

Probabilistic Query Evaluation

Tuple Independent Database (D, π) :

| | | | | |
|---|---|---|---|---|
| S | $.5 \times .2$ | $.5 \times (1 - .2)$ | $(1 - .5) \times .2$ | $(1 - .5) \times (1 - .2)$ |
| represents | S | S | S | S |
| $\begin{array}{ c c } \hline a & a \\ \hline b & c \\ \hline \end{array}$ | $\begin{array}{ c c } \hline a & a \\ \hline b & c \\ \hline \end{array}$ | $\begin{array}{ c c } \hline a & a \\ \hline b & c \\ \hline \end{array}$ | $\begin{array}{ c c } \hline b & c \\ \hline \end{array}$ | $\begin{array}{ c c } \hline b & c \\ \hline \end{array}$ |

Probabilistic Query Evaluation PQE(q):

- **INPUT:** TID (D, π)
- **OUTPUT:** probability that (D, π) satisfies q

$$\Pr((D, \pi) \models \exists x, y S(x, y) \wedge x \neq y) = .5 \times .2 + (1 - .5) \times .2$$

For **Unions of Conjunctive Queries**, [Dalvi & Suciu] imply:

- There is a class of **safe** UCQs such that
 - PQE(q) is **PTIME** if q is **safe**
 - PQE(q) is **#P-hard** if q is **not safe**

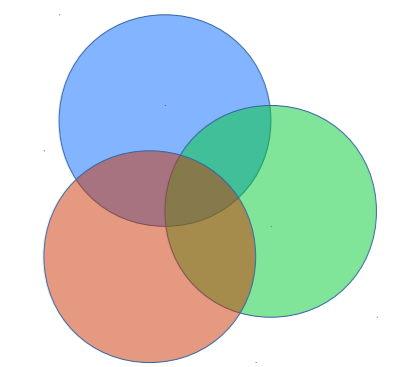
Two Techniques to Compute the Probability

Independence:

- $\Pr(q_1 \wedge q_2) = \Pr(q_1) \times \Pr(q_2)$
 - $\Pr(q_1 \vee q_2) = 1 - (1 - \Pr(q_1)) \times (1 - \Pr(q_2))$
- When q_1 and q_2 depend on disjoint sets of tuples

Inclusion-Exclusion:

- $\Pr(q_1 \wedge q_2 \wedge q_3) = \Pr(q_1) + \Pr(q_2) - \Pr(q_1 \vee q_2) - \Pr(q_2 \vee q_3) - \Pr(q_1 \vee q_3) + \Pr(q_1 \vee q_2 \vee q_3)$



Works for any safe query [Dalvi & Suciu]

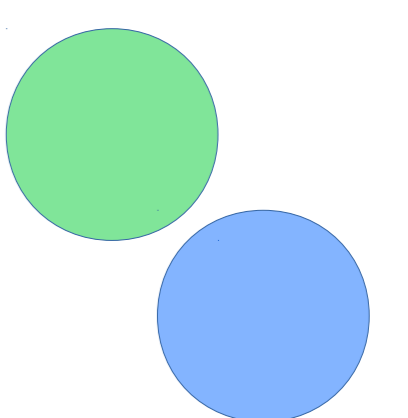
Independence:

- $\Pr(q_1 \wedge q_2) = \Pr(q_1) \times \Pr(q_2)$
- $\Pr(q_1 \vee q_2) = 1 - (1 - \Pr(q_1)) \times (1 - \Pr(q_2))$

Determinism:

- $\Pr(q_1 \vee q_2) = \Pr(q_1) + \Pr(q_2)$

When q_1 and q_2 are mutually exclusive



Question: does this work for any safe query?

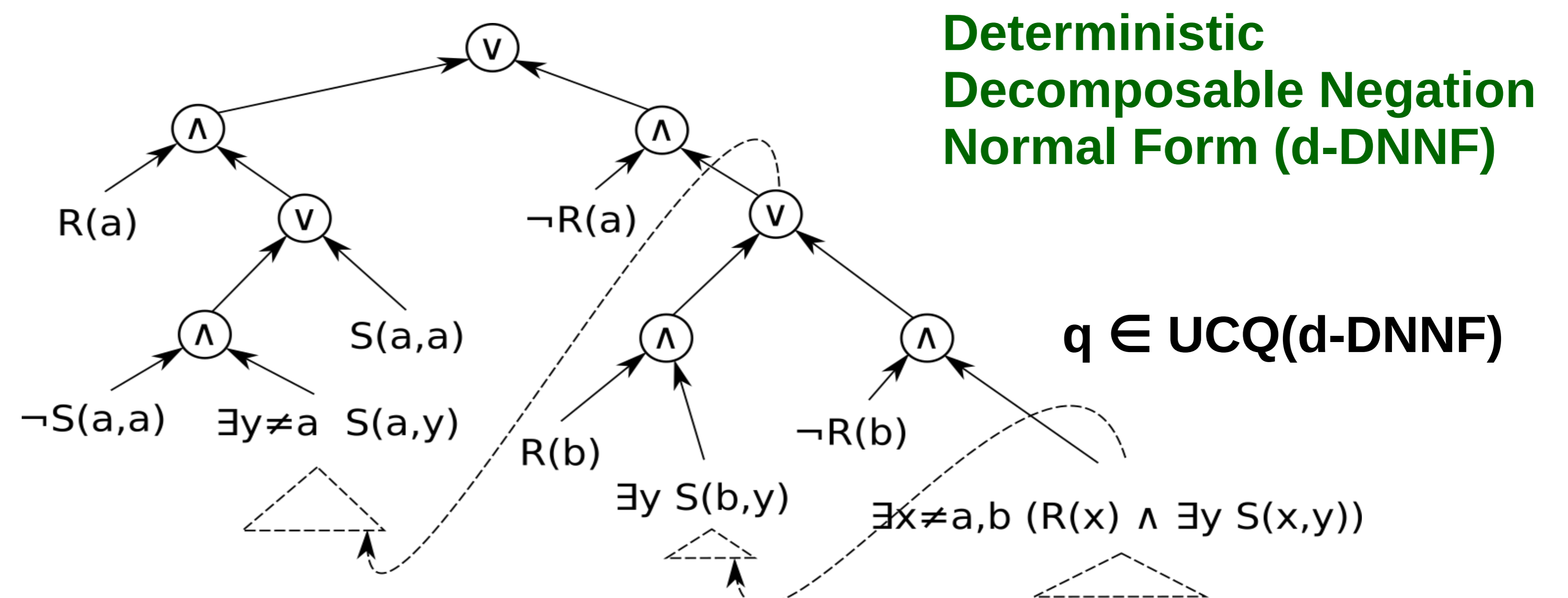
Example: $q = \exists x, y R(x) \wedge S(x, y)$

