# Probabilistic Databases: The Dichotomy of PQE

Antoine Amarilli

What is the complexity of $\mathrm{PQE}(Q)$ depending on the query $Q$?

What is the complexity of $\mathrm{PQE}(Q)$ depending on the query *Q*?

$\rightarrow$ Recall that we study **data complexity**, i.e., *Q* is **fixed** and the input is the **data**

# Research goal: Understanding the complexity of PQE

What is the complexity of $\mathrm{PQE}(Q)$ depending on the query *Q*?

$\rightarrow$ Recall that we study **data complexity**, i.e., *Q* is **fixed** and the input is the **data**

For example:

- For $Q : R(x)$ the problem is **easy** (PTIME)
- For $Q : R(x), S(x,y), T(y)$ the problem is **hard** (#P-hard)

What is the complexity of $\mathrm{PQE}(Q)$ depending on the query *Q*?

$\rightarrow$ Recall that we study **data complexity**, i.e., *Q* is **fixed** and the input is the **data**

For example:

- For *Q* : *R*(*x*) the problem is **easy** (PTIME)
- For *Q* : *R*(*x*), *S*(*x*, *y*), *T*(*y*) the problem is **hard** (#P-hard)

We will present the **dichotomy** of [Dalvi and Suciu, 2007, Dalvi and Suciu, 2012]:

## Research goal: Understanding the complexity of PQE

What is the complexity of $\mathrm{PQE}(Q)$ depending on the query *Q*?

  $\rightarrow$ Recall that we study **data complexity**, i.e., *Q* is **fixed** and the input is the **data**

For example:

- For *Q* : *R*(*x*) the problem is **easy** (PTIME)
- For *Q* : *R*(*x*), *S*(*x*, *y*), *T*(*y*) the problem is **hard** (#P-hard)

We will present the **dichotomy** of [Dalvi and Suciu, 2007, Dalvi and Suciu, 2012]:

- Small dichotomy: **conjunctive queries** that are **self-join-free** and **arity-two**

What is the complexity of $\mathrm{PQE}(Q)$ depending on the query *Q*?

$\rightarrow$ Recall that we study **data complexity**, i.e., *Q* is **fixed** and the input is the **data**

For example:

- For *Q* : *R*(*x*) the problem is **easy** (PTIME)
- For *Q* : *R*(*x*), *S*(*x*, *y*), *T*(*y*) the problem is **hard** (#P-hard)

We will present the **dichotomy** of [Dalvi and Suciu, 2007, Dalvi and Suciu, 2012]:

- Small dichotomy: **conjunctive queries** that are **self-join-free** and **arity-two**
- Large dichotomy: arbitrary **unions of conjunctive queries**

## Research goal: Understanding the complexity of PQE

What is the complexity of $\mathrm{PQE}(Q)$ depending on the query *Q*?

$\rightarrow$ Recall that we study **data complexity**, i.e., *Q* is **fixed** and the input is the **data**

For example:

- For *Q* : *R*(*x*) the problem is **easy** (PTIME)
- For *Q* : *R*(*x*), *S*(*x*, *y*), *T*(*y*) the problem is **hard** (#P-hard)

We will present the **dichotomy** of [Dalvi and Suciu, 2007, Dalvi and Suciu, 2012]:

- Small dichotomy: **conjunctive queries** that are **self-join-free** and **arity-two**
- Large dichotomy: arbitrary **unions of conjunctive queries**

Result of the form:

if *Q has a certain form* then $\mathrm{PQE}(Q)$ is in PTIME, otherwise it is #P-hard

- Conjunctive query (CQ): existentially quantified conjunction of atoms

## The "small" Dalvi and Suciu dichotomy

- **Conjunctive query (CQ)**: existentially quantified conjunction of atoms

- **Arity-two**: all relations are binary
  - We represent the queries as **graphs**: $R(x, y), S(y, z)$ is $x \longrightarrow y \longrightarrow z$

