

A Dichotomy Structure Theorem for the Resilience Problem

Cibele Freire

College of Information and Computer Sciences
University of Massachusetts Amherst

joint with work Wolfgang Gatterbauer
& Neil Immerman & Alexandra Meliou

Imagine a world where..

Imagine a world where..

We could selectively delete memories from the brain and by doing that we could change one's opinions or beliefs.

Imagine a world where..

We could selectively delete memories from the brain and by doing that we could change one's opinions or beliefs.

Would you like to know how difficult it would be to choose what memories to delete?

Imagine a world where..

We could selectively delete memories from the brain and by doing that we could change one's opinions or beliefs.

Would you like to know how difficult it would be to choose what memories to delete?

This is what we investigate!

Imagine a world where..

We could selectively delete memories from the brain and by doing that we could change one's opinions or beliefs.

Would you like to know how difficult it would be to choose what memories to delete?

This is what we investigate! Although, in a simpler, safer and more ethical scenario...

Imagine a world where..

We could selectively delete memories from the brain and by doing that we could change one's opinions or beliefs.

Would you like to know how difficult it would be to choose what memories to delete?

This is what we investigate! Although, in a simpler, safer and more ethical scenario...

- Brains = Databases
- Opinions, beliefs = Query answer
- Delete memories = Delete tuples from the database

Resilience problem

Definition (Resilience)

Given a query q and database D , we say that $(D, k) \in \text{RES}(q)$ if and only if $D \models q$ and there exists some $\Gamma \subseteq D$ such that $D - \Gamma \not\models q$ and $|\Gamma| \leq k$.

Resilience problem

Definition (Resilience)

Given a query q and database D , we say that $(D, k) \in \text{RES}(q)$ if and only if $D \models q$ and there exists some $\Gamma \subseteq D$ such that $D - \Gamma \not\models q$ and $|\Gamma| \leq k$.

$$q_{vc} :- V(x), E(x, y), V(y)$$

Resilience problem

Definition (Resilience)

Given a query q and database D , we say that $(D, k) \in \text{RES}(q)$ if and only if $D \models q$ and there exists some $\Gamma \subseteq D$ such that $D - \Gamma \not\models q$ and $|\Gamma| \leq k$.

$$q_{vc} :- V(x), E(x, y), V(y)$$

What is the complexity of $\text{RES}(q_{vc})$?

Resilience problem

Definition (Resilience)

Given a query q and database D , we say that $(D, k) \in \text{RES}(q)$ if and only if $D \models q$ and there exists some $\Gamma \subseteq D$ such that $D - \Gamma \not\models q$ and $|\Gamma| \leq k$.

$$q_{vc} :- V(x), E(x, y), V(y)$$

What is the complexity of $\text{RES}(q_{vc})$? This is exactly vertex cover.

Resilience problem

Definition (Resilience)

Given a query q and database D , we say that $(D, k) \in \text{RES}(q)$ if and only if $D \models q$ and there exists some $\Gamma \subseteq D$ such that $D - \Gamma \not\models q$ and $|\Gamma| \leq k$.

$$q_{vc} := V(x), E(x, y), V(y)$$

What is the complexity of $\text{RES}(q_{vc})$? This is exactly vertex cover.

Lemma

$\text{RES}(q_{vc})$ is NP-complete.

Resilience problem

Definition (Resilience)

Given a query q and database D , we say that $(D, k) \in \text{RES}(q)$ if and only if $D \models q$ and there exists some $\Gamma \subseteq D$ such that $D - \Gamma \not\models q$ and $|\Gamma| \leq k$.

$$q_{vc} := V(x), E(x, y), V(y)$$

What is the complexity of $\text{RES}(q_{vc})$? This is exactly vertex cover.

Lemma

$\text{RES}(q_{vc})$ is NP-complete.

