

The Main Result

Data complexity of any $1RA^-$ query Q on tuple-independent databases: Polynomial time if Q is **hierarchical** and #P-hard otherwise.

Query Language and Data Model

Relational algebra query language fragment $1RA^-$

- Included: Equi-joins, selections, projections, **difference**
- Excluded: Repeating relation symbols (self-joins), unions

Tuple-independent probabilistic model

- Each tuple associated with a fresh Boolean random variable x .
- $P(x)$ is the probability that the tuple exists in the database.
- Simplest probabilistic model in the literature. Beyond this model, query tractability is quickly lost.
- Used by real-world large-scale probabilistic repositories, e.g., Google Knowledge Vault.

The Hierarchical Property for a Query Q

For every pair of distinct attribute equivalence classes $[A]$, $[B]$ there is no triple of relation symbols R , S , and T in Q such that

- $R^{[A][\neg B]}$ has attributes in $[A]$ and not in $[B]$,
- $S^{[A][B]}$ has attributes in both $[A]$ and $[B]$, and
- $T^{[\neg A][B]}$ has attributes in $[B]$ and not in $[A]$.

Non-hierarchical queries	Hierarchical queries
$\pi_{\emptyset}[R(A) \bowtie S(A, B) \bowtie T(B)]$	$\pi_{\emptyset}[(R(A) \bowtie S(A, B)) - T(A, B)]$
$\pi_{\emptyset}[\pi_B(R(A) \bowtie S(A, B)) - T(B)]$	$\pi_{\emptyset}[(R(A) \times T(B)) - (U(A) \times V(B))]$
$\pi_{\emptyset}[T(B) - \pi_B(R(A) \bowtie S(A, B))]$	$\pi_{\emptyset}[(M(A) \times N(B)) - [(R(A) \times T(B)) - (U(A) \times V(B))]]$

The hierarchical property can be recognized in LOGSPACE.

The Hard Queries

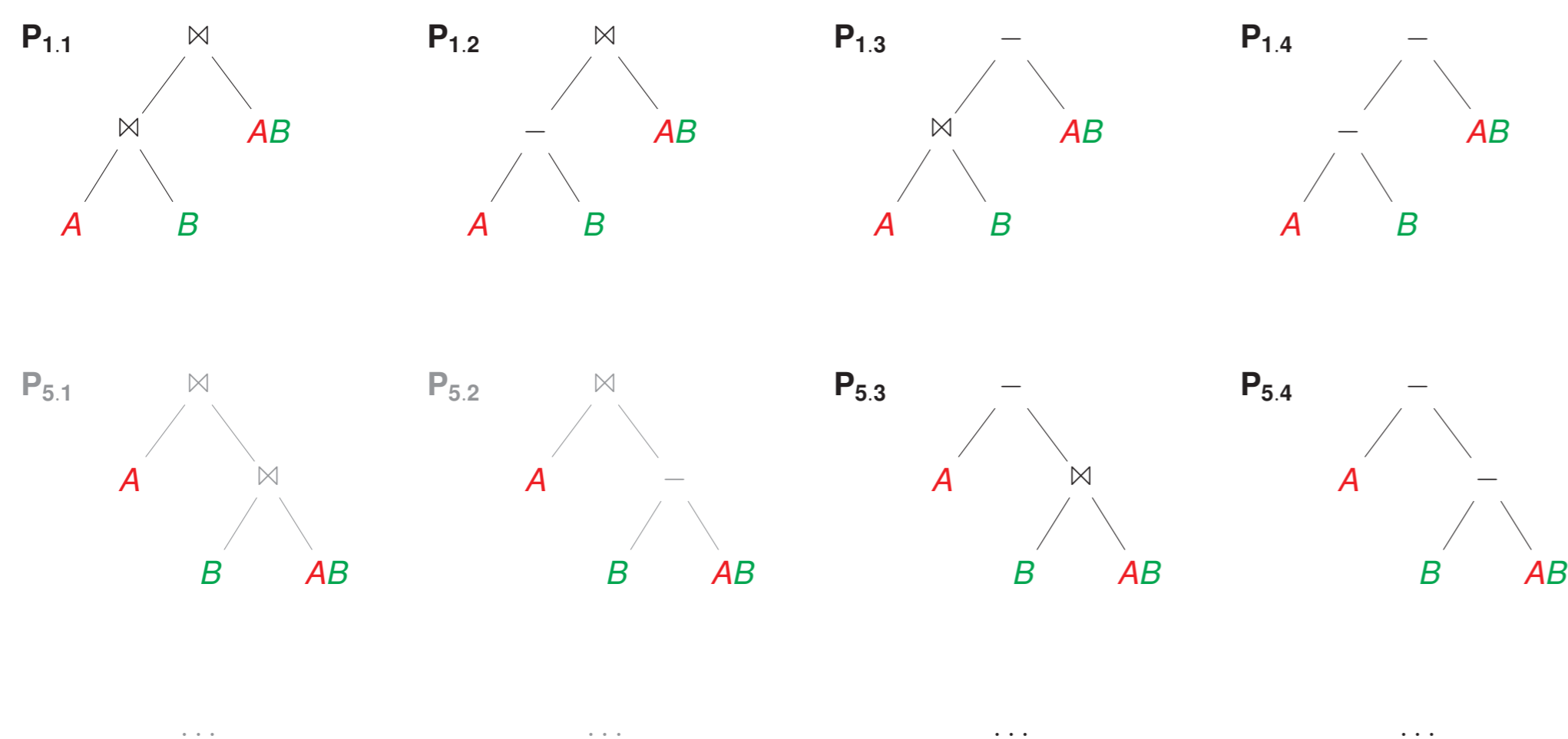
Reduction from the #P-hard problem #SAT for positive 2DNF.

