# The Complexity of Resilience and Responsibility for Conjunctive Queries

Neil Immerman

College of Information and Computer Sciences

UMass Amherst

The **resilience** of a boolean query with respect to a database, $D$, is the minimum number of tuples that must be removed from $D$ to make the query false.

The **resilience** of a boolean query with respect to a database, $D$, is the minimum number of tuples that must be removed from $D$ to make the query false.

Resilience is crucial to figuring out **why** a certain tuple, **t**, occurs in the answer to a query or view, $q$, on a database, $D$.

The **resilience** of a boolean query with respect to a database, $D$, is the minimum number of tuples that must be removed from $D$ to make the query false.

Resilience is crucial to figuring out **why** a certain tuple, **t**, occurs in the answer to a query or view, $q$, on a database, $D$. and to computing the **minimum change** needed to remove **t** from the view.

The **resilience** of a boolean query with respect to a database, $D$, is the minimum number of tuples that must be removed from $D$ to make the query false.

Resilience is crucial to figuring out **why** a certain tuple, **t**, occurs in the answer to a query or view, $q$, on a database, $D$. and to computing the **minimum change** needed to remove **t** from the view.

Often $D = D^x \cup D^n$ is partly **exogenous** and partly **endogenous**.

The **resilience** of a boolean query with respect to a database, $D$, is the minimum number of tuples that must be removed from $D$ to make the query false.

Resilience is crucial to figuring out **why** a certain tuple, **t**, occurs in the answer to a query or view, $q$, on a database, $D$. and to computing the **minimum change** needed to remove **t** from the view.

Often $D = D^x \cup D^n$ is partly **exogenous** and partly **endogenous**.

Treat exogenous part as fixed, beyond our control; only consider possible changes to the endogenous part.

# Resilience as a decision problem

$$\mathrm{RES}(q) \quad = \quad \left\{ (D, k) \mid \exists \Gamma \subseteq D^n \, (D - \Gamma) \not\models q \,\&\, |\Gamma| \leq k \right\}$$

## Resilience as a decision problem

$$\text{RES}(q) \quad = \quad \big\{ (D, k) \mid \exists \Gamma \subseteq D^n \, (D - \Gamma) \not\models q \,\&\, |\Gamma| \leq k \big\}$$

**Example:** $\quad q_{vc} :- \; V(x) \, E(x, y) \, V(y)$

## Resilience as a decision problem

$$\text{RES}(q) \quad = \quad \big\{(D, k) \mid \exists \Gamma \subseteq D^n \ (D - \Gamma) \not\models q \ \& \ |\Gamma| \leq k\big\}$$

**Example:** $q_{\text{vc}} :- \ V(x) \, E(x, y) \, V(y)$

**Prop:** $\text{RES}(q_{\text{vc}})$ is NP complete.

## Resilience as a decision problem

$$\text{RES}(q) \quad = \quad \left\{ (D, k) \mid \exists \Gamma \subseteq D^n \, (D - \Gamma) \not\models q \ \& \ |\Gamma| \leq k \right\}$$

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

**Prop:** $\text{RES}(q_{vc})$ is NP complete.

Proof.

$\text{RES}(q_{vc})$ is exactly the vertex cover problem: how many vertices need we remove so that no edges remain. $\qquad\square$

# Resilience as a decision problem

$$\text{RES}(q) \quad = \quad \big\{ (D, k) \mid \exists \Gamma \subseteq D^n \, (D - \Gamma) \not\models q \,\&\, |\Gamma| \leq k \big\}$$

**Example:** $q_{\text{vc}} :- V(x) \, E(x, y) \, V(y)$

**Prop:** $\text{RES}(q_{\text{vc}})$ is NP complete.

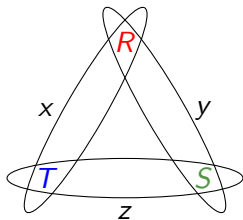## Proof.

$\text{RES}(q_{\text{vc}})$ is exactly the vertex cover problem: how many vertices need we remove so that no edges remain. $\qquad\qquad\qquad\square$

$q_{\text{vc}}$ has a **self join**.

