

Leveraging the Structure of Uncertain Data

Tirer parti de la structure des données incertaines

Antoine Amarilli

Télécom ParisTech, DBWeb

March 14th, 2016



Databases

Computers often use **databases** to **store** data and **query** it

Databases

Computers often use **databases** to **store** data and **query** it

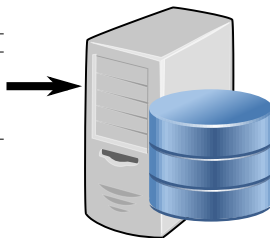


Databases

Computers often use **databases** to **store** data and **query** it

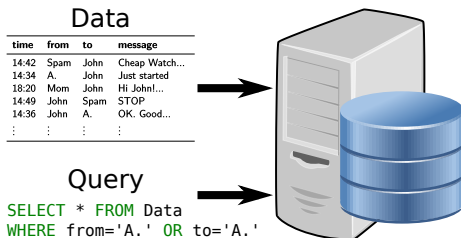
Data

time	from	to	message
14:42	Spam	John	Cheap Watch...
14:34	A.	John	Just started
18:20	Mom	John	Hi John!...
14:49	John	Spam	STOP
14:36	John	A.	OK, Good...
⋮	⋮	⋮	⋮



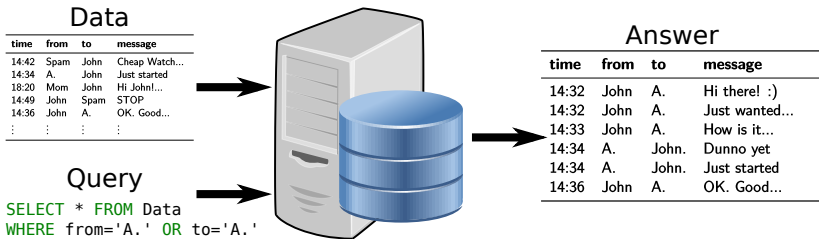
Databases

Computers often use **databases** to **store** data and **query** it



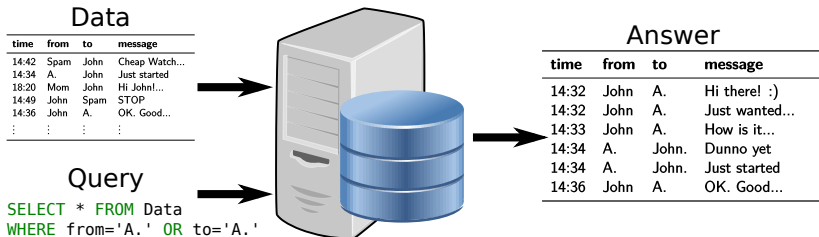
Databases

Computers often use **databases** to **store** data and **query** it



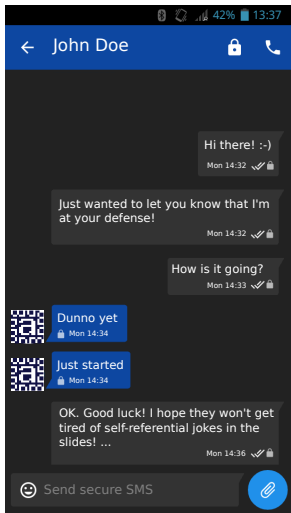
Databases

Computers often use **databases** to **store** data and **query** it

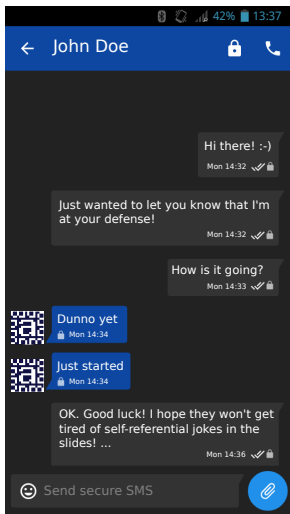


→ Let's see a few **examples**...

Database example: SMS on Android



Database example: SMS on Android



time	from	to	message
14:32	John	A.	Hi there! :-)
14:32	John	A.	Just wanted...
14:33	John	A.	How is it...
14:34	A.	John	Dunno yet
14:34	A.	John	Just started
14:36	John	A.	OK. Good...

In reality...

```
CREATE TABLE sms (_id INTEGER, thread_id INTEGER,
  address TEXT, address_device_id INTEGER, person INTEGER,
  date INTEGER, date_sent INTEGER, protocol INTEGER,
  read INTEGER, status INTEGER, type INTEGER,
  reply_path_present INTEGER,
  delivery_receipt_count INTEGER, subject TEXT, body TEXT,
  mismatched_identities TEXT, service_center TEXT,
  date_delivery_received INTEGER);
```

In reality...

```
CREATE TABLE sms (_id INTEGER, thread_id INTEGER,
  address TEXT, address_device_id INTEGER, person INTEGER,
  date INTEGER, date_sent INTEGER, protocol INTEGER,
  read INTEGER, status INTEGER, type INTEGER,
  reply_path_present INTEGER,
  delivery_receipt_count INTEGER, subject TEXT, body TEXT,
  mismatched_identities TEXT, service_center TEXT,
  date_delivery_received INTEGER);
```

```
INSERT INTO sms VALUES(
  14041,224,'+33611210549',1,NULL,1451921855098,
  1451921849000,0,1,-1,-2147483628,0,0,NULL,
  'Hi there!',NULL,'+33609002960',0);
```

```
INSERT INTO sms VALUES(
  14042,224,'+33611210549',1,NULL,1451921945081,
  1451921945081,NULL,1,-1,-2147483561,NULL,0,NULL,
  'Just wanted...',NULL,NULL,0);
```

Database example: Wikipedia

Recent changes

- [Naza](#); 14:48 . . (-59) . . [98.115.58.241](#)
- [HK Olimpija Ljubljana \(2004\)](#); 14:48 . . (+4) . . [86.58.36.235](#)
- [Monster High](#); 14:48 . . (+18) . . [66.244.123.117](#)
- [List of songs recorded by Celine Dion](#); 14:48 . . (+25) . . [79.94.26.185](#)
- [Biodegradable waste](#); 14:48 . . (+5) . . [59.90.26.215](#)

Database example: Wikipedia

Recent changes

- [Naza](#); 14:48 .. (-59) .. [98.115.58.241](#)
- [HK Olimpija Ljubljana \(2004\)](#); 14:48 .. (+4) .. [86.58.36.235](#)
- [Monster High](#); 14:48 .. (+18) .. [66.244.123.117](#)
- [List of songs recorded by Celine Dion](#); 14:48 .. (+25) .. [79.94.26.185](#)
- [Biodegradable waste](#); 14:48 .. (+5) .. [59.90.26.215](#)

title	time	size	user
Naza	14:48	-59	92.115.58.241
HK Olimpija Ljubljana (2004)	14:48	+4	86.58.36.235
Monster High	14:48	+18	66.244.123.117
List of songs recorded by Celine Dion	14:48	+25	79.94.26.185
Biodegradable waste	14:48	+5	59.90.26.215

In reality...

```
CREATE TABLE mw_recentchanges (rc_id INT(8),
  rc_timestamp VARCHAR(14), rc_cur_time VARCHAR(14),
  rc_user INT(10), rc_user_text VARCHAR(255),
  rc_namespace INT(11), rc_title VARCHAR(255),
  rc_comment VARCHAR(255), rc_minor TINYINT(3),
  rc_bot TINYINT(3), rc_new TINYINT(3),
  rc_cur_id INT(10), rc_this_oldid INT(10),
  rc_last_oldid INT(10), rc_type TINYINT(3),
  rc_moved_to_ns TINYINT(3), rc_moved_to_title VARCHAR(255),
  rc_patrolled TINYINT(3), rc_ip CHAR(15),
  rc_old_len INT(10), rc_new_len INT(10),
  rc_deleted TINYINT(1), rc_logid INT(10),
  rc_log_type VARCHAR(255), rc_log_action VARCHAR(255),
  rc_params BLOB,
);
```