## The "small" Dalvi and Suciu dichotomy

- **Conjunctive query (CQ)**: existentially quantified conjunction of atoms

- **Arity-two**: all relations are binary
  - We represent the queries as **graphs**: $R(x, y), S(y, z)$ is $x \longrightarrow y \longrightarrow z$

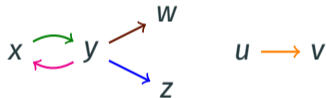- **Self-join-free CQ**: only one edge of each color (no repeated color)

# The "small" Dalvi and Suciu dichotomy

- **Conjunctive query (CQ)**: existentially quantified conjunction of atoms
- **Arity-two**: all relations are binary
  - We represent the queries as **graphs**: $R(x,y), S(y,z)$ is $x \longrightarrow y \longrightarrow z$
- **Self-join-free CQ**: only one edge of each color (no repeated color)

---

**Theorem ([Dalvi and Suciu, 2007])**

*Let $Q$ be an arity-two self-join-free CQ:*

- *If $Q$ is a conjunction of stars, then* $\mathrm{PQE}(Q)$ *is in PTIME*
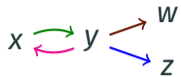- *Otherwise,* $\mathrm{PQE}(Q)$ *is #P-hard*

- A **star** is a CQ with a **separator variable** that occurs in all edges
- A **conjunction of stars** is a conjunction of one or several stars



The following is **not a star**: $x \longrightarrow y \longrightarrow z \longrightarrow w$

$$x \rightleftarrows y \begin{smallmatrix} \nearrow & w \\ \searrow & z \end{smallmatrix} \qquad u \longrightarrow v \qquad \text{How to solve } \mathrm{PQE}(Q) \text{ for } Q \text{ a conjunction of stars?}$$

$x \rightleftarrows y \diagdown \nearrow^{w}_{\searrow z}$   $u \longrightarrow v$   How to solve $\mathrm{PQE}(Q)$ for $Q$ a conjunction of stars?

$x \rightleftarrows y \nearrow^{w}_{\searrow z}$

- We consider each connected component separately
$\rightarrow$ **Independent conjunction** over the connected components

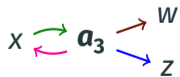$x \underset{\longleftarrow}{\overset{\longrightarrow}{}} y \overset{w}{\underset{z}{\nearrow \searrow}}$       $u \longrightarrow v$       How to solve $\mathrm{PQE}(Q)$ for $Q$ a conjunction of stars?

$x \underset{\longleftarrow}{\overset{\longrightarrow}{}} y \overset{w}{\underset{z}{\nearrow \searrow}}$

- We consider each connected component separately
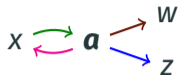- $\rightarrow$ **Independent conjunction** over the connected components

$x \underset{\longleftarrow}{\overset{\longrightarrow}{}} a_1 \overset{w}{\underset{z}{\nearrow \searrow}}$

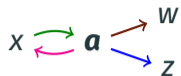$x \underset{\longleftarrow}{\overset{\longrightarrow}{}} a_2 \overset{w}{\underset{z}{\nearrow \searrow}}$

- We can test all possible values of the **separator variable**
- $\rightarrow$ **Independent disjunction** over the values of the separator

$x \underset{\longleftarrow}{\overset{\longrightarrow}{}} a_3 \overset{w}{\underset{z}{\nearrow \searrow}}$

$\vdots$

$x \rightleftarrows a \longrightarrow w$
$\phantom{x \rightleftarrows a} \searrow z$

$x \rightleftarrows a$

- For every match, we consider every **other variable** separately
$\rightarrow$ **Independent conjunction** over the variables

$x \rightleftarrows a \xrightarrow{} w$
$\searrow z$

$x \rightleftarrows a$

- For every match, we consider every **other variable** separately
$\rightarrow$ **Independent conjunction** over the variables

$b_1 \rightleftarrows a$
$b_2 \rightleftarrows a$
$b_3 \rightleftarrows a$
$\vdots$

- We consider every value for the **other variable**
$\rightarrow$ **Independent disjunction** over the possible assignments

$x \rightleftarrows a \begin{array}{c} \nearrow w \\ \searrow z \end{array}$

$x \rightleftarrows a$

- For every match, we consider every **other variable** separately
- $\rightarrow$ **Independent conjunction** over the variables

$b_1 \rightleftarrows a$
$b_2 \rightleftarrows a$
$b_3 \rightleftarrows a$
$\vdots$

- We consider every value for the **other variable**
- $\rightarrow$ **Independent disjunction** over the possible assignments

$b \rightarrow a$

- We consider every fact
- $\rightarrow$ **Independent conjunction** over the facts
- $\rightarrow$ Just **read the probability** of the ground fact $R(b, a)$.