- Conjunctive queries **without self-joins**
- Deleted tuples \rightarrow Contingency set Γ
- Endogenous (changeable) vs exogenous (non changeable) tuples

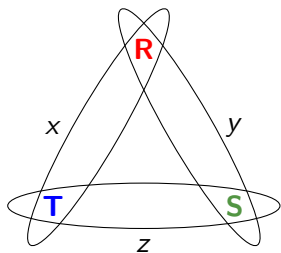
Dual hypergraph

$$q_{\Delta} := R(x, y), S(y, z), T(z, x)$$

$$q_{\Gamma} := A(x), B(y), C(z), W(x, y, z)$$

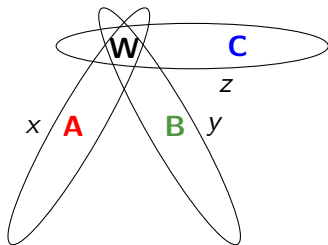
Dual hypergraph

$$q_{\Delta} := R(x, y), S(y, z), T(z, x)$$



Triangle query

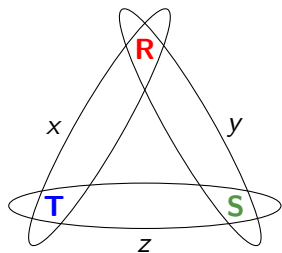
$$q_T := A(x), B(y), C(z), W(x, y, z)$$



Tripod query

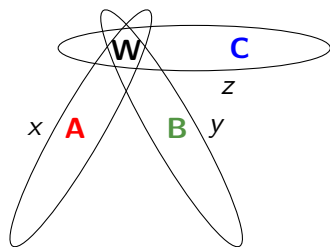
Dual hypergraph

$$q_{\Delta} := R(x, y), S(y, z), T(z, x)$$



Triangle query

$$q_T := A(x), B(y), C(z), W(x, y, z)$$

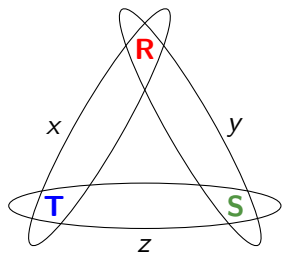


Tripod query

What is the complexity of resilience for those queries?

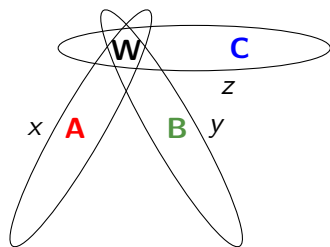
Dual hypergraph

$$q_{\Delta} := R(x, y), S(y, z), T(z, x)$$



Triangle query

$$q_T := A(x), B(y), C(z), W(x, y, z)$$



Tripod query

Lemma

$\text{RES}(q_{\Delta})$ and $\text{RES}(q_T)$ are NP-complete.

RES(q_Δ) is NP-complete.

3SAT \leq RES(q_Δ). Let $\psi = C_1 \wedge \dots \wedge C_m$ be a 3-CNF formula,
 $\text{var}(\psi) = \{v_1, \dots, v_n\}$

Map $\psi \mapsto (D_\psi, k_\psi)$ s.t. $\psi \in \text{3SAT} \Leftrightarrow (D_\psi, k_\psi) \in \text{RES}(q_\Delta)$

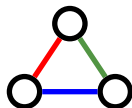
RES(q_Δ) is NP-complete.

3SAT \leq RES(q_Δ). Let $\psi = C_1 \wedge \dots \wedge C_m$ be a 3-CNF formula,
 $\text{var}(\psi) = \{v_1, \dots, v_n\}$

Map $\psi \mapsto (D_\psi, k_\psi)$ s.t. $\psi \in \text{3SAT} \Leftrightarrow (D_\psi, k_\psi) \in \text{RES}(q_\Delta)$

$$q_\Delta :- R(x, y), S(y, z), T(z, x)$$

$(D_\psi, k_\psi) \in \text{RES}(q_\Delta) \Leftrightarrow \exists \Gamma (|\Gamma| = k_\psi) \wedge (D_\psi - \Gamma)$ has no



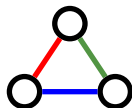
RES(q_Δ) is NP-complete.