In reality...

```
CREATE TABLE mw_recentchanges (rc_id INT(8),
  rc_timestamp VARCHAR(14), rc_cur_time VARCHAR(14),
  rc_user INT(10), rc_user_text VARCHAR(255),
  rc_namespace INT(11), rc_title VARCHAR(255),
  rc_comment VARCHAR(255), rc_minor TINYINT(3),
  rc_bot TINYINT(3), rc_new TINYINT(3),
  rc_cur_id INT(10), rc_this_oldid INT(10),
  rc_last_oldid INT(10), rc_type TINYINT(3),
  rc_moved_to_ns TINYINT(3), rc_moved_to_title VARCHAR(255),
  rc_patrolled TINYINT(3), rc_ip CHAR(15),
  rc_old_len INT(10), rc_new_len INT(10),
  rc_deleted TINYINT(1), rc_logid INT(10),
  rc_log_type VARCHAR(255), rc_log_action VARCHAR(255),
  rc_params BLOB,
);

INSERT INTO mw_recentchanges VALUES
  (1, '20160314144837', '20160314144827', 1, '92.115.58.241', 0,
  'Naza', '', 0, 0, 0, 1, 2, 1, 0, 0, '', 1, '92.115.58.241',
  559, 500, 0, 0, NULL, NULL, '');

INSERT INTO mw_recentchanges VALUES
  (2, '20160314144842', '20160314144842', 1, '66.244.123.117', 2,
  'Monster High', '', 0, 0, 1, 2, 3, 0, 1, 0, '', 1, '66.244.123.117',
  102, 120, 0, 0, NULL, NULL, '');
```

Uncertainty



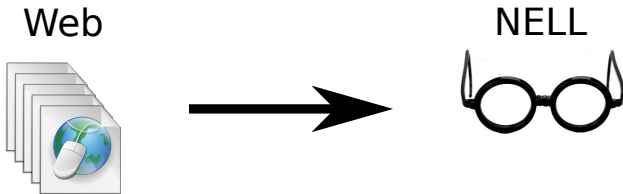
- Databases usually assume that data is
 - complete
 - crisp
 - certain
 - correct
- In many situations, this is not the case...

Example: Never-Ending Language Learning

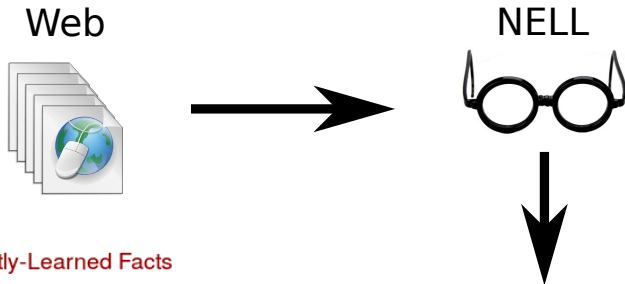
Web



Example: Never-Ending Language Learning



Example: Never-Ending Language Learning



Recently-Learned Facts

instance	iteration	date learned	confidence
kampioenschap van zwitserland is a sports race	955	20-oct-2015	95.0
cochran mill nature center is an aquarium	955	20-oct-2015	96.9
kozy shack chocolate pudding is a kind of candy	956	23-oct-2015	90.3
red delicious apple tree is a plant	955	20-oct-2015	92.8
sale miami dade county is a sport	955	20-oct-2015	99.1
chicken001 eat black beans	955	20-oct-2015	100.0
wrigley field is the home venue for the sports team chicago cubs	959	07-nov-2015	100.0
lorena ochoa is a person who has residence in the geopolitical location mexico	958	03-nov-2015	100.0
umass lowell river hawks hired john calipari	955	20-oct-2015	98.4
nuggets participated in the event games	955	20-oct-2015	100.0

Many sources of uncertainty

- Errors in sources:



This article's **factual accuracy is disputed**. Please help to ensure that disputed statements are **reliably sourced**. See the relevant discussion on the [talk page](#). *(November 2015)*

Many sources of uncertainty

- Errors in sources:



This article's **factual accuracy is disputed**. Please help to ensure that disputed statements are **reliably sourced**. See the relevant discussion on the **talk page**. *(November 2015)*

- Entity disambiguation:

*“The place and **function of Venus** in Ovid...”*

*“Computed backscattering **function of Venus** and the moon...”*

Many sources of uncertainty

- Errors in sources:



This article's **factual accuracy is disputed**. Please help to ensure that disputed statements are **reliably sourced**. See the relevant discussion on the **talk page**. *(November 2015)*

- Entity disambiguation:

*“The place and **function of Venus** in Ovid...”*

*“Computed backscattering **function of Venus** and the moon...”*

- Anaphora resolution:

*“Obama told Hollande that **he** was not a spying target”*

Many uncertain data applications

- Information extraction
- Machine learning
- Speech recognition
- Data integration
- Crowdsourcing
- ...
- PhD defense scheduling

Journal Articles

OCR → ... The Namurian Tsingyuan Formation from Ningxia, China, is divided into three members ...

NLP → The Namurian Tsingyuan Formation from Ningxia

Relational Features

Entity1	Entity2	Feature
Namurian	Tsingyuan Fm.	nn
Silesian	Tsingyuan	SameRow

SQL+Python → Existing Tools

(a) Possible Mapping

Mapping	Prob
$m_1 = \{(pname, name), (email-addr, email), (current-addr, mailing-addr), (permanent-addr, home-addr)\}$	0.5
$m_2 = \{(pname, name), (email-addr, email), (permanent-addr, mailing-addr), (current-addr, home-addr)\}$	0.4
$m_3 = \{(pname, name), (email-addr, mailing-addr), (current-addr, home-addr)\}$	0.1

(b)

pname	email-addr	current-addr	permanent-addr
Alice	alice@0	Mountain View	Sunnyvale
Bob	bob@0	Sunnyvale	Sunnyvale

(c)

Tuple (mailing-addr)	Prob
("Sunnyvale")	0.9
("Mountain View")	0.5
("alice@0")	0.1
("bob@0")	0.1

(d)

$Z \rightarrow \mu^+ \mu^- + jets$

● Data
● Sherpa
● Madgraph

ATLAS Preliminary
13 TeV, 85 pb⁻¹

anti- k_t , $R=0.4$
 $p_t^j > 30$ GeV
 $|y^j| < 2.5$

$\sigma(Z \rightarrow \mu^+ \mu^- + jets) / (N_{jets} + 1) N_{jets}$

Uncertainty applied to PhD defenses

Who will attend this PhD defense?

Statistics

Number of people invited

Uncertainty applied to PhD defenses

Who will attend this PhD defense?

Statistics

Number of people invited 79

Uncertainty applied to PhD defenses

Who will attend this PhD defense?

Statistics

Number of people invited 79

Number of definite **yes** answers

Uncertainty applied to PhD defenses

Who will attend this PhD defense?

Statistics

Number of people invited	79
Number of definite <i>yes</i> answers	37

Uncertainty applied to PhD defenses

Who will attend this PhD defense?

Statistics

Number of people invited	79
Number of definite yes answers	37
Number of definite no answers	

Uncertainty applied to PhD defenses

Who will attend this PhD defense?

Statistics

Number of people invited	79
Number of definite yes answers	37
Number of definite no answers	13

Uncertainty applied to PhD defenses

Who will attend this PhD defense?

Statistics

Number of people invited	79
Number of definite yes answers	37
Number of definite no answers	13
Number of uncertain answers	

Uncertainty applied to PhD defenses

Who will attend this PhD defense?

Statistics

Number of people invited	79
Number of definite yes answers	37
Number of definite no answers	13
Number of uncertain answers	29

Uncertainty applied to PhD defenses

Who will attend this PhD defense?

Statistics

Number of people invited	79
Number of definite yes answers	37
Number of definite no answers	13
Number of uncertain answers	29
Number of additional people showing up	

Uncertainty applied to PhD defenses

Who will attend this PhD defense?

Statistics

Number of people invited	79
Number of definite yes answers	37
Number of definite no answers	13
Number of uncertain answers	29
Number of additional people showing up	??