# Resilience as a decision problem

$$\text{RES}(q) \quad = \quad \left\{ (D, k) \mid \exists \Gamma \subseteq D^n \, (D - \Gamma) \not\models q \,\&\, |\Gamma| \leq k \right\}$$

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

**Prop:** $\text{RES}(q_{vc})$ is NP complete.

Proof.

$\text{RES}(q_{vc})$ is exactly the vertex cover problem: how many vertices need we remove so that no edges remain. □
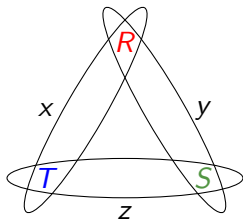
$q_{vc}$ has a **self join**.

**Goal:** Characterize the **complexity** of **resilience** for **sj-free conjunctive queries**.

# Triangle Query



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

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

Query **hypergraph**: relations are vertices;
variables are hyperedges

**Prop.**   RES($q_\triangle$) is NP-complete.

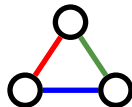**Prop.** $\text{RES}(q_\triangle)$ is NP-complete.

**Proof:** Reduce 3SAT to $\text{RES}(q_\triangle)$. Let $\psi = C_1 \wedge \cdots \wedge C_m$ be a 3-CNF formula, $\text{var}(\psi) = \{v_1, \ldots, v_n\}$

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

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

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

$$\psi = C_1 \wedge \cdots \wedge C_m \quad \text{var}(\psi) = \{v_1, \ldots, v_n\} \quad \psi \mapsto (D_\psi, k_\psi)$$
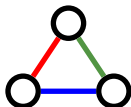
$$q_\triangle \ :- \ R(x,y), S(y,z), T(z,x)$$

$$(D_\psi, k_\psi) \in \text{RES}(q_\triangle) \ \Leftrightarrow \ \exists \Gamma \ |\Gamma| = k_\psi \wedge D_\psi - \Gamma \text{ has no}$$
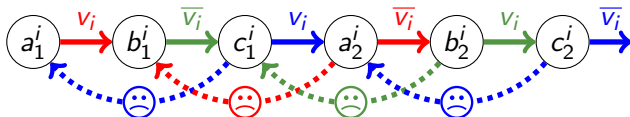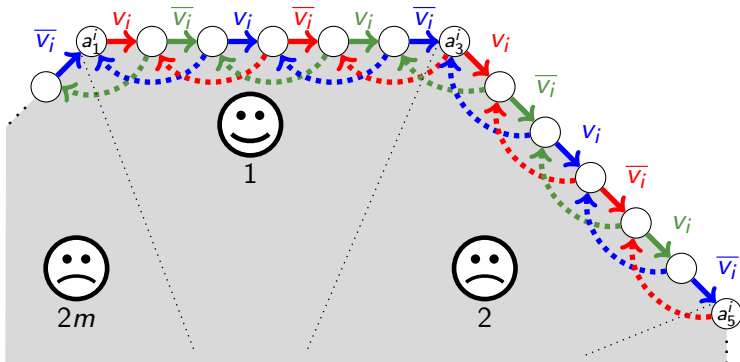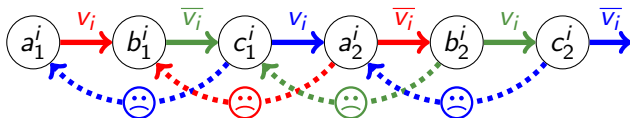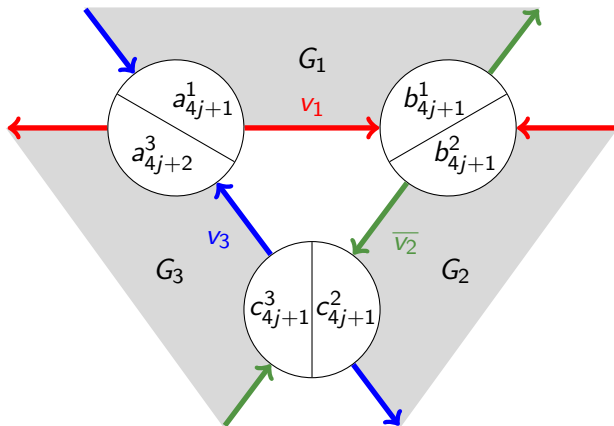
$D_\psi$ has one circular gadget $G_i$ for each variable $v_i$.