$3SAT \leq RES(q_\Delta)$. Let $\psi = C_1 \wedge \dots \wedge C_m$ be a 3-CNF formula,
 $\text{var}(\psi) = \{v_1, \dots, v_n\}$

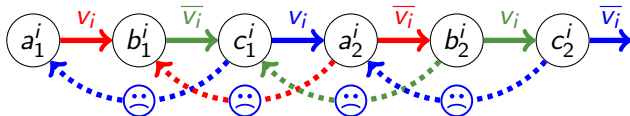
Map $\psi \mapsto (D_\psi, k_\psi)$ s.t. $\psi \in 3SAT \Leftrightarrow (D_\psi, k_\psi) \in RES(q_\Delta)$

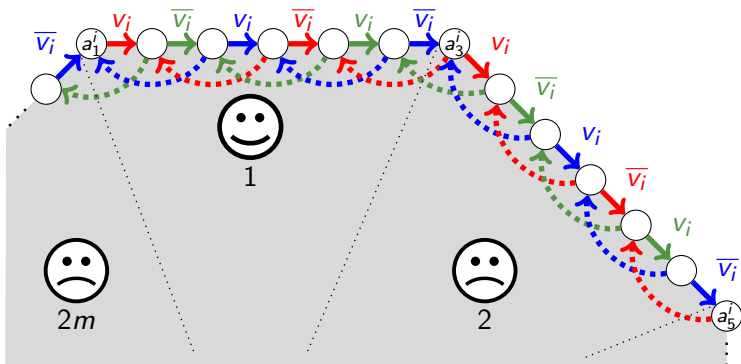
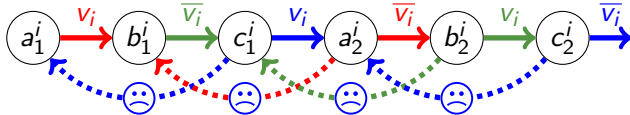
$$q_\Delta := R(x, y), S(y, z), T(z, x)$$

$(D_\psi, k_\psi) \in RES(q_\Delta) \Leftrightarrow \exists \Gamma (|\Gamma| = k_\psi) \wedge (D_\psi - \Gamma)$ has no

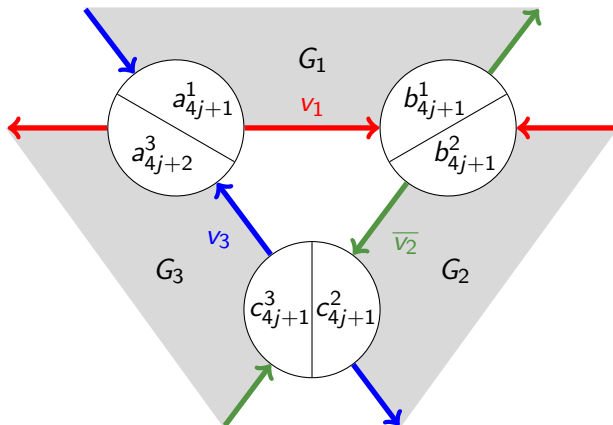


D_ψ has one circular gadget G_i for each variable v_i .