- Input formula and query: $\Psi = x_1 y_1 \vee x_1 y_2$, $Q = \pi_{\emptyset}[R(A) - \pi_A(T(B) \bowtie S(A, B))]$
- Construct database such that Ψ annotates Q 's result:
 - $S(a, b, \phi)$: Clause a has variable b exactly when ϕ is true.
 - $R(a, \top)$ and $T(b, \neg b)$: a is a clause and b is a variable in Ψ .

R	T	S	$T \bowtie S$	$\pi_A(T \bowtie S)$	$R - \pi_A(T \bowtie S)$
$A \ \phi$	$B \ \phi$	$A \ B \ \phi$	$A \ B \ \phi$	$A \ \phi$	$A \ \phi$
1 \top	$x_1 \neg x_1$	1 $x_1 \top$	1 $x_1 \neg x_1$	1 $\neg x_1 \vee \neg y_1$	1 $x_1 y_1$
2 \top	$y_1 \neg y_1$	1 $y_1 \top$	1 $y_1 \neg y_1$	2 $\neg x_1 \vee \neg y_2$	2 $x_1 y_2$
	$y_2 \neg y_2$	1 $y_2 \perp$	1 $y_2 \perp$		
		2 $x_1 \top$	2 $x_1 \neg x_1$		
		2 $y_1 \perp$	2 $y_1 \perp$		
		2 $y_2 \top$	2 $y_2 \neg y_2$		

There are 48 (!) minimal non-hierarchical query patterns.

- Binary trees with leaves A , AB , and B and inner nodes \bowtie or $-$.



- There is a database construction scheme for each pattern.

Each non-hierarchical query Q matches a pattern $P_{x,y}$:

- There is a total mapping from $P_{x,y}$ to Q 's parse tree that
 - is identity on inner nodes \bowtie and $-$,
 - preserves ancestor-descendant relationships,
 - maps leaves A , B , AB to relations $R^{[A][\neg B]}$, $S^{[A][B]}$, $T^{[\neg A][B]}$.
- The match preserves the annotation of the query pattern: Q and $P_{x,y}$ have the same annotation for any input database.

The Evaluation Algorithm for Hierarchical Queries

- For any database D , the probability $P_{Q(D)}$ of a $1RA^-$ query Q is the probability P_{Ψ} of the query annotation Ψ .

$$Q = \pi_{\emptyset}(R(A) \times T(B)) - (U(A) \times V(B))$$

R	T	U	V	$R \bowtie T$	$R \bowtie T - U \bowtie V$
$A \ \phi$	$B \ \phi$	$A \ \phi$	$B \ \phi$	$A \ B \ \phi$	$A \ B \ \phi$
1 r_1	1 t_1	1 u_1	1 v_1	1 1 $r_1 t_1$	1 1 $r_1 t_1 \neg(u_1 v_1)$
2 r_2	2 t_2	2 u_2	2 v_2	1 2 $r_1 t_2$	1 2 $r_1 t_2 \neg(u_1 v_2)$
				2 1 $r_2 t_1$	2 1 $r_2 t_1 \neg(u_2 v_1)$
				2 2 $r_2 t_2$	2 2 $r_2 t_2 \neg(u_2 v_2)$