# In $G_i$ must choose all $v_i$'s or all $\overline{v_i}$'s

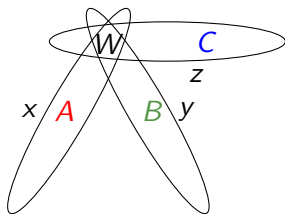For each clause, e.g., $C_j = (v_1 \vee \overline{v_2} \vee v_3)$, pick the $j$th occurrences of $v_1 \in G_1$, $\overline{v_2} \in G_2$ and $v_3 \in G_3$. Identify head of $v_1$ with tail of $\overline{v_2}$, head of $\overline{v_2}$ with tail of $v_3$, head of $v_3$ with tail of $v_1$



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

$$q_{\mathrm{T}} \; :- \; A(x), B(y), C(z), W(x, y, z)$$

# Tripod Query



$$q_{\mathrm{T}} \; :- \; A(x), B(y), C(z), W(x,y,z)$$

**Prop.** RES($q_{\mathrm{T}}$) is NP complete.

$$q_{\mathrm{T}} :\!- A(x), B(y), C(z), W(x, y, z)$$

# RES($q_{\mathrm{T}}$) is NP complete.

$$q_{\mathrm{T}} \; :- \; A(x), B(y), C(z), W(x, y, z)$$

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

$A$ **dominates** $W$.

# RES($q_T$) is NP complete.

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

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

$A$ **dominates** $W$.

**Prop.** 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** $W$.

**Prop.** 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**.

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

# $\mathrm{RES}(q_{\mathrm{T}})$ is NP complete.

$$q_\triangle \quad :- \quad R(x,y), S(y,z), T(z,x)$$
$$q_{\mathrm{T}} \quad :- \quad A(x), B(y), C(z), W^x(x,y,z)$$

**Proof:** Show $\mathrm{RES}(q_\triangle) \leq \mathrm{RES}(q_{\mathrm{T}})$

# RES($q_\mathrm{T}$) is NP complete.

$$
\begin{aligned}
q_\triangle \quad &:- \quad R(x,y), S(y,z), T(z,x) \\
q_\mathrm{T} \quad &:- \quad A(x), B(y), C(z), W^x(x,y,z)
\end{aligned}
$$

**Proof:** Show RES($q_\triangle$) $\leq$ RES($q_\mathrm{T}$)

Let $(D, k)$ be an instance of RES($q_\triangle$).

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

# RES($q_T$) is NP complete.

$$q_\triangle \quad :- \quad R(x,y), S(y,z), T(z,x)$$
$$q_T \quad :- \quad A(x), B(y), C(z), W^x(x,y,z)$$

**Proof:** Show RES($q_\triangle$) $\leq$ RES($q_T$)

Let $(D, k)$ be an instance of RES($q_\triangle$).

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

$$A = \big\{ \langle ab \rangle \mid R(a,b) \in D \big\}$$
$$B = \big\{ \langle bc \rangle \mid S(b,c) \in D \big\}$$
$$C = \big\{ \langle ca \rangle \mid T(c,a) \in D \big\}$$
$$W^x = \big\{ (\langle ab \rangle, \langle bc \rangle, \langle ca \rangle) \mid a, b, c \in \text{dom}(D) \big\}$$

# RES($q_{\mathrm{T}}$) is NP complete.

$$q_{\triangle} \quad :- \quad R(x,y), S(y,z), T(z,x)$$
$$q_{\mathrm{T}} \quad :- \quad A(x), B(y), C(z), W^{\times}(x,y,z)$$

**Proof:** Show RES($q_{\triangle}$) $\leq$ RES($q_{\mathrm{T}}$)

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

$$(D, k) \mapsto (D', k) \qquad D' \stackrel{\mathrm{def}}{=} (A, B, C, W^{\times})$$

$$A = \{\langle ab \rangle \mid R(a,b) \in D\}$$
$$B = \{\langle bc \rangle \mid S(b,c) \in D\}$$
$$C = \{\langle ca \rangle \mid T(c,a) \in D\}$$
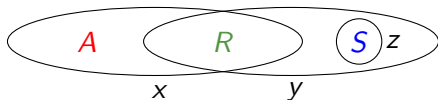$$W^{\times} = \{(\langle ab \rangle, \langle bc \rangle, \langle ca \rangle) \mid a, b, c \in \mathrm{dom}(D)\}$$

**Claim** $\quad (D, k) \in$ RES($q_{\triangle}$) $\quad \Leftrightarrow \quad (D', k) \in$ RES($q_{\mathrm{T}}$). $\qquad \square$

**Def.** A query is **linear** if all of the vertices of its hypergraph can be drawn along a straight line with all of its hyperedges convex.