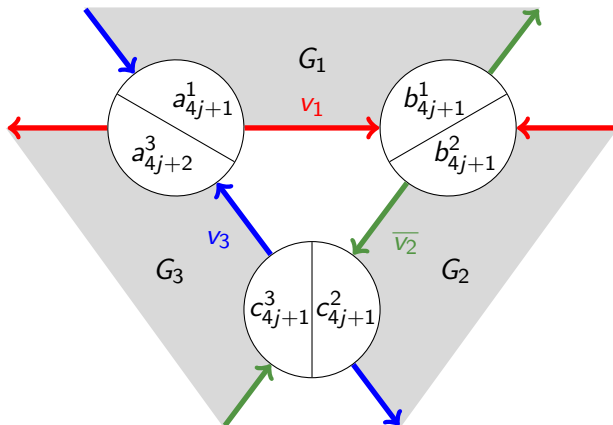




For each clause, e.g., $C_j = (v_1 \vee \overline{v_2} \vee v_3)$, pick the j th occurrences of $v_1 \in G_1$, $\overline{v_2} \in G_2$ and $v_3 \in G_3$. Identify head of v_1 with tail of $\overline{v_2}$, head of $\overline{v_2}$ with tail of v_3 , head of v_3 with tail of v_1



For each clause, e.g., $C_j = (v_1 \vee \overline{v_2} \vee v_3)$, pick the j th occurrences of $v_1 \in G_1$, $\overline{v_2} \in G_2$ and $v_3 \in G_3$. Identify head of v_1 with tail of $\overline{v_2}$, head of $\overline{v_2}$ with tail of v_3 , head of v_3 with tail of v_1



This new RGB triangle is automatically removed iff one of the literals in C_j is chosen true. \square

$\text{RES}(q_T)$ is NP-complete.

$$q_T := A(x), B(y), C(z), W(x, y, z)$$

$\text{RES}(q_T)$ is NP-complete.

$$q_T := A(x), B(y), C(z), W(x, y, z)$$

$$\text{var}(A) \subseteq \text{var}(W).$$

$\text{RES}(q_T)$ is NP-complete.

$$q_T := A(x), B(y), C(z), W(x, y, z)$$

$\text{var}(A) \subseteq \text{var}(W)$.

A dominates W.

$\text{RES}(q_T)$ is NP-complete.

$$q_T := A(x), B(y), C(z), W(x, y, z)$$

$\text{var}(A) \subseteq \text{var}(W)$.

A **dominates** W .

Proposition

If A dominates W , then we can assume that W is **exogenous**, i.e., rewrite as W^x , tuples from W^x are **never chosen**.

RES(q_T) is NP-complete.

$$q_T := A(x), B(y), C(z), W(x, y, z)$$

$$\text{var}(A) \subseteq \text{var}(W).$$

A dominates W.

Proposition

If A dominates W, then we can assume that W is **exogenous**, i.e., rewrite as W^x , tuples from W^x are **never chosen**.

$$q_T := A(x), B(y), C(z), W^x(x, y, z)$$

$\text{RES}(q_T)$ is NP-complete.

$$q_{\Delta} \quad :- \quad R(x, y), S(y, z), T(z, x)$$

$$q_T \quad :- \quad A(x), B(y), C(z), W^x(x, y, z)$$

Show $\text{RES}(q_{\Delta}) \leq \text{RES}(q_T)$

$\text{RES}(q_T)$ is NP-complete.

$$q_{\Delta} \quad :- \quad R(x, y), S(y, z), T(z, x)$$

$$q_T \quad :- \quad A(x), B(y), C(z), W^x(x, y, z)$$

Show $\text{RES}(q_{\Delta}) \leq \text{RES}(q_T)$

Let (D, k) be an instance of $\text{RES}(q_{\Delta})$.

$$(D, k) \mapsto (D', k) \quad D' \stackrel{\text{def}}{=} (A, B, C, W^x)$$

RES(q_T) is NP-complete.

$$q_{\Delta} \quad :- \quad R(x, y), S(y, z), T(z, x)$$

$$q_T \quad :- \quad A(x), B(y), C(z), W^x(x, y, z)$$

Show $\text{RES}(q_{\Delta}) \leq \text{RES}(q_T)$

Let (D, k) be an instance of $\text{RES}(q_{\Delta})$.