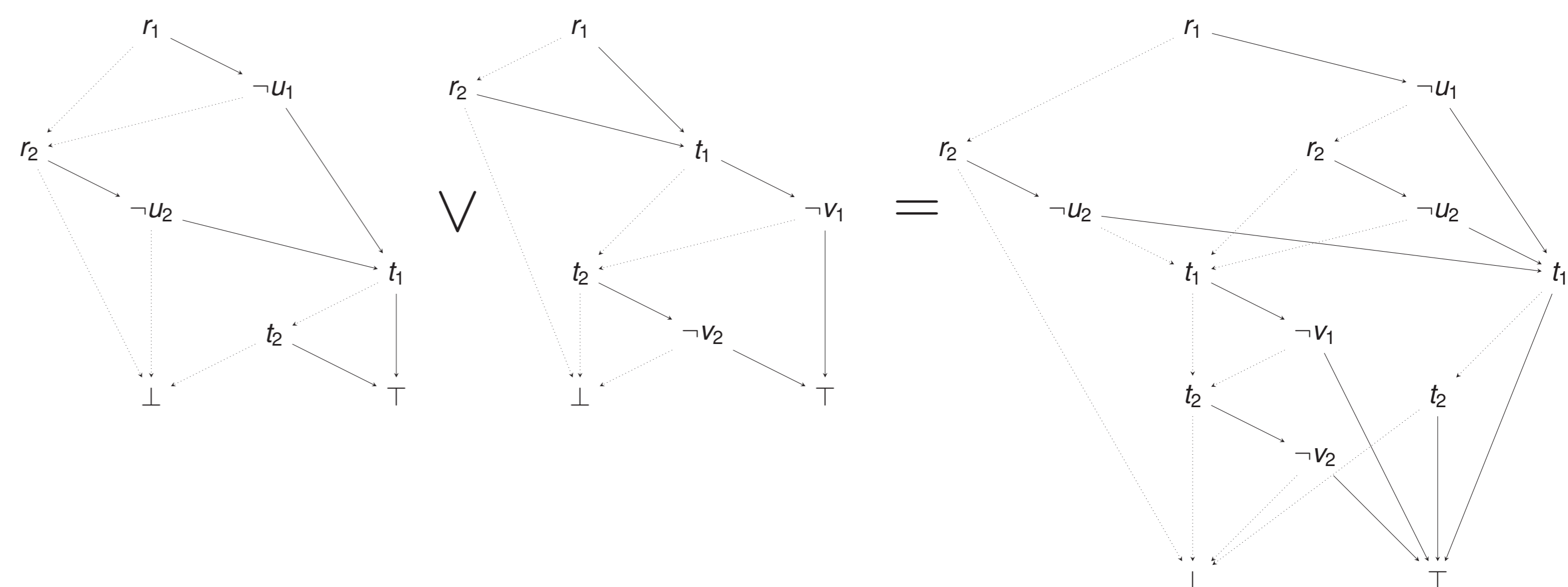
- Translate query Q into equivalent RC^{\exists} : A disjunction of disjunction-free existential relational calculus queries.

$$Q_{RC} = \underbrace{\exists_A(R(A) \wedge \neg U(A)) \wedge \exists_B T(B)}_{Q_1} \vee \underbrace{\exists_A R(A) \wedge \exists_B (T(B) \wedge \neg V(B))}_{Q_2}$$

- **RC-hierarchical**: For each quantifier $\exists_X(Q)$, every relation symbol in Q has variable X .
- \exists -consistent: All disjuncts have the same nesting order of \exists s.
- Compile query annotation into OBDD

$$\Psi = \underbrace{(r_1 \neg u_1 \vee r_2 \neg u_2) \wedge (t_1 \vee t_2)}_{\Psi_1} \vee \underbrace{(r_1 \vee r_2) \wedge (t_1 \neg v_1 \vee t_2 \neg v_2)}_{\Psi_2}$$

- **RC-hierarchical**: Each disjunct gives rise to a poly-size OBDD.
- \exists -consistent: All OBDDs have compatible variable orders.
- The OBDD width grows exponentially with the number of disjuncts, while its height stays linear in the database size.



Dichotomies Beyond $1RA^-$

Some known dichotomies

- Conjunctive queries w/o self-joins, unions of conjunctive queries [Dalvi & Suciu 2004-2010], quantified queries [F.&O.& Rath 2011]

Full relational algebra

- seems unattainable since it is undecidable whether the union of two equivalent queries, one hard and one tractable, is tractable.

Non-repeating relational algebra = $1RA^-$ + union.

- Hierarchical property not enough.
- $\pi_{\emptyset}[(R(A) \bowtie S_1(A, B) \cup T(B) \bowtie S_2(A, B)) - S(A, B)]$ is hard, though it is equivalent to a union of two hierarchical $1RA^-$ queries.

Non-repeating relational calculus

- $S(x, y) \wedge \neg R(x)$ is tractable, $S(x, y) \wedge (R(x) \vee T(y))$ is hard. Both are non-repeatable, yet not expressible in $1RA^-$.
- Possible (though expensive) approach: Translate to RC^{\exists} and check RC-hierarchical and \exists -consistency.