Uncertainty applied to PhD defenses

Who will attend this PhD defense?

Statistics

Number of people invited	79
Number of definite yes answers	37
Number of definite no answers	13
Number of uncertain answers	29
Number of additional people showing up	??

Why is uncertainty challenging?

- Data is **uncertain** if we don't know its exact state
- A **possible world** is an actual outcome

Why is uncertainty challenging?

- Data is **uncertain** if we don't know its exact state
- A **possible world** is an actual outcome
- **Simplest method**: write out all possible worlds

Why is uncertainty challenging?

- Data is **uncertain** if we don't know its exact state
- A **possible world** is an actual outcome
- **Simplest method**: write out all possible worlds

List of the people
who **may** show up:

- Flo
- Guy
- Tat
- ...
- more?

Why is uncertainty challenging?

- Data is **uncertain** if we don't know its exact state
- A **possible world** is an actual outcome
- **Simplest method**: write out all possible worlds

List of the people
who **may** show up:

- Flo → 29 uncertain people
- Guy
- Tat
- ...
- more?

Why is uncertainty challenging?

- Data is **uncertain** if we don't know its exact state
- A **possible world** is an actual outcome
- **Simplest method**: write out all possible worlds

List of the people
who **may** show up:

- Flo → **29** uncertain people
- Guy → **536 870 912** possibilities
- Tat
- ...
- more?

Why is uncertainty challenging?

- Data is **uncertain** if we don't know its exact state
- A **possible world** is an actual outcome
- **Simplest method**: write out all possible worlds

List of the people
who **may** show up:

- Flo → 29 uncertain people
- Guy → 536 870 912 possibilities
- Tat → If the list of people is **incomplete**,
infinitely many possible completions
- ...
- more?

Uncertainty representation and semantics

Uncertain databases represent **implicitly** the possible worlds

Uncertainty representation and semantics

Uncertain databases represent **implicitly** the possible worlds

→ Probabilities

Flo	0.4
Guy	0.3
Tat	0.2
⋮	

Uncertainty representation and semantics

Uncertain databases represent **implicitly** the possible worlds

→ Probabilities

Flo	0.4
Guy	0.3
Tat	0.2
⋮	

→ Correlations

- Only one of Isa and Pal can come
- Mat and Val either come **together** or **not**
- Nell will probably come **if** Mike does

Uncertainty representation and semantics

Uncertain databases represent **implicitly** the possible worlds

→ Probabilities

Flo	0.4
Guy	0.3
Tat	0.2
⋮	

→ Correlations

- Only one of Isa and Pal can come
- Mat and Val either come **together** or **not**
- Nell will probably come **if** Mike does

→ Logical rules

If someone comes to the defense **then** they will also come to the drinks

Summary of uncertainty goals

- **Representing** our knowledge about the data
- **Computing** numerical probabilities
- **Reasoning** with logical constraints

Summary of uncertainty goals

- **Representing** our knowledge about the data
 - **Computing** numerical probabilities
 - **Reasoning** with logical constraints
- **End goal:** A database system with **first-class** uncertainty
- Feed uncertain data to the system
 - Get uncertain query results

Summary of uncertainty goals

- **Representing** our knowledge about the data
 - **Computing** numerical probabilities
 - **Reasoning** with logical constraints
- **End goal:** A database system with **first-class** uncertainty
- Feed uncertain data to the system
 - Get uncertain query results



Summary of uncertainty goals

- **Representing** our knowledge about the data
- **Computing** numerical probabilities
- **Reasoning** with logical constraints

→ **End goal:** A database system with **first-class** uncertainty

- Feed uncertain data to the system
- Get uncertain query results

Uncertain data

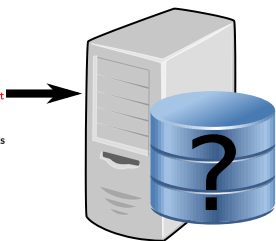
→ Probas

Flo	0.4
Guy	0.3
Tat	0.2
⋮	

→ Correlations

- Only one of Isa and Pal can come
- Mat and Val either come **together** or **not**
- Nell will probably come **if** Mike does

→ Logical rules If defense **then** drinks



Summary of uncertainty goals

- **Representing** our knowledge about the data
- **Computing** numerical probabilities
- **Reasoning** with logical constraints

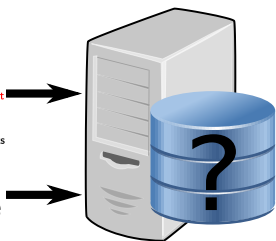
→ **End goal:** A database system with **first-class** uncertainty

- Feed uncertain data to the system
- Get uncertain query results

Uncertain data

	→ Probas	→ Correlations
Flo	0.4	● Only one of Isa and Pal can come
Guy	0.3	● Mat and Val either come together or not
Tat	0.2	● Nell will probably come if Mike does
⋮		
→ Logical rules		If defense then drinks

Query
How many people
to the drinks?



Summary of uncertainty goals

- **Representing** our knowledge about the data
- **Computing** numerical probabilities
- **Reasoning** with logical constraints

→ **End goal:** A database system with **first-class** uncertainty

- Feed uncertain data to the system
- Get uncertain query results

Uncertain data

→ Probas	→ Correlations
Flo 0.4	● Only one of Isa and Pal can come
Guy 0.3	● Mat and Val either come together or not
Tat 0.2	● Nell will probably come if Mike does
⋮	
→ Logical rules	If defense then drinks

Query

How many people to the drinks?



Uncertain answer

42 ±5 with 80% confidence

Why are uncertainty and probabilities challenging?

Uncertain attendees

Flo	0.4
Guy	0.3
Tat	0.2
Ell	0.1
⋮	

Why are uncertainty and probabilities challenging?

Uncertain attendees

Flo	0.4
Guy	0.3
Tat	0.2
Ell	0.1
⋮	

People who should meet

Flo	Guy
Ell	Tat
Ell	Guy

Why are uncertainty and probabilities challenging?

Uncertain attendees

Flo	0.4
Guy	0.3
Tat	0.2
Ell	0.1
⋮	

People who should meet

Flo	Guy
Ell	Tat
Ell	Guy

What is the probability that one of the pairs can meet?

Computing probabilities

Ell Tat
0.1 0.2



Guy Flo
0.3 0.4

Computing probabilities

Ell Tat
0.1 0.2



$$0.1 \times 0.2$$

Guy Flo
0.3 0.4

Computing probabilities

Ell Tat
0.1 0.2



$$0.1 \times 0.2 = 0.02$$

Guy Flo
0.3 0.4

Computing probabilities

Ell Tat
0.1 0.2



Guy Flo
0.3 0.4



Computing probabilities

Ell Tat
0.1 0.2



$$0.1 \times 0.2$$

Guy Flo
0.3 0.4



Computing probabilities

Ell Tat
0.1 0.2



$$0.1 \times 0.2 = 0.02$$

Guy Flo
0.3 0.4



Computing probabilities

Ell Tat
0.1 0.2

$$0.1 \times 0.2 = 0.02$$

$$0.3 \times 0.4$$

Guy Flo
0.3 0.4

Computing probabilities

Ell Tat
0.1 0.2

$$0.1 \times 0.2 = 0.02$$

$$0.3 \times 0.4 = 0.12$$

Guy Flo
0.3 0.4

Computing probabilities

Ell Tat
0.1 0.2

$$0.1 \times 0.2 = 0.02$$

$$0.3 \times 0.4 = 0.12$$

$$1 - (1 - 0.02) \times (1 - 0.12)$$

Guy Flo
0.3 0.4

Computing probabilities

Ell _____ Tat
0.1 0.2

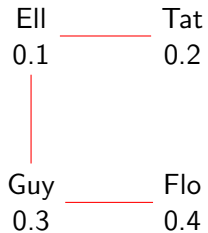
$$0.1 \times 0.2 = 0.02$$

$$0.3 \times 0.4 = 0.12$$

$$1 - (1 - 0.02) \times (1 - 0.12) = 0.1376$$

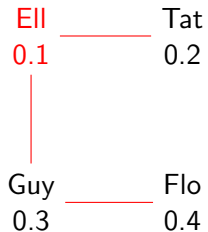
Guy _____ Flo
0.3 0.4

Computing probabilities

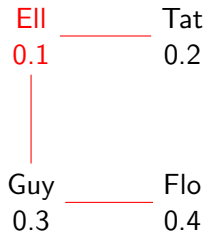


Computing probabilities

If **EII** is missing:



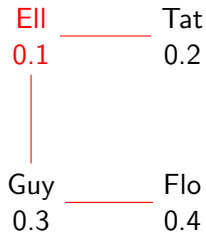
Computing probabilities



If **EII** is missing:

$$0.3 \times 0.4$$

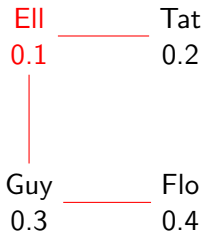
Computing probabilities



If **EII** is missing:

$$0.3 \times 0.4 = 0.12$$

Computing probabilities

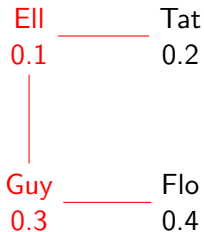


If **EII** is missing:

$$0.3 \times 0.4 = 0.12$$

If **EII** is here:

Computing probabilities



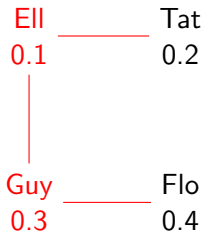
If **EII** is missing:

$$0.3 \times 0.4 = 0.12$$

If **EII** is here:

If **Guy** is missing:

Computing probabilities



If **EII** is missing:

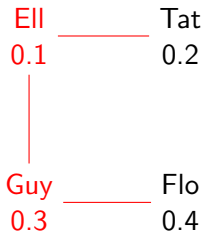
$$0.3 \times 0.4 = 0.12$$

If **EII** is here:

If **Guy** is missing:

We need Tat: 0.2

Computing probabilities



If **EII** is missing:

$$0.3 \times 0.4 = 0.12$$

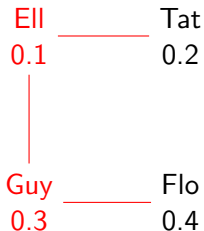
If **EII** is here:

If **Guy** is missing:

We need Tat: 0.2

If **Guy** is here:

Computing probabilities



If **EII** is missing:

$$0.3 \times 0.4 = 0.12$$

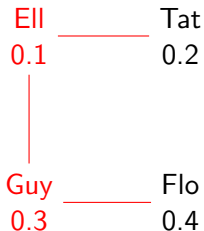
If **EII** is here:

If **Guy** is missing:

We need Tat: 0.2

If **Guy** is here: **success!**

Computing probabilities



If **EII** is missing:

$$0.3 \times 0.4 = 0.12$$

If **EII** is here:

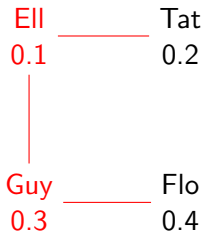
If **Guy** is missing:

We need Tat: 0.2

If **Guy** is here: **success!**

Total:

Computing probabilities



If **EII** is missing:

$$0.3 \times 0.4 = 0.12$$

If **EII** is here:

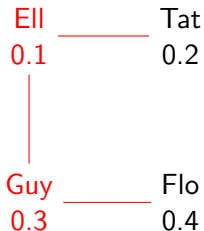
If **Guy** is missing:

We need Tat: 0.2

If **Guy** is here: **success!**

$$\text{Total: } (1 - 0.1) \times 0.12$$

Computing probabilities



If **EII** is missing:

$$0.3 \times 0.4 = 0.12$$

If **EII** is here:

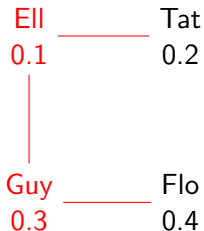
If **Guy** is missing:

We need Tat: 0.2

If **Guy** is here: **success!**

$$\text{Total: } (1 - 0.1) \times 0.12 \\ + 0.1 \times$$

Computing probabilities



If **EII** is missing:

$$0.3 \times 0.4 = 0.12$$

If **EII** is here:

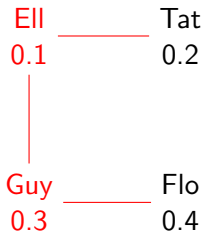
If **Guy** is missing:

We need Tat: 0.2

If **Guy** is here: **success!**

$$\begin{aligned} \text{Total: } & (1 - 0.1) \times 0.12 \\ & + 0.1 \times (0.3 + (1 - 0.3) \times 0.2) \end{aligned}$$

Computing probabilities



If **EII** is missing:

$$0.3 \times 0.4 = 0.12$$

If **EII** is here:

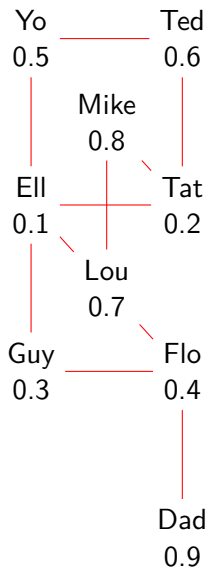
If **Guy** is missing:

We need Tat: 0.2

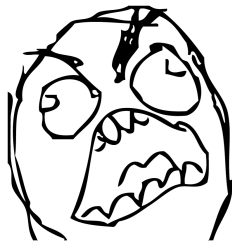
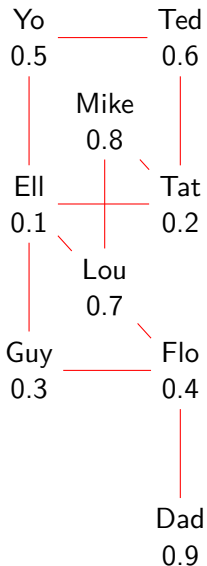
If **Guy** is here: **success!**

$$\begin{aligned} \text{Total: } & (1 - 0.1) \times 0.12 \\ & + 0.1 \times (0.3 + (1 - 0.3) \times 0.2) \\ & = 0.152 \end{aligned}$$

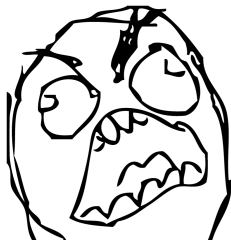
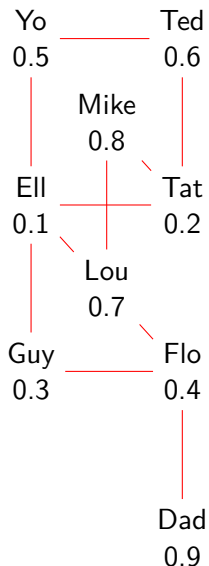
Computing probabilities



Computing probabilities



Computing probabilities



- This task is **intractable** (#P-hard)
- Many other tasks on uncertain data are intractable or even **undecidable**

My PhD topic

- Make it easier to use uncertain data
by making assumptions on the **structure** of data

My PhD topic

- Make it easier to use uncertain data
by making assumptions on the **structure** of data

0.1 ——— 0.2

0.3 ——— 0.4

0.5 0.6
| |
| |
| |
| |
0.7 0.8

My PhD topic

→ Make it easier to use uncertain data
by making assumptions on the **structure** of data

0.1 ——— 0.2

● $0.1 \times 0.2 = 0.02$

0.3 ——— 0.4

0.5
|
0.7

0.6
|
0.8

My PhD topic

→ Make it easier to use uncertain data
by making assumptions on the **structure** of data

0.1 ——— 0.2

- $0.1 \times 0.2 = 0.02$

- $0.3 \times 0.4 = 0.12$

0.3 ——— 0.4

0.5 0.6

| |

0.7 0.8

My PhD topic

→ Make it easier to use uncertain data
by making assumptions on the **structure** of data