$$(D, k) \mapsto (D', k) \quad D' \stackrel{\text{def}}{=} (A, B, C, W^x)$$

$$A = \{ \langle ab \rangle \mid R(a, b) \in D \}$$

$$B = \{ \langle bc \rangle \mid S(b, c) \in D \}$$

$$C = \{ \langle ca \rangle \mid T(c, a) \in D \}$$

$$W^x = \{ (\langle ab \rangle, \langle bc \rangle, \langle ca \rangle) \mid a, b, c \in \text{dom}(D) \}$$

RES(q_T) is NP-complete.

$$q_\Delta \quad :- \quad R(x, y), S(y, z), T(z, x)$$

$$q_T \quad :- \quad A(x), B(y), C(z), W^x(x, y, z)$$

Show $\text{RES}(q_\Delta) \leq \text{RES}(q_T)$

Let (D, k) be an instance of $\text{RES}(q_\Delta)$.

$$(D, k) \mapsto (D', k) \quad D' \stackrel{\text{def}}{=} (A, B, C, W^x)$$

$$A = \{ \langle ab \rangle \mid R(a, b) \in D \}$$

$$B = \{ \langle bc \rangle \mid S(b, c) \in D \}$$

$$C = \{ \langle ca \rangle \mid T(c, a) \in D \}$$

$$W^x = \{ (\langle ab \rangle, \langle bc \rangle, \langle ca \rangle) \mid a, b, c \in \text{dom}(D) \}$$

Claim

$$(D, k) \in \text{RES}(q_\Delta) \quad \Leftrightarrow \quad (D', k) \in \text{RES}(q_T).$$



Triads

What do the triangle query and the tripod query have in common?

Triads

What do the triangle query and the tripod query have in common?

Definition (triad)

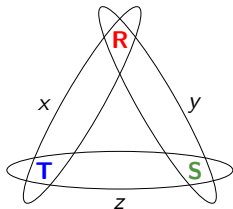
A **triad** is a set of three *endogenous atoms*, $\mathcal{T} = \{S_0, S_1, S_2\}$ such that for every pair i, j , there is a path from S_i to S_j that uses no variable occurring in the other atom of \mathcal{T} .

Triads

What do the triangle query and the tripod query have in common?

Definition (triad)

A **triad** is a set of three *endogenous atoms*, $\mathcal{T} = \{S_0, S_1, S_2\}$ such that for every pair i, j , there is a path from S_i to S_j that uses no variable occurring in the other atom of \mathcal{T} .

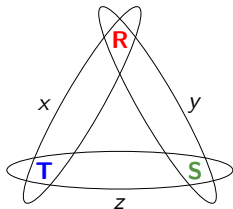


Triads

What do the triangle query and the tripod query have in common?

Definition (triad)

A **triad** is a set of three *endogenous atoms*, $\mathcal{T} = \{S_0, S_1, S_2\}$ such that for every pair i, j , there is a path from S_i to S_j that uses no variable occurring in the other atom of \mathcal{T} .



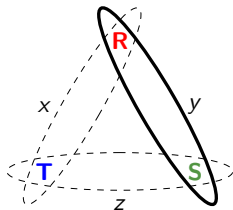
Where is the triad?

Triads

What do the triangle query and the tripod query have in common?

Definition (triad)

A **triad** is a set of three *endogenous atoms*, $\mathcal{T} = \{S_0, S_1, S_2\}$ such that for every pair i, j , there is a path from S_i to S_j that uses no variable occurring in the other atom of \mathcal{T} .



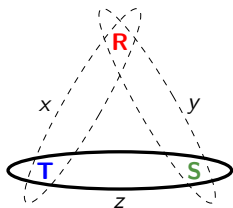
Where is the triad?
Path from R to S .

Triads

What do the triangle query and the tripod query have in common?

Definition (triad)

A **triad** is a set of three *endogenous atoms*, $\mathcal{T} = \{S_0, S_1, S_2\}$ such that for every pair i, j , there is a path from S_i to S_j that uses no variable occurring in the other atom of \mathcal{T} .