Every arity-two self-join-free CQ which is **not a conjunction of stars** contains a pattern essentially like:

$$x \xrightarrow{\hspace{1cm}} y \xrightarrow{\hspace{1cm}} z \xrightarrow{\hspace{1cm}} w$$

Every arity-two self-join-free CQ which is **not a conjunction of stars** contains a pattern essentially like:

$$x \xrightarrow{\hspace{1cm}} y \xrightarrow{\hspace{1cm}} z \xrightarrow{\hspace{1cm}} w$$

We can **add facts with probability 1** to instances so the other facts are always satisfied, and focus on **only these three facts**

$\rightarrow$ Let us show #P-hardness of this query

Let us show that $\mathrm{PQE}(Q)$ is **#P-hard** for the CQ $Q :$ $x \longrightarrow y \longrightarrow z \longrightarrow w$

Let us show that $\mathrm{PQE}(Q)$ is **#P-hard** for the CQ $Q$ : $x \longrightarrow y \longrightarrow z \longrightarrow w$

- Reduce from the problem of **counting satisfying valuations** of a Boolean formula
  - e.g., given $(x \lor y) \land z$, compute that it has **3** satisfying valuations

Let us show that $\mathrm{PQE}(Q)$ is **#P-hard** for the CQ $Q$ : $x \longrightarrow y \longrightarrow z \longrightarrow w$

- Reduce from the problem of **counting satisfying valuations** of a Boolean formula
  - e.g., given $(x \lor y) \land z$, compute that it has **3** satisfying valuations

- This problem is already **#P-hard** for so-called **PP2DNF formulas**:

Let us show that $\mathrm{PQE}(Q)$ is **#P-hard** for the CQ $Q$ : $x \xrightarrow{\quad\quad} y \xrightarrow{\quad\quad} z \xrightarrow{\quad\quad} w$

- Reduce from the problem of **counting satisfying valuations** of a Boolean formula
  - e.g., given $(x \vee y) \wedge z$, compute that it has **3** satisfying valuations

- This problem is already **#P-hard** for so-called **PP2DNF formulas**:
  - **Positive** (no negation) and **Partitioned variables**: $X_1, \dots, X_n$ and $Y_1, \dots, Y_m$

Let us show that $\mathrm{PQE}(Q)$ is **#P-hard** for the CQ $Q$ : $x \longrightarrow y \longrightarrow z \longrightarrow w$

- Reduce from the problem of **counting satisfying valuations** of a Boolean formula
  - e.g., given $(x \vee y) \wedge z$, compute that it has **3** satisfying valuations

- This problem is already **#P-hard** for so-called **PP2DNF formulas**:
  - **Positive** (no negation) and **Partitioned variables**: $X_1, \ldots, X_n$ and $Y_1, \ldots, Y_m$
  - **2-DNF**: disjunction of clauses like $X_i \wedge Y_j$

Let us show that $\mathrm{PQE}(Q)$ is **#P-hard** for the CQ $Q$ : $x \xrightarrow{\hspace{1cm}} y \xrightarrow{\hspace{1cm}} z \xrightarrow{\hspace{1cm}} w$

- Reduce from the problem of **counting satisfying valuations** of a Boolean formula
  - e.g., given $(x \lor y) \land z$, compute that it has **3** satisfying valuations

- This problem is already **#P-hard** for so-called **PP2DNF formulas**:
  - **Positive** (no negation) and **Partitioned variables**: $X_1, \ldots, X_n$ and $Y_1, \ldots, Y_m$
  - **2-DNF**: disjunction of clauses like $X_i \land Y_j$