0.1 ——— 0.2

● $0.1 \times 0.2 = 0.02$

● $0.3 \times 0.4 = 0.12$

● $0.5 \times 0.7 = 0.35$

0.3 ——— 0.4

0.5 0.6

| |
| |
| |

0.7 0.8

My PhD topic

→ Make it easier to use uncertain data
by making assumptions on the **structure** of data

0.1 ——— 0.2

● $0.1 \times 0.2 = 0.02$

● $0.3 \times 0.4 = 0.12$

● $0.5 \times 0.7 = 0.35$

0.3 ——— 0.4

● $0.6 \times 0.8 = 0.48$

0.5

0.6

0.7

0.8

My PhD topic

→ Make it easier to use uncertain data
by making assumptions on the **structure** of data

0.1 ——— 0.2

- $0.1 \times 0.2 = 0.02$

- $0.3 \times 0.4 = 0.12$

- $0.5 \times 0.7 = 0.35$

0.3 ——— 0.4

- $0.6 \times 0.8 = 0.48$

→ $1 - (1 - 0.02) \times \dots \times (1 - 0.48)$

0.5

0.6

0.7

0.8

My PhD topic

→ Make it easier to use uncertain data
by making assumptions on the **structure** of data

0.1 ——— 0.2

- $0.1 \times 0.2 = 0.02$

- $0.3 \times 0.4 = 0.12$

- $0.5 \times 0.7 = 0.35$

0.3 ——— 0.4

- $0.6 \times 0.8 = 0.48$

→ $1 - (1 - 0.02) \times \dots \times (1 - 0.48)$
 $= 0.7085088$

0.5

0.6

0.7

0.8

My PhD topic

→ Make it easier to use uncertain data
by making assumptions on the **structure** of data

0.1 ——— 0.2

- $0.1 \times 0.2 = 0.02$

- $0.3 \times 0.4 = 0.12$

- $0.5 \times 0.7 = 0.35$

0.3 ——— 0.4

- $0.6 \times 0.8 = 0.48$

→ $1 - (1 - 0.02) \times \dots \times (1 - 0.48)$
 $= 0.7085088$

0.5

0.6

0.7

0.8

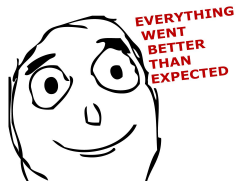


Table of contents

- 1 Databases
- 2 Uncertainty
- 3 Overview of my PhD Research**
- 4 Probabilities and Provenance on Trees and Treelike Instances
- 5 Conclusion

Roadmap

I studied different questions related to uncertainty:

Roadmap

I studied different questions related to uncertainty:

- Representing and querying uncertain **ordered data**
 - Possibility and certainty on **ordered relations**
Preprint: A., Ba, Deutch, Senellart 2016
 - Completing uncertain ordered **numerical values**
Preprint: A., Amsterdamer, Milo, Senellart 2016

Roadmap

I studied different questions related to uncertainty:

- Representing and querying uncertain **ordered data**
 - Possibility and certainty on **ordered relations**
Preprint: A., Ba, Deutch, Senellart 2016
 - Completing uncertain ordered **numerical values**
Preprint: A., Amsterdamer, Milo, Senellart 2016
- Reasoning on **incomplete data** under constraints
 - Combining several decidable **reasoning languages**
A., Benedikt 2015a, IJCAI'15
 - Addressing the **finiteness** hypothesis
A., Benedikt 2015b, LICS'15; Thesis Part II

Roadmap

I studied different questions related to uncertainty:

- Representing and querying uncertain **ordered data**
 - Possibility and certainty on **ordered relations**
Preprint: A., Ba, Deutch, Senellart 2016
 - Completing uncertain ordered **numerical values**
Preprint: A., Amsterdamer, Milo, Senellart 2016
- Reasoning on **incomplete data** under constraints
 - Combining several decidable **reasoning languages**
A., Benedikt 2015a, IJCAI'15
 - Addressing the **finiteness** hypothesis
A., Benedikt 2015b, LICS'15; Thesis Part II
- Query evaluation on **treelike** probabilistic data
A., Bourhis, Senellart 2015, 2016, ICALP'15, PODS'16; Thesis Part I

Roadmap

I studied different questions related to uncertainty:

- Representing and querying uncertain **ordered data**
 - Possibility and certainty on **ordered relations**
Preprint: A., Ba, Deutch, Senellart 2016
 - Completing uncertain ordered **numerical values**
Preprint: A., Amsterdamer, Milo, Senellart 2016
- Reasoning on **incomplete data** under constraints
 - Combining several decidable **reasoning languages**
A., Benedikt 2015a, IJCAI'15
 - Addressing the **finiteness** hypothesis
A., Benedikt 2015b, LICS'15; Thesis Part II
- Query evaluation on **treelike** probabilistic data
A., Bourhis, Senellart 2015, 2016, ICALP'15, PODS'16; Thesis Part I

Other work: (A. 2014, 2015a,b; A., Allauzen, Mohri 2015; A., Amsterdamer, Milo 2014a,b; A., Maniu, Senellart 2015; A., Galárraga, Preda, Suchanek 2014; Talaika, Biega, A., Suchanek 2015; Tang, A., Senellart, Bressan 2014a,b)

Uncertain ordered relations

Food

tiramisu kougelhopf

bretzel

munster

Drinks

champagne

riesling

Uncertain ordered relations

Food

tiramisu kougelhopf ● I partially know guest preferences

bretzel

munster

Drinks

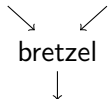
champagne

riesling

Uncertain ordered relations

Food

tiramisu kougelhopf



munster

- I **partially know** guest preferences

Drinks

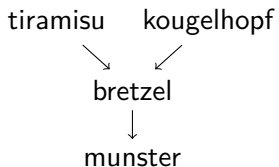
champagne



riesling

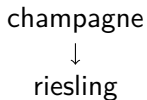
Uncertain ordered relations

Food



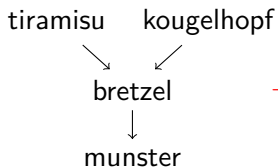
- I **partially know** guest preferences
- What should my parents bring?

Drinks



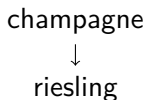
Uncertain ordered relations

Food



- I **partially know** guest preferences
 - What should my parents bring?
- What are the top two **Alsatian products**?

Drinks



Uncertain ordered relations

Food

tiramisu

kougelhopf



bretzel

munster

- I **partially know** guest preferences
 - What should my parents bring?
- What are the top two **Alsatian products**?

Drinks

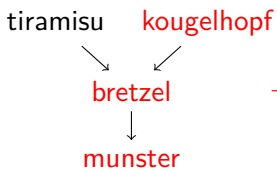
champagne



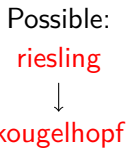
riesling

Uncertain ordered relations

Food



- I **partially know** guest preferences
 - What should my parents bring?
- What are the top two **Alsatian products**?

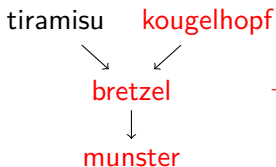


Drinks



Uncertain ordered relations

Food



- I **partially know** guest preferences
 - What should my parents bring?
- What are the top two **Alsatian products**?

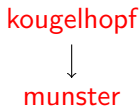
Drinks



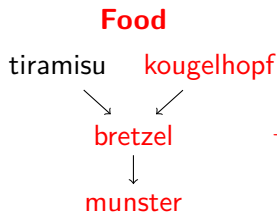
Possible:



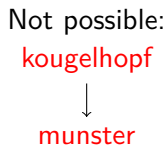
Not possible:



Uncertain ordered relations



- I **partially know** guest preferences
 - What should my parents bring?
- What are the top two **Alsatian products**?



- On which queries and data can we **efficiently** find possible and certain answers?

Uncertain numerical values

- How much **food** do people eat?