$$q = \bigvee_{a \in \text{Dom}(D)} [R(a) \wedge \exists y S(a, y)]$$

$$\Pr(q) = 1 - \prod_{a \in \text{Dom}(D)} [1 - \Pr(R(a) \wedge \exists y S(a, y))]$$

$$\Pr(R(a) \wedge \exists y S(a, y)) = \Pr(R(a)) \times \Pr(\exists y S(a, y))$$

$$\Pr(\exists y S(a, y)) = 1 - \prod_{b \in \text{Dom}(D)} [1 - \Pr(S(a, b))]$$



Deterministic Decomposable Negation Normal Form (d-DNNF)

$q \in \text{UCQ}(d\text{-DNNF})$

Problem

$$\text{Safe UCQs} \stackrel{?}{=} \text{UCQ}(d\text{-DNNF})$$

The H_k query classes

Fix $k \geq 1$, define:

- $0 = \exists x, y R(x) \wedge S(x, y)$
 - $i = \exists x, y S_i(x) \wedge S_{i+1}(x, y)$ for $1 \leq i \leq k$
 - $k = \exists x, y S_k(x, y) \wedge T(y)$
- H_k query class:** monotone Boolean combination of the queries $\{0, 1, \dots, k\}$
- Example: take $k = 3$ and define $q' = (2 \vee 3) \wedge (1 \vee 3) \wedge (0 \vee 3) \wedge (0 \vee 1 \vee 2)$
 - $\Pr(q') = \Pr(2 \vee 3) + \Pr(1 \vee 3) + \Pr(0 \vee 3) + \Pr(0 \vee 1 \vee 2) - \Pr(2 \vee 3 \vee 1) - \Pr(2 \vee 3 \vee 0) - \Pr(2 \vee 3 \vee 0 \vee 1) - \Pr(1 \vee 3 \vee 0) - \Pr(1 \vee 3 \vee 0 \vee 2) - \Pr(0 \vee 3 \vee 1 \vee 2) + \Pr(2 \vee 3 \vee 1 \vee 0) + \Pr(2 \vee 3 \vee 0 \vee 1) + \Pr(2 \vee 3 \vee 0 \vee 1) + \Pr(1 \vee 3 \vee 0 \vee 2) - \Pr(2 \vee 3 \vee 1 \vee 0)$

If the term $0 \vee 1 \vee \dots \vee k$ cancels out, the query is **safe**, otherwise it is **not safe**

If the Boolean function representing the query does not depend on all the variables $\{0, \dots, k\}$, then the query is in **UCQ(d-DNNF)**

Our Method

We **rewrite** q' as $q_0 \vee q_1 \vee q_2 \vee q_3$, where:

- $q_0 = 0 \wedge \neg 2 \wedge 3$ which is in UCQ(d-DNNF)
- $q_1 = \neg 1 \wedge 2 \wedge 3$ which is in UCQ(d-DNNF)
- $q_2 = \neg 0 \wedge 1 \wedge 3$ which is in UCQ(d-DNNF)
- $q_3 = 0 \wedge 1 \wedge 2$ which is in UCQ(d-DNNF)

The disjunctions are deterministic, hence q' is in UCQ(d-DNNF)

We say that q' is **nice**

Results

- We **generated** all the non-trivial safe queries in H_k for $1 \leq k \leq 6$ (about 20M queries)
- We determined if they are **nice** using a **sat-solver**
- Turns out they are almost all nice, so in **UCQ(d-DNNF)**!
- The 69 queries that are not nice are **co-nice** (their negation is nice), so they are in **UCQ(d-D)**
- This is encouraging experiments in favor of **Safe UCQs = UCQ(d-D)**
- Could the queries that are not nice but co-nice separate UCQ(d-DNNF) from UCQ(d-D)?