- Example: $\phi$ : $(X_1 \land Y_1) \lor (X_1 \land Y_2) \lor (X_2 \land Y_2) \lor (X_3 \land Y_1) \lor (X_3 \land Y_2)$

Reduce from **#PP2DNF** to $\mathrm{PQE}(Q)$ for CQ $Q$ : $x \xrightarrow{\hspace{0.8cm}} y \xrightarrow{\hspace{0.8cm}} z \xrightarrow{\hspace{0.8cm}} w$

Example: $\phi : (X_1 \wedge Y_1) \vee (X_1 \wedge Y_2) \vee (X_2 \wedge Y_2) \vee (X_3 \wedge Y_1) \vee (X_3 \wedge Y_2)$

Reduce from **#PP2DNF** to $\mathrm{PQE}(Q)$ for CQ $Q$ : $x \xrightarrow{\hspace{0.8cm}} y \xrightarrow{\hspace{0.8cm}} z \xrightarrow{\hspace{0.8cm}} w$

Example: $\phi : (X_1 \wedge Y_1) \vee (X_1 \wedge Y_2) \vee (X_2 \wedge Y_2) \vee (X_3 \wedge Y_1) \vee (X_3 \wedge Y_2)$

Build an **instance** $I_\phi$ from $\phi$:

Reduce from **#PP2DNF** to $\mathrm{PQE}(Q)$ for CQ $Q$ : $x \xrightarrow{\quad} y \xrightarrow{\quad} z \xrightarrow{\quad} w$

Example: $\phi : (X_1 \wedge Y_1) \vee (X_1 \wedge Y_2) \vee (X_2 \wedge Y_2) \vee (X_3 \wedge Y_1) \vee (X_3 \wedge Y_2)$

Build an **instance** $I_\phi$ from $\phi$:

$$a'_1 \xrightarrow{\ 1/2\ } a_1$$

$$a'_2 \xrightarrow{\ 1/2\ } a_2$$

$$a'_3 \xrightarrow{\ 1/2\ } a_3$$

Reduce from **#PP2DNF** to $\mathrm{PQE}(Q)$ for CQ $Q$ : $x \xrightarrow{\hspace{1cm}} y \xrightarrow{\hspace{1cm}} z \xrightarrow{\hspace{1cm}} w$

Example: $\phi : (X_1 \wedge Y_1) \vee (X_1 \wedge Y_2) \vee (X_2 \wedge Y_2) \vee (X_3 \wedge Y_1) \vee (X_3 \wedge Y_2)$

Build an **instance** $I_\phi$ from $\phi$:

$$a_1' \xrightarrow{1/2} a_1 \qquad\qquad\qquad b_1 \xrightarrow{1/2} b_1'$$

$$a_2' \xrightarrow{1/2} a_2$$

$$a_3' \xrightarrow{1/2} a_3 \qquad\qquad\qquad b_2 \xrightarrow{1/2} b_2'$$

Reduce from **#PP2DNF** to $\mathrm{PQE}(Q)$ for CQ $Q: x \longrightarrow y \longrightarrow z \longrightarrow w$

Example: $\phi: (X_1 \wedge Y_1) \vee (X_1 \wedge Y_2) \vee (X_2 \wedge Y_2) \vee (X_3 \wedge Y_1) \vee (X_3 \wedge Y_2)$

Build an **instance** $I_\phi$ from $\phi$:

Reduce from **#PP2DNF** to $\mathrm{PQE}(Q)$ for CQ $Q$ : $x \longrightarrow y \longrightarrow z \longrightarrow w$

Example: $\phi : (X_1 \wedge Y_1) \vee (X_1 \wedge Y_2) \vee (X_2 \wedge Y_2) \vee (X_3 \wedge Y_1) \vee (X_3 \wedge Y_2)$
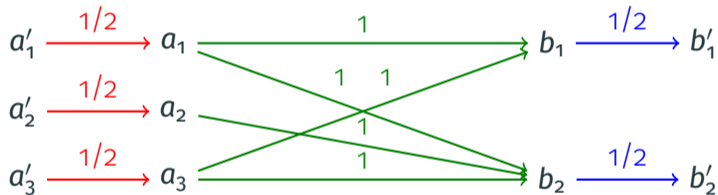
Build an **instance** $I_\phi$ from $\phi$:



**Idea:**

- **Valuations** of $\phi$ correspond to **possible worlds** of $I_\phi$

Reduce from **#PP2DNF** to $\mathrm{PQE}(Q)$ for CQ $Q$ : $x \xrightarrow{\quad} y \xrightarrow{\quad} z \xrightarrow{\quad} w$

Example: $\phi : (X_1 \wedge Y_1) \vee (X_1 \wedge Y_2) \vee (X_2 \wedge Y_2) \vee (X_3 \wedge Y_1) \vee (X_3 \wedge Y_2)$
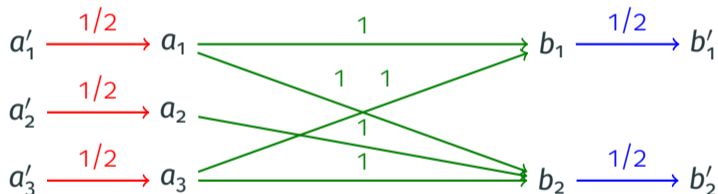
Build an **instance** $I_\phi$ from $\phi$:



**Idea:**

- **Valuations** of $\phi$ correspond to **possible worlds** of $I_\phi$
- A valuation **satisfies** $\phi$ iff the corresponding possible world **satisfies** $Q$

Reduce from **#PP2DNF** to $\mathrm{PQE}(Q)$ for CQ $Q$ : $x \longrightarrow y \longrightarrow z \longrightarrow w$

Example: $\phi : (X_1 \wedge Y_1) \vee (X_1 \wedge Y_2) \vee (X_2 \wedge Y_2) \vee (X_3 \wedge Y_1) \vee (X_3 \wedge Y_2)$
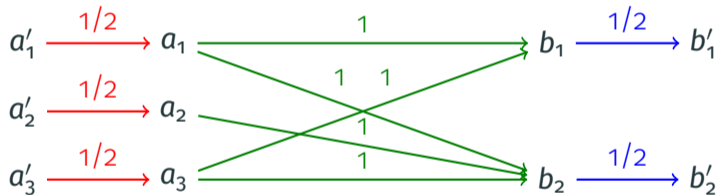
Build an **instance** $I_\phi$ from $\phi$:



**Idea:**

- **Valuations** of $\phi$ correspond to **possible worlds** of $I_\phi$
- A valuation **satisfies** $\phi$ iff the corresponding possible world **satisfies** $Q$
- $\rightarrow$ The **probability** of $Q$ on $I_\phi$ is the **number of accepting valuations of** $\phi$, divided by the number of valuations ($2^{-|\mathrm{Vars}|}$)

How can we extend beyond **arity-two queries**?

> **Theorem ([Dalvi and Suciu, 2007])**
>
> *Let $Q$ be a ~~arity-two~~ **self-join-free CQ**:*
>
> - *If $Q$ is ~~a conjunction of stars~~ **hierarchical**, then $\mathrm{PQE}(Q)$ is in **PTIME***
> - *Otherwise, $\mathrm{PQE}(Q)$ is **#P-hard***

## Extending beyond arity-two (2)

Class of **Hierarchical** CQs defined inductively:

- A query with **no variables** is hierarchical

## Extending beyond arity-two (2)

Class of **Hierarchical** CQs defined inductively:

- A query with **no variables** is hierarchical
- A **conjunction** of hierarchical connected components is hierarchical

## Extending beyond arity-two (2)

Class of **Hierarchical** CQs defined inductively:

- A query with **no variables** is hierarchical
- A **conjunction** of hierarchical connected components is hierarchical
- Induction case: for a connected CQ:
  - It must have a **separator variable** occurring in all atoms
  - If we remove this separator variable, the query must be hierarchical