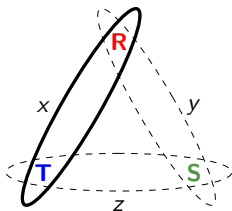
Where is the triad?
Path from S to T .

Triads

What do the triangle query and the tripod query have in common?

Definition (triad)

A **triad** is a set of three *endogenous atoms*, $\mathcal{T} = \{S_0, S_1, S_2\}$ such that for every pair i, j , there is a path from S_i to S_j that uses no variable occurring in the other atom of \mathcal{T} .



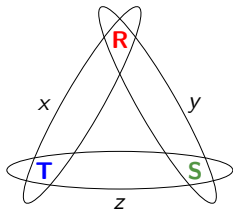
Where is the triad?
Path from T to R .

Triads

What do the triangle query and the tripod query have in common?

Definition (triad)

A **triad** is a set of three *endogenous atoms*, $\mathcal{T} = \{S_0, S_1, S_2\}$ such that for every pair i, j , there is a path from S_i to S_j that uses no variable occurring in the other atom of \mathcal{T} .



Where is the triad?

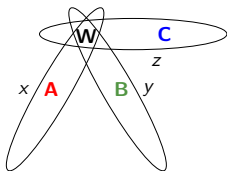
$$\mathcal{T} = \{R, S, T\}$$

Triads

What do the triangle query and the tripod query have in common?

Definition (triad)

A **triad** is a set of three *endogenous atoms*, $\mathcal{T} = \{S_0, S_1, S_2\}$ such that for every pair i, j , there is a path from S_i to S_j that uses no variable occurring in the other atom of \mathcal{T} .



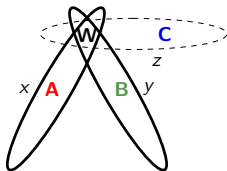
Where is the triad?

Triads

What do the triangle query and the tripod query have in common?

Definition (triad)

A **triad** is a set of three *endogenous atoms*, $\mathcal{T} = \{S_0, S_1, S_2\}$ such that for every pair i, j , there is a path from S_i to S_j that uses no variable occurring in the other atom of \mathcal{T} .



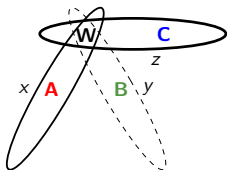
Where is the triad?
Path from A to B.

Triads

What do the triangle query and the tripod query have in common?

Definition (triad)

A **triad** is a set of three *endogenous atoms*, $\mathcal{T} = \{S_0, S_1, S_2\}$ such that for every pair i, j , there is a path from S_i to S_j that uses no variable occurring in the other atom of \mathcal{T} .



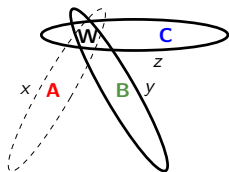
Where is the triad?
Path from A to C.

Triads

What do the triangle query and the tripod query have in common?

Definition (triad)

A **triad** is a set of three *endogenous atoms*, $\mathcal{T} = \{S_0, S_1, S_2\}$ such that for every pair i, j , there is a path from S_i to S_j that uses no variable occurring in the other atom of \mathcal{T} .



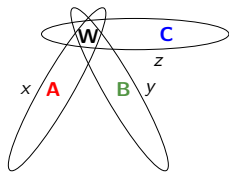
Where is the triad?
Path from B to C .

Triads

What do the triangle query and the tripod query have in common?

Definition (triad)

A **triad** is a set of three *endogenous atoms*, $\mathcal{T} = \{S_0, S_1, S_2\}$ such that for every pair i, j , there is a path from S_i to S_j that uses no variable occurring in the other atom of \mathcal{T} .



Where is the triad?

$$\mathcal{T} = \{A, B, C\}$$

Triads \rightarrow hardness

Lemma

If q has a triad, then $\text{RES}(q)$ is NP-complete.