Uncertain numerical values

- How much **food** do people eat?
- Let's ask **friends** who defeneded recently

Uncertain numerical values

small		small
sweet	tiny	salty
	both	

medium		medium
sweet	small	salty
	both	

large		large
sweet	medium	salty
	both	

large
both

- How much **food** do people eat?
- Let's ask **friends** who defended recently

Uncertain numerical values

small

sweet

tiny

both

small

salty

medium

sweet

small

both

medium

salty

- How much **food** do people eat?
- Let's ask **friends** who defended recently

large

sweet

medium

both

large

salty

large

both

Uncertain numerical values

small		small
sweet	tiny	salty
	both	

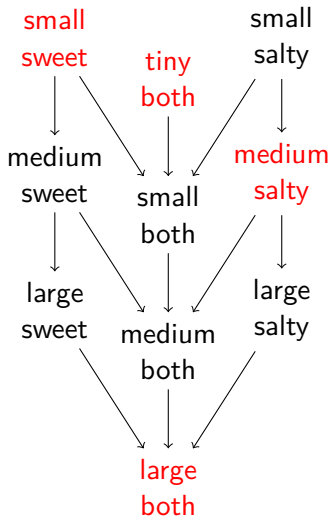
medium		medium
sweet	small	salty
	both	

large		large
sweet	medium	salty
	both	

large
both

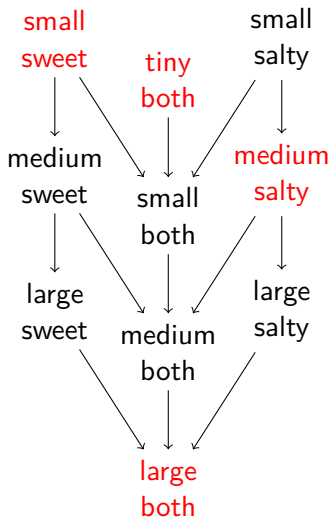
- How much **food** do people eat?
- Let's ask **friends** who defended recently
- Some **order relations** are implied

Uncertain numerical values



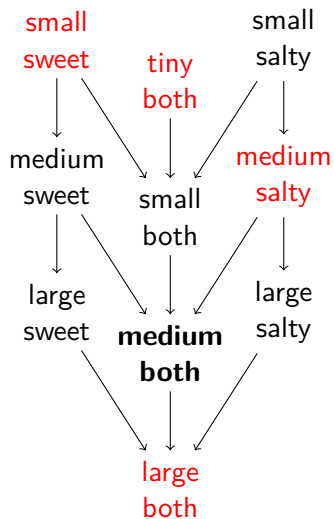
- How much **food** do people eat?
- Let's ask **friends** who defended recently
- Some **order relations** are implied

Uncertain numerical values



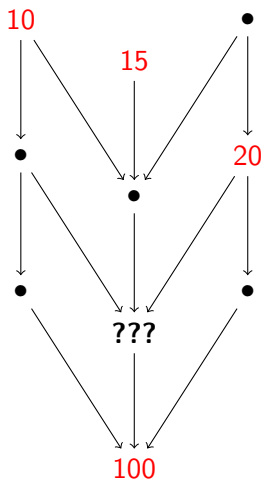
- How much **food** do people eat?
- Let's ask **friends** who defended recently
- Some **order relations** are implied
- How to **estimate** for my own defense?

Uncertain numerical values



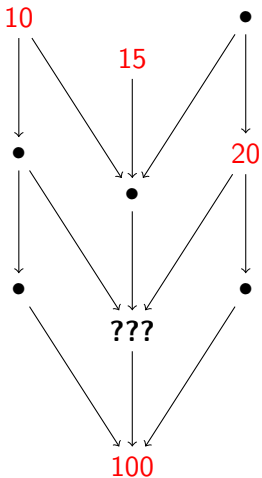
- How much **food** do people eat?
- Let's ask **friends** who defended recently
- Some **order relations** are implied
- How to **estimate** for my own defense?

Uncertain numerical values



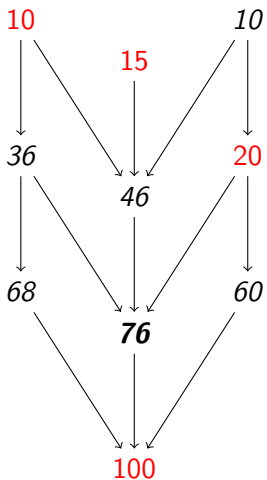
- How much **food** do people eat?
- Let's ask **friends** who defended recently
- Some **order relations** are implied
- How to **estimate** for my own defense?

Uncertain numerical values



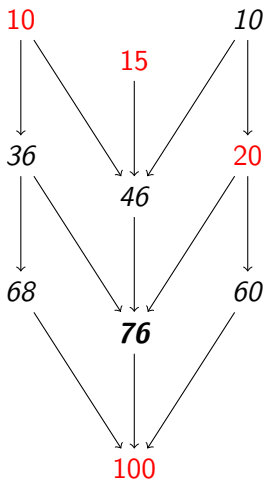
- How much **food** do people eat?
 - Let's ask **friends** who defended recently
 - Some **order relations** are implied
 - How to **estimate** for my own defense?
- Interpolation scheme for **posets** based on integration on polytopes

Uncertain numerical values



- How much **food** do people eat?
 - Let's ask **friends** who defended recently
 - Some **order relations** are implied
 - How to **estimate** for my own defense?
- Interpolation scheme for **posets** based on integration on polytopes

Uncertain numerical values



- How much **food** do people eat?
 - Let's ask **friends** who defended recently
 - Some **order relations** are implied
 - How to **estimate** for my own defense?
- Interpolation scheme for **posets** based on integration on polytopes
- **Complexity study** and **tractable cases**

Open-world query answering

Incomplete data:

- Fabian **supervises** Luis
- Fabian is **at the defense**
- Fabian is **in DBWeb**

Open-world query answering



Incomplete data:

- Fabian **supervises** Luis
- Fabian is **at the defense**
- Fabian is **in DBWeb**



Logical constraints:

- People at the defense will have drinks
- All DBWeb students will have drinks
- If your advisor is in DBWeb then you are a DBWeb student

Open-world query answering



Incomplete data:

- Fabian **supervises** Luis
- Fabian is **at the defense**
- Fabian is **in DBWeb**



Logical constraints:

- People at the defense will have drinks
- All DBWeb students will have drinks
- If your advisor is in DBWeb then you are a DBWeb student



Is the following query **certain**?

→ Will a DBWeb student meet their advisor at the drinks?

Open-world query answering



Incomplete data:

- Fabian **supervises** Luis
- Fabian is **at the defense**
- Fabian is **in DBWeb**
- **Fabian comes to the drinks**



Logical constraints:

- People at the defense will have drinks
- All DBWeb students will have drinks
- If your advisor is in DBWeb then you are a DBWeb student



Is the following query **certain**?

→ Will a DBWeb student meet their advisor at the drinks?

Open-world query answering



Incomplete data:

- Fabian **supervises** Luis
- Fabian is **at the defense**
- Fabian is **in DBWeb**
- **Fabian comes to the drinks**
- **Luis is a DBWeb student**



Logical constraints:

- People at the defense will have drinks
- All DBWeb students will have drinks
- If your advisor is in DBWeb then you are a DBWeb student



Is the following query **certain**?

→ Will a DBWeb student meet their advisor at the drinks?

Open-world query answering



Incomplete data:

- Fabian **supervises** Luis
- Fabian is **at the defense**
- Fabian is **in DBWeb**
- **Fabian comes to the drinks**
- Luis is a DBWeb student
- Luis comes to the drinks



Logical constraints:

- People at the defense will have drinks
- All DBWeb students will have drinks
- If your advisor is in DBWeb then you are a DBWeb student



Is the following query **certain**?

→ Will a DBWeb student meet their advisor at the drinks?

Open-world query answering



Incomplete data:

- Fabian **supervises** Luis
- Fabian is **at the defense**
- Fabian is **in DBWeb**
- **Fabian comes to the drinks**
- **Luis is a DBWeb student**
- **Luis comes to the drinks**



Logical constraints:

- People at the defense will have drinks
- All DBWeb students will have drinks
- If your advisor is in DBWeb then you are a DBWeb student



Is the following query **certain**?

