# Efficient Query Computing for Uncertain Possibilistic Databases with Provenance

Angelos Vasilakopoulos     Verena Kantere

National Technical University of Athens

Cyprus University of Technology

June 20, 2011

# The problem we investigate:

How to compute answers to queries for uncertain data with attached "confidence values"?

## Problem 1

Existing models for uncertain data (e.g., possibility theory) are not closed for SPJ queries.

## Problem 2

Computing probabilities of SPJ answers in models that combine uncertainty and provenance is a problem with #P complexity.

# We propose:
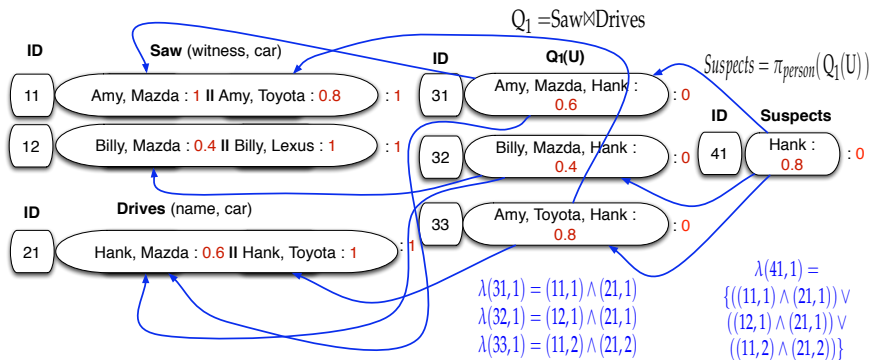
A data model that combines uncertainty, provenance and possibilities.

**Benefits of the proposed model:**

- Closed for SPJ queries
- Computing possibilities of SPJ answers has PTIME complexity.

# Motivating Example



$Q_1 = Saw \bowtie Drives$

**ID**     **Saw** (witness, car)

| 11 | Amy, Mazda : 1 **II** Amy, Toyota : 0.8 | : 1 |

| 12 | Billy, Mazda : 0.4 **II** Billy, Lexus : 1 | : 1 |

**ID**     **Drives** (name, car)

| 21 | Hank, Mazda : 0.6 **II** Hank, Toyota : 1 | : 1 |

**ID**     **Q₁(U)**

| 31 | Amy, Mazda, Hank : 0.6 | : 0 |
| 32 | Billy, Mazda, Hank : 0.4 | : 0 |
| 33 | Amy, Toyota, Hank : 0.8 | : 0 |

$Suspects = \pi_{person}\big(Q_1(U)\big)$

**ID**     **Suspects**

| 41 | Hank : 0.8 | : 0 |

$\lambda(31,1) = (11,1) \wedge (21,1)$
$\lambda(32,1) = (12,1) \wedge (21,1)$
$\lambda(33,1) = (11,2) \wedge (21,2)$

$\lambda(41,1) =$
$\{((11,1) \wedge (21,1)) \vee$
$((12,1) \wedge (21,1)) \vee$
$((11,2) \wedge (21,2))\}$

# Possible Worlds:

# Possibilistic Databases are not closed



| **ID** | $Q_1 =$ Saw⋈Drives | |
|---|---|---|
| 31 | Amy, Mazda, Hank : 0.6 | : 0 |
| 32 | Billy, Mazda, Hank : 0.4 | : 0 |
| 33 | Amy, Toyota, Hank : 0.8 | : 0 |

- Tuples 31 and 33 cannot coexist in any possible world.
- Adding provenance (lineage) makes possibilistic uncertain data model closed for SPJ queries.

A. Das Sarma et al. "Representing Uncertain Data: Models, Properties, and Algorithms". VLDB Journal, October, 2009.
P. Bosc and O. Pivert. "About PSJ queries addressed to possibilistic relational databases". IEEE T. Fuzzy Systems, 2005.

# Computing probabilities is intractable

$$Suspects = \pi_{person}\big(Q_1(U)\big)$$

| ID | Suspects |
|----|----------|
| 41 | Hank |

$$\lambda(41,1) =$$
$$\{((11,1) \wedge (21,1)) \vee$$
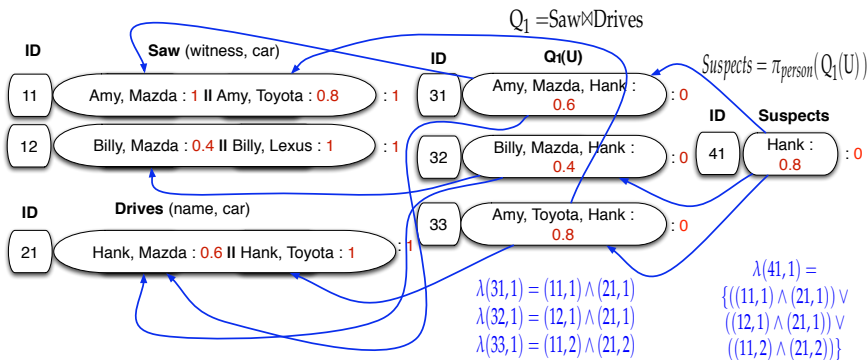$$((12,1) \wedge (21,1)) \vee$$
$$((11,2) \wedge (21,2))\}$$

: Computing probabilities has #P complexity

N. Dalvi and D. Suciu. "Efficient query evaluation on probabilistic databases." In VLDB 2004.

## Axioms of Possibilistic theory:

- $\Pi(X) = 1$
- $\Pi(\emptyset) = 0$
- $\Pi(E_1 \cup E_2) = \max(\Pi(E_1), \Pi(E_2))$
- $\Pi(E_1 \cap E_2) \leq \min(\Pi(E_1), \Pi(E_2))$
- $\Pi(E_1 \cap E_2) = \min(\Pi(E_1), \Pi(E_2))$
  (for not-interactive events)
- $\max\{\Pi(E), \Pi(\bar{E})\} = 1$
- $N(E) = 1 - \Pi(\bar{E})$

# Running Example



$Q_1 =$ Saw$\bowtie$Drives

$Suspects = \pi_{person}(Q_1(U))$

**ID** — **Saw** (witness, car)

| 11 | Amy, Mazda : 1 **II** Amy, Toyota : 0.8 | : 1 |
| 12 | Billy, Mazda : 0.4 **II** Billy, Lexus : 1 | : 1 |

**ID** — **Drives** (name, car)

| 21 | Hank, Mazda : 0.6 **II** Hank, Toyota : 1 | : 1 |

**ID** — $Q_1(U)$

| 31 | Amy, Mazda, Hank : 0.6 | : 0 |
| 32 | Billy, Mazda, Hank : 0.4 | : 0 |
| 33 | Amy, Toyota, Hank : 0.8 | : 0 |

**ID** — **Suspects**

| 41 | Hank : 0.8 | : 0 |

$\lambda(31,1) = (11,1) \wedge (21,1)$
$\lambda(32,1) = (12,1) \wedge (21,1)$
$\lambda(33,1) = (11,2) \wedge (21,2)$

$\lambda(41,1) =$
$\{((11,1) \wedge (21,1)) \vee$
$((12,1) \wedge (21,1)) \vee$
$((11,2) \wedge (21,2))\}$

# Possible Worlds:

# Ongoing Work:

- Extend the query language to extentions of conjunctive queries.
- Find for which class of query languages the problem remains in polynomial time.
- Find for which class of query languages the problem becomes intractable.
- Study complexity of new query languages that can query over uncertainty and provenance.

**Thank you**