Triads \rightarrow hardness

Lemma

If q has a triad, then $\text{RES}(q)$ is NP-complete.

Important remark

It is easy to check if a query has a triad!

Triads \rightarrow hardness

Lemma

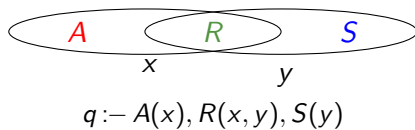
If q has a triad, then $\text{RES}(q)$ is NP-complete.

Important remark

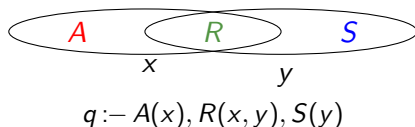
It is easy to check if a query has a triad!

What if a query does not have a triad?

Linear queries



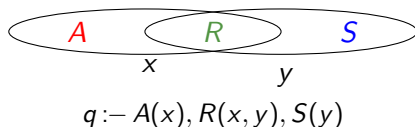
Linear queries



Definition

A query q is **linear** if its atoms may be arranged in a linear order such that each variable occurs in a contiguous sequence of atoms.

Linear queries



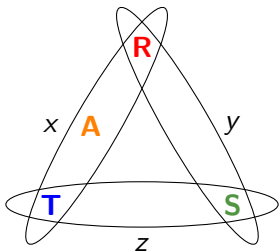
Definition

A query q is **linear** if its atoms may be arranged in a linear order such that each variable occurs in a contiguous sequence of atoms.

Fact [Meliou et.al., VLDB10]

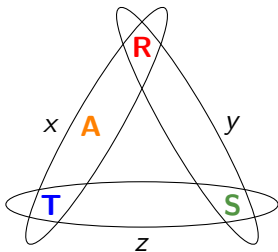
For any linear sj-free CQ q , $\text{RES}(q)$ is in P. (Reduction to network flow.)

Is there a triad in the following query?



$q_{\text{rats}} :- A(x), R(x, y), S(y, z), T(z, x)$

Is there a triad in the following query?



$q_{rats} :- A(x), R(x, y), S(y, z), T(z, x)$

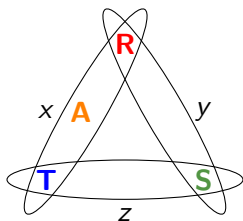
RES(q_{rats}) is in P!

Domination

T and R are **dominated** by A in q_{rats} . This guarantees we don't need tuples from R, T in a minimum contingency set.

Domination

T and R are **dominated** by A in q_{rats} . This guarantees we don't need tuples from R, T in a minimum contingency set.



$$q_{\text{rats}} := A(x), R^x(x, y), S(y, z), T^x(z, x)$$

No triads \rightarrow easy

Lemma

If q is a query in normal form with no triads, we can transform it into a linear query q' , such that $\text{RES}(q) \leq \text{RES}(q')$. Therefore $\text{RES}(q)$ is in P.

Dichotomy for Resilience - sj-free CQ

Theorem

Let q be an sj-free CQ and $\text{nf}(q)$ its normal form.

- If $\text{nf}(q)$ has a triad, then $\text{RES}(q)$ is NP-complete
- If $\text{nf}(q)$ does **not** have a triad, then $\text{RES}(q)$ is in P.

Adding functional dependencies

$q_T :- A(x), B(y), C(z), W(x, y, z)$

A		B		C		W			
	X		Y		Z	X	Y	Z	
a_1	1	b_1	2	c_1	5	w_1	1	2	5
a_2	3	b_2	4	c_2	6	w_2	1	2	6
						w_3	3	4	5
						w_4	1	4	5

Adding functional dependencies

$q_T :- A(x), B(y), C(z), W(x, y, z)$, $fd = W : x \rightarrow y$

	A	B	C		W		
	X	Y	Z		X	Y	Z
a_1	1	b_1 2	c_1 5	w_1	1	2	5
a_2	3	b_2 4	c_2 6	w_2	1	2	6
				w_3	3	4	5
				w_4	1	4	5