→ Will a DBWeb student meet their advisor at the drinks?

→ **Yes!**

Open-world query answering



Incomplete data:

- Fabian **supervises** Luis
- Fabian is **at the defense**
- Fabian is **in DBWeb**
- **Fabian comes to the drinks**
- **Luis is a DBWeb student**
- **Luis comes to the drinks**



Logical constraints:

- People at the defense will have drinks
- All DBWeb students will have drinks
- If your advisor is in DBWeb then you are a DBWeb student



Is the following query **certain**?

→ Will a DBWeb student meet their advisor at the drinks?

→ **Yes!**

→ For which **rule languages** is this task decidable?

Expressive open-world query answering

Different **communities** use different kinds of **constraints**:

- Constraints with facts of **arity** > 2

Expressive open-world query answering

Different **communities** use different kinds of **constraints**:

- Constraints with facts of **arity** > 2
 - **Fabian** supervises **Luis**: arity 2
 - **Antoine**'s defense is in **B312** on **Monday**: arity 3

Expressive open-world query answering

Different **communities** use different kinds of **constraints**:

- Constraints with facts of **arity** > 2
 - **Fabian** supervises **Luis**: arity 2
 - **Antoine**'s defense is in **B312** on **Monday**: arity 3
- Constraints with **number restrictions**

Expressive open-world query answering

Different **communities** use different kinds of **constraints**:

- Constraints with facts of **arity > 2**
 - **Fabian** supervises **Luis**: arity 2
 - **Antoine**'s defense is in **B312** on **Monday**: arity 3
- Constraints with **number restrictions**
 - Everyone can invite **at most one** person
 - Students have **at most two** advisors

Expressive open-world query answering

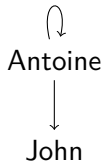
Different **communities** use different kinds of **constraints**:

- Constraints with facts of **arity** > 2
 - **Fabian** supervises **Luis**: arity 2
 - **Antoine's** defense is in **B312** on **Monday**: arity 3
 - Constraints with **number restrictions**
 - Everyone can invite **at most one** person
 - Students have **at most two** advisors
- I show that those can be **combined** under some restrictions to obtain **decidable** query answering

Query answering assuming finiteness

Consider the guests to the defense, \longrightarrow shows who invites whom

Data:

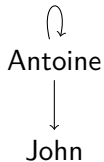


Rules:

Query answering assuming finiteness

Consider the guests to the defense, \longrightarrow shows who invites whom

Data:



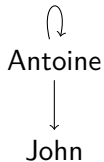
Rules:

- Each guest invites **someone**

Query answering assuming finiteness

Consider the guests to the defense, \longrightarrow shows who invites whom

Data:



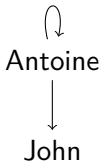
Rules:

- Each guest invites **someone**
- Nobody is invited by two **people**

Query answering assuming finiteness

Consider the guests to the defense, \longrightarrow shows who invites whom

Data:



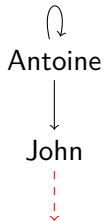
Rules:

- Each guest invites **someone**
 - Nobody is invited by two **people**
- \rightarrow Is this **sensible**?

Query answering assuming finiteness

Consider the guests to the defense, \longrightarrow shows who invites whom

Data:



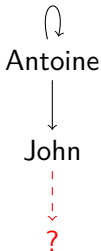
Rules:

- Each guest invites **someone**
 - Nobody is invited by two **people**
- \rightarrow Is this **sensible**?

Query answering assuming finiteness

Consider the guests to the defense, \longrightarrow shows who invites whom

Data:



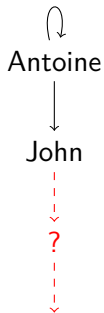
Rules:

- Each guest invites **someone**
 - Nobody is invited by two **people**
- \rightarrow Is this **sensible**?

Query answering assuming finiteness

Consider the guests to the defense, \longrightarrow shows who invites whom

Data:



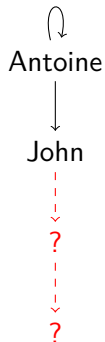
Rules:

- Each guest invites **someone**
 - Nobody is invited by two **people**
- \rightarrow Is this **sensible**?

Query answering assuming finiteness

Consider the guests to the defense, \longrightarrow shows who invites whom

Data:



Rules:

- Each guest invites **someone**
 - Nobody is invited by two **people**
- \rightarrow Is this **sensible**?

Query answering assuming finiteness

Consider the guests to the defense, \longrightarrow shows who invites whom

Data:



Rules:

- Each guest invites **someone**
 - Nobody is invited by two **people**
- \rightarrow Is this **sensible**?

Query answering assuming finiteness

Consider the guests to the defense, \longrightarrow shows who invites whom

Data:



Rules:

- Each guest invites **someone**
 - Nobody is invited by two **people**
- \rightarrow Is this **sensible**?

Query answering assuming finiteness

Consider the guests to the defense, \longrightarrow shows who invites whom

Data:



Rules:

- Each guest invites **someone**
 - Nobody is invited by two **people**
- \longrightarrow Is this **sensible? No!**

Query answering assuming finiteness

Consider the guests to the defense, \longrightarrow shows who invites whom

Data:



Rules:

- Each guest invites **someone**
- Nobody is invited by two **people**
- There are **finitely many guests!**

Query answering assuming finiteness

Consider the guests to the defense, \longrightarrow shows who invites whom

Data:



Rules:

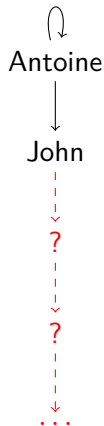
- Each guest invites **someone**
- Nobody is invited by two **people**
- There are **finitely many guests!**

\longrightarrow Can we do reasoning assuming **finiteness?**

Query answering assuming finiteness

Consider the guests to the defense, \longrightarrow shows who invites whom

Data:



Rules:

- Each guest invites **someone**
 - Nobody is invited by two **people**
 - There are **finitely many guests!**
- \longrightarrow Can we do reasoning assuming **finiteness?**
- \longrightarrow What **difference** does it make?

Finite open-world query answering

- I study the following constraints on **arbitrary arity**:
 - **Inclusion dependencies** with **one** exported element

Finite open-world query answering

- I study the following constraints on **arbitrary arity**:
 - **Inclusion dependencies** with **one** exported element
 - If x invites y **then** y invites some z

Finite open-world query answering

- I study the following constraints on **arbitrary arity**:
 - **Inclusion dependencies** with **one** exported element
 - If x invites y **then** y invites some z
 - **Functional dependencies**

Finite open-world query answering

- I study the following constraints on **arbitrary arity**:
 - **Inclusion dependencies** with **one** exported element
 - If x invites y **then** y invites some z
 - **Functional dependencies**
 - If x and y invite z **then** $x = y$

Finite open-world query answering

- I study the following constraints on **arbitrary arity**:
 - **Inclusion dependencies** with **one** exported element
 - If x invites y **then** y invites some z
 - **Functional dependencies**
 - If x and y invite z **then** $x = y$
- We can **compute** new constraints implied by finiteness using (Cosmadakis, Kanellakis, Vardi 1990)

Finite open-world query answering

- I study the following constraints on **arbitrary arity**:
 - **Inclusion dependencies** with **one** exported element
 - If x invites y **then** y invites some z
 - **Functional dependencies**
 - If x and y invite z **then** $x = y$
- We can **compute** new constraints implied by finiteness using (Cosmadakis, Kanellakis, Vardi 1990)
- With the new constraints, we can **forget** finiteness

Finite open-world query answering

- I study the following constraints on **arbitrary arity**:
 - **Inclusion dependencies** with **one** exported element
 - If x invites y **then** y invites some z
 - **Functional dependencies**
 - If x and y invite z **then** $x = y$
- We can **compute** new constraints implied by finiteness using (Cosmadakis, Kanellakis, Vardi 1990)
- With the new constraints, we can **forget** finiteness
- First techniques for open-world query answering with **arbitrary arity** signatures and **functional dependencies** where assuming **finiteness** makes a difference

Table of contents

- 1 Databases
- 2 Uncertainty
- 3 Overview of my PhD Research
- 4 Probabilities and Provenance on Trees and Treelike Instances**
- 5 Conclusion

Tuple-independent databases

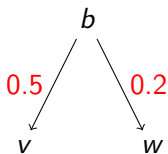
We fix a **relational signature** σ (here: S , arity 2).