For example, the following query is linear:

$$q :- A(x), R(x, y), S(y, z)$$

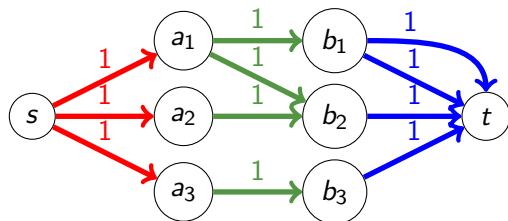**Prop.** For any linear sj-free conjuctive query $q$, $\mathrm{RES}(q) \in \mathrm{P}$.

# Linear Queries are Easy

**Prop.** For any linear sj-free conjunctive query $q$, $\mathrm{RES}(q) \in \mathrm{P}$.

**Proof:** Use Network Flow.

$\mathrm{RES}(D, q)$ is the min cut of corresponding network.
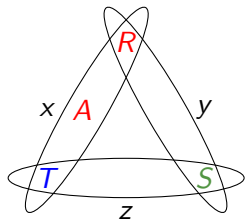


$$q :- \quad\quad A(x) \quad\quad R(x, y) \quad\quad S(y, z)$$

$$q_{\mathrm{rats}} \quad :- \quad A(x), R(x, y), S(y, z), T(z, x)$$
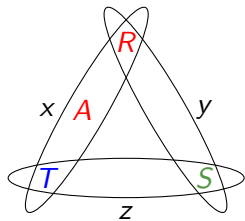
Is Rats **hard** or **easy** ?

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

Is Rats **hard** or **easy** ?

$$
\begin{aligned}
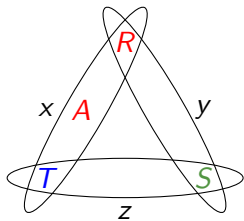q_{\text{rats}} &:- & A(x), R(x,y), S(y,z), T(z,x) & \\
q_1 &\equiv & A(x), R^{\times}(x,y), S(y,z), T^{\times}(z,x) & \quad \text{Domination} \\
\text{RES}(q_{\text{rats}}) &\equiv & \text{RES}(q_1) &
\end{aligned}
$$

Is Rats **hard** or **easy** ?

$$
\begin{aligned}
q_{\mathrm{rats}} \quad &:- \quad A(x), R(x,y), S(y,z), T(z,x) \\
q_1 \quad &\equiv \quad A(x), R^{\times}(x,y), S(y,z), T^{\times}(z,x) \qquad \text{Domination} \\
\mathrm{RES}(q_{\mathrm{rats}}) \quad &\equiv \quad \mathrm{RES}(q_1) \\
q_2 \quad &\equiv \quad A(x), R^{\times}(x,y,z), S(y,z), T^{\times}(z,x) \qquad \text{Dissociation} \\
\mathrm{RES}(q_1) \quad &\leq \quad \mathrm{RES}(q_2)
\end{aligned}
$$

Is Rats **hard** or **easy** ?

$$
\begin{aligned}
q_{\mathrm{rats}} \quad &:- \quad A(x), R(x,y), S(y,z), T(z,x) \\
q_1 \quad &\equiv \quad A(x), R^{\times}(x,y), S(y,z), T^{\times}(z,x) \qquad \text{Domination} \\
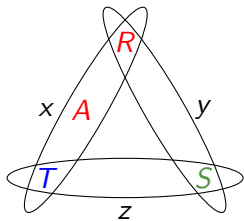\mathrm{RES}(q_{\mathrm{rats}}) \quad &\equiv \quad \mathrm{RES}(q_1) \\
q_2 \quad &\equiv \quad A(x), R^{\times}(x,y,z), S(y,z), T^{\times}(z,x) \qquad \text{Dissociation} \\
\mathrm{RES}(q_1) \quad &\leq \quad \mathrm{RES}(q_2) \\
q_3 \quad &\equiv \quad A(x), R^{\times}(x,y,z), S(y,z), T^{\times}(z,x,y) \qquad \text{Dissociation} \\
\mathrm{RES}(q_2) \quad &\leq \quad \mathrm{RES}(q_3)
\end{aligned}
$$