Adding functional dependencies

$q_T :- A(x), B(y), C(z), W(x, y, z)$, $fd = W : x \rightarrow y$

	A	B	C		W		
	X	Y	Z		X	Y	Z
a_1	1	b_1 2	c_1 5	w_1	1	2	5
a_2	3	b_2 4	c_2 6	w_2	1	2	6
				w_3	3	4	5
				w_4	1	4	5

- FDs constrain the databases we can consider

Adding functional dependencies

$q_T :- A(x), B(y), C(z), W(x, y, z)$, $fd = W : x \rightarrow y$

	A	B	C		W		
	X	Y	Z		X	Y	Z
a_1	1	b_1 2	c_1 5	w_1	1	2	5
a_2	3	b_2 4	c_2 6	w_2	1	2	6
				w_3	3	4	5
				w_4	1	4	5

- FDs constrain the databases we can consider
- FDs can reduce the complexity of resilience

Adding functional dependencies

Transform the query based on FDs

Adding functional dependencies

Transform the query based on FDs \rightarrow *induced rewrites*

Adding functional dependencies

Transform the query based on FDs \rightarrow *induced rewrites*

$$q_T = A(x), B(y), C(z), W(x, y, z), W : x \rightarrow y$$



$$q_T^* = \mathbf{A'(x, y)}, B(y), C(z), W(x, y, z), W : x \rightarrow y$$



$$q_T^* = A'(x, y), B(y), C(z), W(x, y, z)$$

Adding functional dependencies

Transform the query based on FDs \rightarrow *induced rewrites*

$$q_T = A(x), B(y), C(z), W(x, y, z), W : x \rightarrow y$$



$$q_T^* = \mathbf{A'}(x, y), B(y), C(z), W(x, y, z), W : x \rightarrow y$$



$$q_T^* = A'(x, y), B(y), C(z), W(x, y, z)$$

$\text{RES}(q_T, \varphi)$ is in P.

Adding functional dependencies

Induced rewrites

Let q be a query and $\bar{v} \rightarrow u \in \Phi$ be a functional dependency. We write $(q; \Phi) \rightsquigarrow (q'; \Phi)$ to mean that q' is the result of adding the dependent variable u to some relation that contains all the determinant variables \bar{v} . After applying all possible rewrites, we obtain query q^* , which we call **closed query**.

Adding functional dependencies

Induced rewrites

Let q be a query and $\bar{v} \rightarrow u \in \Phi$ be a functional dependency. We write $(q; \Phi) \rightsquigarrow (q'; \Phi)$ to mean that q' is the result of adding the dependent variable u to some relation that contains all the determinant variables \bar{v} . After applying all possible rewrites, we obtain query q^* , which we call **closed query**.

Lemma

Let q^ be q after all possible induced rewrites have been applied. Then $\text{RES}(q; \Phi) \equiv \text{RES}(q^*; \Phi) \equiv \text{RES}(q^*)$.*

Dichotomy for resilience - sj-free CQ + FD

Theorem

Let $(q; \Phi)$ be a sj-free CQ with functional dependencies. Let (q^*, Φ) be its closure under induced rewrites, and such that all dominated atoms of q^* are exogenous.

- If q^* has a triad then $\text{RES}(q; \Phi)$ is NP-complete.
- If q^* does **not** have a triad, then $\text{RES}(q; \Phi)$ is in P.

Future Directions

- Resilience for CQ with self-joins - (q_{vc})
- Deletion propagation: view side-effects for CQ with self-joins
 - Dichotomy results for CQ without self-joins [Kimelfeld et.al., PODS11]
 - Extended to functional dependencies [Kimelfeld, PODS12]
- Characterize the complexity of the parts of the problem that are in P , cf. [Allender, et. al.]