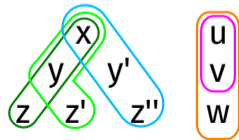Class of **Hierarchical** CQs defined inductively:

- A query with **no variables** is hierarchical
- A **conjunction** of hierarchical connected components is hierarchical
- Induction case: for a connected CQ:
  - It must have a **separator variable** occurring in all atoms
  - If we remove this separator variable, the query must be hierarchical

$$\exists x \left(\exists y \left(\exists z\ R_1(x, y, z)\right) \land \left(\exists z'\ R_2(x, y, z')\right)\right) \land \left(\exists y' \exists z''\ R_3(x, y', z'')\right)$$
$$\land \left(\exists u \left(\exists v\ R_4(u, v)\right) \land \left(\exists w\ R_5(u, v, w)\right)\right)$$

How does the proof change?

## Extending beyond arity-two (3)

How does the proof change?

- **Upper bound**: we can generalize the algorithm
  - **Independent AND** of connected components
  - **Independent OR** of possible choices for the separator variable
  - Both cases use self-join-freeness!

How does the proof change?

- **Upper bound**: we can generalize the algorithm
  - **Independent AND** of connected components
  - **Independent OR** of possible choices for the separator variable
  - **Both cases use self-join-freeness!**

- **Lower bound**: a non-hierarchical expression contains a pattern like

  $x \longrightarrow y \longrightarrow z \longrightarrow w$

## Extending beyond arity-two (3)

How does the proof change?

- **Upper bound**: we can generalize the algorithm
  - **Independent AND** of connected components
  - **Independent OR** of possible choices for the separator variable
  - **Both cases use self-join-freeness!**

- **Lower bound**: a non-hierarchical expression contains a pattern like

  $x \longrightarrow y \longrightarrow z \longrightarrow w$

Via **equivalent characterization**: a non-hierarchical query has two variables *x* and *y* and:

- One atom containing *x* and *y*

- One atom containing *x* but not *y*

- One atom containing *y* but not *x*

# The "big" Dalvi and Suciu dichotomy

Full dichotomy on the **unions of conjunctive queries** (UCQs):

### Theorem ([Dalvi and Suciu, 2012])

*Let **Q** be a UCQ:*

- *If **Q** is handled by a complicated algorithm then* $\mathrm{PQE}(Q)$ *is in **PTIME***
- *Otherwise,* $\mathrm{PQE}(Q)$ *is **#P-hard***

# The "big" Dalvi and Suciu dichotomy

Full dichotomy on the **unions of conjunctive queries** (UCQs):

> **Theorem ([Dalvi and Suciu, 2012])**
>
> *Let $Q$ be a UCQ:*
>
> - *If $Q$ is handled by a complicated algorithm then* $\mathrm{PQE}(Q)$ *is in **PTIME***
> - *Otherwise,* $\mathrm{PQE}(Q)$ *is **#P-hard***

This result is **far more challenging**:

- **Upper bound:**
  - an algorithm generalizing the previous case with **inclusion-exclusion**
  - many **unpleasant details** (e.g., a ranking transformation)

# The "big" Dalvi and Suciu dichotomy

Full dichotomy on the **unions of conjunctive queries** (UCQs):

## Theorem ([Dalvi and Suciu, 2012])

*Let $Q$ be a UCQ:*

- *If $Q$ is handled by a complicated algorithm then* $\mathrm{PQE}(Q)$ *is in PTIME*
- *Otherwise,* $\mathrm{PQE}(Q)$ *is #P-hard*

This result is **far more challenging**:

- **Upper bound:**
  - an algorithm generalizing the previous case with **inclusion-exclusion**
  - many **unpleasant details** (e.g., a ranking transformation)
- **Lower bound:** hardness proof on minimal cases where the algorithm does not work (very challenging)

📄 Dalvi, N. and Suciu, D. (2007).
**The dichotomy of conjunctive queries on probabilistic structures.**
In *Proc. PODS*.

📄 Dalvi, N. and Suciu, D. (2012).
**The dichotomy of probabilistic inference for unions of conjunctive queries.**
*J. ACM*, 59(6).