Is Rats **hard** or **easy** ?

$$
\begin{aligned}
q_{\mathrm{rats}} \quad &:- \quad A(x), R(x,y), S(y,z), T(z,x) \\
q_1 \quad &\equiv \quad A(x), R^{\times}(x,y), S(y,z), T^{\times}(z,x) \qquad && \text{Domination} \\
\mathrm{RES}(q_{\mathrm{rats}}) \quad &\equiv \quad \mathrm{RES}(q_1) \\
q_2 \quad &\equiv \quad A(x), R^{\times}(x,y,z), S(y,z), T^{\times}(z,x) \qquad && \text{Dissociation} \\
\mathrm{RES}(q_1) \quad &\leq \quad \mathrm{RES}(q_2) \\
q_3 \quad &\equiv \quad A(x), R^{\times}(x,y,z), S(y,z), T^{\times}(z,x,y) \qquad && \text{Dissociation} \\
\mathrm{RES}(q_2) \quad &\leq \quad \mathrm{RES}(q_3) \\
q_4 \quad &\equiv \quad A(x), R^{\times}(x,y,z), S(y,z) \qquad && \text{Repetition} \\
\mathrm{RES}(q_3) \quad &\equiv \quad \mathrm{RES}(q_4)
\end{aligned}
$$

Is Rats **hard** or **easy** ?

$$
\begin{array}{rcll}
q_{\mathrm{rats}} & :- & A(x), R(x,y), S(y,z), T(z,x) & \\
q_1 & \equiv & A(x), R^{\times}(x,y), S(y,z), T^{\times}(z,x) & \text{Domination} \\
\mathrm{RES}(q_{\mathrm{rats}}) & \equiv & \mathrm{RES}(q_1) & \\
q_2 & \equiv & A(x), R^{\times}(x,y,z), S(y,z), T^{\times}(z,x) & \text{Dissociation} \\
\mathrm{RES}(q_1) & \leq & \mathrm{RES}(q_2) & \\
q_3 & \equiv & A(x), R^{\times}(x,y,z), S(y,z), T^{\times}(z,x,y) & \text{Dissociation} \\
\mathrm{RES}(q_2) & \leq & \mathrm{RES}(q_3) & \\
q_4 & \equiv & A(x), R^{\times}(x,y,z), S(y,z) & \text{Repetition} \\
\mathrm{RES}(q_3) & \equiv & \mathrm{RES}(q_4) & \\
& & q_4 \text{ is } \textbf{linear} \text{ and therefore } \textbf{easy}! &
\end{array}
$$

Is Rats **hard** or **easy** ?

$$
\begin{aligned}
q_{\mathrm{rats}} \quad &:- \quad A(x), R(x,y), S(y,z), T(z,x) \\
q_1 \quad &\equiv \quad A(x), R^{\times}(x,y), S(y,z), T^{\times}(z,x) && \text{Domination} \\
\mathrm{RES}(q_{\mathrm{rats}}) \quad &\equiv \quad \mathrm{RES}(q_1) \\
q_2 \quad &\equiv \quad A(x), R^{\times}(x,y,z), S(y,z), T^{\times}(z,x) && \text{Dissociation} \\
\mathrm{RES}(q_1) \quad &\leq \quad \mathrm{RES}(q_2) \\
q_3 \quad &\equiv \quad A(x), R^{\times}(x,y,z), S(y,z), T^{\times}(z,x,y) && \text{Dissociation} \\
\mathrm{RES}(q_2) \quad &\leq \quad \mathrm{RES}(q_3) \\
q_4 \quad &\equiv \quad A(x), R^{\times}(x,y,z), S(y,z) && \text{Repetition} \\
\mathrm{RES}(q_3) \quad &\equiv \quad \mathrm{RES}(q_4) \\
& \quad\;\; q_4 \text{ is } \textbf{linear} \text{ and therefore } \textbf{easy}! \\
\mathrm{RES}(q_{\mathrm{rats}}) \quad &\leq \quad \mathrm{RES}(q_5) && q_{\mathrm{rats}} \text{ is } \textbf{easy}!
\end{aligned}
$$

