_4 = RXXRYRY \quad RXXRYRXXRYRY \in \text{NFA}(q_4)$$

C_3 : q is a factor of every word in $\text{NFA}(q)$

PTIME-complete

$$q_3 = RXYRY$$

C_2 : Whenever $q = uRvRw$, q is a factor of $uRvRvRw$; and whenever $q = uRv_1Rv_2Rw$ for consecutive occurrences of R , $v_1 = v_2$ or Rw is a prefix of Rv_1 .

NL-complete

$$q_2 = RXY \quad RXXRY \in \text{NFA}(q_2)$$

C_1 : q is a prefix of every word in $\text{NFA}(q)$

FO-rewritable

$$q_1 = RXX \quad RXXRX \in \text{NFA}(q_1)$$