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



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

$q_\mathrm{T} :- A(x), B(y), C(z), W^\times(x, y, z)$

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



$q_\triangle :- R(x, y), S(y, z), T(z, x)$      $q_\mathrm{T} :- A(x), B(y), C(z), W^\times(x, y, z)$

**Def.** 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 $\mathcal{T}$.
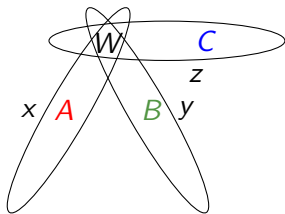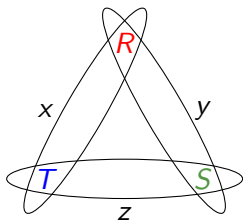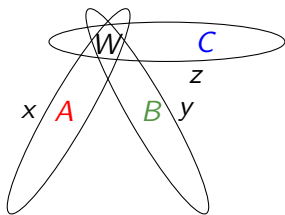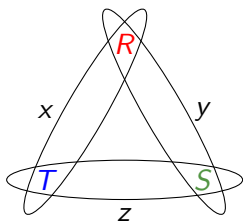
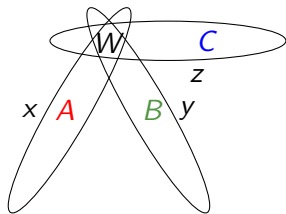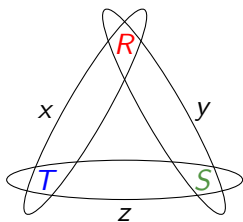# What do the triangle and the tripod have in common?



$q_\triangle :- R(x, y), S(y, z), T(z, x)$     $q_T :- A(x), B(y), C(z), W^x(x, y, z)$

**Def.** 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 $\mathcal{T}$.

$\{R, S, T\}$ is a triad in $q_\triangle$.

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



$q_{\triangle} :\text{--} R(x,y), S(y,z), T(z,x)$ $\qquad q_{\mathrm{T}} :\text{--} A(x), B(y), C(z), W^{\times}(x,y,z)$

**Def.** 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 $\mathcal{T}$.

$\{R, S, T\}$ is a triad in $q_{\triangle}$.

$\{A, B, C\}$ is a triad in $q_{\mathrm{T}}$.

**Lemma** Let $q$ be an sj-free conjunctive query where all dominated atoms are exogenous. If $q$ has a triad, then $\text{RES}(q)$ is $\text{NP}$-complete.

**Lemma** Let $q$ be an sj-free conjunctive query where all dominated atoms are exogenous. If $q$ has a triad, then $\mathrm{RES}(q)$ is $\mathrm{NP}$-complete.

**Proof:** Show $\mathrm{RES}(q_\triangle) \leq \mathrm{RES}(q)$

$\square$

**Lemma** Let $q$ be an sj-free conjunctive query that has no triad. Then $\mathrm{RES}(q) \in \mathrm{P}$.

**Lemma** Let $q$ be an sj-free conjunctive query that has no triad. Then $\text{RES}(q) \in \text{P}$.

**Proof:** By induction on the number of *endogenous* atoms in $q$ that we can transform it into a linear query by using dissociations.

**Lemma** Let $q$ be an sj-free conjunctive query that has no triad. Then $\mathrm{RES}(q) \in \mathrm{P}$.

**Proof:** By induction on the number of *endogenous* atoms in $q$ that we can transform it into a linear query by using dissociations.

*Inductive case*: assume true for triad-free queries with $n$ endogenous atoms. Let $q_{n+1}$ be triad-free and have $n + 1$ endogenous atoms.

**Lemma** Let $q$ be an sj-free conjunctive query that has no triad. Then $\mathrm{RES}(q) \in \mathrm{P}$.

**Proof:** By induction on the number of *endogenous* atoms in $q$ that we can transform it into a linear query by using dissociations.
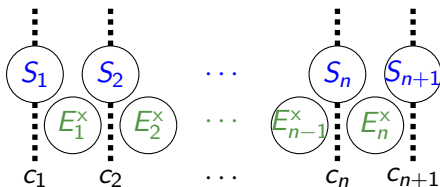
*Inductive case*: assume true for triad-free queries with $n$ endogenous atoms. Let $q_{n+1}$ be triad-free and have $n + 1$ endogenous atoms.

Since there is no triad, we can linearize the endogenous atoms:



$\square$

**Dichotomy Theorem for Resilience:** Let $q$ be a sj-free conjunctive query all of whose dominated atoms are exogenous. If $q$ has a triad then $\text{RES}(q)$ is NP complete. Otherwise, $\text{RES}(q) \in \text{P}$.

# Extend to Databases with Functional Dependencies

# Extend to Databases with Functional Dependencies

**induced rewrites** preserve complexity of resilience:

$q :- R(x, y), S(y, z), T(z, x); \ x \mapsto y$

$q^* :- R(x, y), S(y, z), T(z, x, y); \ x \mapsto y$

Let $q^*$ be $q$ after all possible induced rewrites have been applied.

**induced rewrites** preserve complexity of resilience:

$q :\!- R(x,y), S(y,z), T(z,x); \ x \mapsto y$

$q^* :\!- R(x,y), S(y,z), T(z,x,y); \ x \mapsto y$

Let $q^*$ be $q$ after all possible induced rewrites have been applied.

**Lemma:** $\quad \mathrm{RES}(q) \equiv \mathrm{RES}(q^*)$

# Extend to Databases with Functional Dependencies

**induced rewrites** preserve complexity of resilience:

$q :\!- R(x,y), S(y,z), T(z,x); \; x \mapsto y$

$q^* :\!- R(x,y), S(y,z), T(z,x,y); \; x \mapsto y$

Let $q^*$ be $q$ after all possible induced rewrites have been applied.

**Lemma:** $\text{RES}(q) \equiv \text{RES}(q^*)$

**Dichotomy Theorem for Resilience with FD's** Let $q^*$ be an sj-free conjunctive query with FD's, all possible induced rewrites applied and all dominated atoms are exogenous. If $q^*$ has a triad then $\text{RES}(q)$ is NP complete. Otherwise, $\text{RES}(q) \in \mathrm{P}$.

# Extend to Databases with Functional Dependencies

**induced rewrites** preserve complexity of resilience:

$q :\!- R(x, y), S(y, z), T(z, x); \ x \mapsto y$

$q^* :\!- R(x, y), S(y, z), T(z, x, y); \ x \mapsto y$

Let $q^*$ be $q$ after all possible induced rewrites have been applied.

**Lemma:** $\quad \mathrm{RES}(q) \equiv \mathrm{RES}(q^*)$

**Dichotomy Theorem for Resilience with FD's** Let $q^*$ be an sj-free conjunctive query with FD's, all possible induced rewrites applied and all dominated atoms are exogenous. If $q^*$ has a triad then $\mathrm{RES}(q)$ is NP complete. Otherwise, $\mathrm{RES}(q) \in \mathrm{P}$.

**Corollary** Induced rewrites characterized the effect of FD's:

$$\mathrm{RES}(q; \Phi) \quad \equiv \quad \mathrm{RES}(q^*; \Phi) \quad \equiv \quad \mathrm{RES}(q^*)$$

# Future Directions

- Extend characterization of complexity of resilience to conjunctive queries with self joins.

# Future Directions

- Extend characterization of complexity of resilience to conjunctive queries with self joins.
- Extend to sj's with FD's.

# Future Directions

- Extend characterization of complexity of resilience to conjunctive queries with self joins.
- Extend to sj's with FD's.
- Extend to the complexity of "view side-effects" problem.

# Future Directions

- Extend characterization of complexity of resilience to conjunctive queries with self joins.
- Extend to sj's with FD's.
- Extend to the complexity of "view side-effects" problem.
- Characterize the complexity of the parts of the problem that are in $P$, cf. [Allender, et. al.]

# Future Directions

- Extend characterization of complexity of resilience to conjunctive queries with self joins.
- Extend to sj's with FD's.
- Extend to the complexity of "view side-effects" problem.
- Characterize the complexity of the parts of the problem that are in $P$, cf. [Allender, et. al.]
- **Understand & explain Dichotomy Phenomenon**