<hr/>		
S		
<hr/>		
a	a	1
b	v	0.5
b	w	0.2
<hr/>		

Tuple-independent databases

We fix a **relational signature** σ (here: S , arity 2).

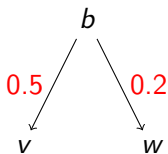
S		
a	a	1
b	v	0.5
b	w	0.2



Tuple-independent databases

We fix a **relational signature** σ (here: S , arity 2).

S		
a	a	1
b	v	0.5
b	w	0.2

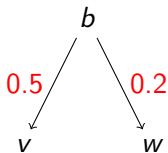


This TID instance represents the following **probability distribution**:

Tuple-independent databases

We fix a **relational signature** σ (here: S , arity 2).

S		
a	a	1
b	v	0.5
b	w	0.2



This TID instance represents the following **probability distribution**:

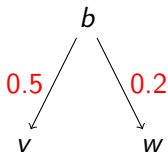
$$0.5 \times 0.2$$

S	
a	a
b	v
b	w

Tuple-independent databases

We fix a **relational signature** σ (here: S , arity 2).

S		
a	a	1
b	v	0.5
b	w	0.2



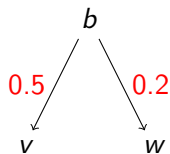
This TID instance represents the following **probability distribution**:

0.5×0.2		$0.5 \times (1 - 0.2)$	
S		S	
a	a	a	a
b	v	b	v
b	w		

Tuple-independent databases

We fix a **relational signature** σ (here: S , arity 2).

S		
a	a	1
b	v	0.5
b	w	0.2



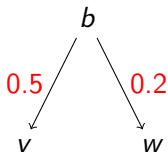
This TID instance represents the following **probability distribution**:

0.5×0.2		$0.5 \times (1 - 0.2)$		$(1 - 0.5) \times 0.2$	
S		S		S	
a	a	a	a	a	a
b	v	b	v		
b	w			b	w

Tuple-independent databases

We fix a **relational signature** σ (here: S , arity 2).

S		
a	a	1
b	v	0.5
b	w	0.2



This TID instance represents the following **probability distribution**:

0.5×0.2	$0.5 \times (1 - 0.2)$	$(1 - 0.5) \times 0.2$	$(1 - 0.5) \times (1 - 0.2)$
S	S	S	S
a a	a a	a a	a a
b v	b v		
b w		b w	

Query evaluation on probabilistic instances

We want to evaluate the probability of a **query** on a TID instance

Query evaluation on probabilistic instances

We want to evaluate the probability of a **query** on a TID instance

$$q : \exists xy R(x) \wedge S(x, y) \wedge T(y)$$

Query evaluation on probabilistic instances

We want to evaluate the probability of a **query** on a TID instance

$$q : \exists xy R(x) \wedge S(x, y) \wedge T(y)$$

R	
<i>a</i>	1
<i>b</i>	0.4
<i>c</i>	0.6

Query evaluation on probabilistic instances

We want to evaluate the probability of a **query** on a TID instance

$$q : \exists xy R(x) \wedge S(x, y) \wedge T(y)$$

R		S		
<i>a</i>	1	<i>a</i>	<i>a</i>	1
<i>b</i>	0.4	<i>b</i>	<i>v</i>	0.5
<i>c</i>	0.6	<i>b</i>	<i>w</i>	0.2

Query evaluation on probabilistic instances

We want to evaluate the probability of a **query** on a TID instance

$$q : \exists xy R(x) \wedge S(x, y) \wedge T(y)$$

R	S	T
<i>a</i> 1	<i>a a</i> 1	<i>v</i> 0.3
<i>b</i> 0.4	<i>b v</i> 0.5	<i>w</i> 0.7
<i>c</i> 0.6	<i>b w</i> 0.2	<i>b</i> 1

Query evaluation on probabilistic instances

We want to evaluate the probability of a **query** on a TID instance

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

R	S	T
<i>a</i> 1	<i>a a</i> 1	<i>v</i> 0.3
<i>b</i> 0.4	<i>b v</i> 0.5	<i>w</i> 0.7
<i>c</i> 0.6	<i>b w</i> 0.2	<i>b</i> 1

Query evaluation on probabilistic instances

We want to evaluate the probability of a **query** on a TID instance

$$q : \exists xy R(x) \wedge S(x, y) \wedge T(y)$$

R	S	T
<i>a</i> 1	<i>a a</i> 1	<i>v</i> 0.3
<i>b</i> 0.4	<i>b v</i> 0.5	<i>w</i> 0.7
<i>c</i> 0.6	<i>b w</i> 0.2	<i>b</i> 1

Query evaluation on probabilistic instances

We want to evaluate the probability of a **query** on a TID instance

$$q : \exists xy R(x) \wedge S(x, y) \wedge T(y)$$

R		S			T	
<i>a</i>	1	<i>a</i>	<i>a</i>	1	<i>v</i>	0.3
<i>b</i>	0.4	<i>b</i>	<i>v</i>	0.5	<i>w</i>	0.7
<i>c</i>	0.6	<i>b</i>	<i>w</i>	0.2	<i>b</i>	1

- The query is true iff $R(b)$ is here and one of:

Query evaluation on probabilistic instances

We want to evaluate the probability of a **query** on a TID instance

$$q : \exists xy R(x) \wedge S(x, y) \wedge T(y)$$

R		S			T	
<i>a</i>	1	<i>a</i>	<i>a</i>	1	<i>v</i>	0.3
<i>b</i>	0.4	<i>b</i>	<i>v</i>	0.5	<i>w</i>	0.7
<i>c</i>	0.6	<i>b</i>	<i>w</i>	0.2	<i>b</i>	1

- The query is true iff $R(b)$ is here and one of:
 - $S(b, v)$ and $T(v)$ are here

Query evaluation on probabilistic instances

We want to evaluate the probability of a **query** on a TID instance

$$q : \exists xy R(x) \wedge S(x, y) \wedge T(y)$$

R		S			T	
<i>a</i>	1	<i>a</i>	<i>a</i>	1	<i>v</i>	0.3
<i>b</i>	0.4	<i>b</i>	<i>v</i>	0.5	<i>w</i>	0.7
<i>c</i>	0.6	<i>b</i>	<i>w</i>	0.2	<i>b</i>	1

- The query is true iff $R(b)$ is here and one of:
 - $S(b, v)$ and $T(v)$ are here
 - $S(b, w)$ and $T(w)$ are here

Query evaluation on probabilistic instances

We want to evaluate the probability of a **query** on a TID instance

$$q : \exists xy R(x) \wedge S(x, y) \wedge T(y)$$

R		S			T	
<i>a</i>	1	<i>a</i>	<i>a</i>	1	<i>v</i>	0.3
<i>b</i>	0.4	<i>b</i>	<i>v</i>	0.5	<i>w</i>	0.7
<i>c</i>	0.6	<i>b</i>	<i>w</i>	0.2	<i>b</i>	1

- The query is true iff $R(b)$ is here and one of:
 - $S(b, v)$ and $T(v)$ are here
 - $S(b, w)$ and $T(w)$ are here

→ Probability:

Query evaluation on probabilistic instances

We want to evaluate the probability of a **query** on a TID instance

$$q : \exists xy R(x) \wedge S(x, y) \wedge T(y)$$

R		S			T	
<i>a</i>	1	<i>a</i>	<i>a</i>	1	<i>v</i>	0.3
<i>b</i>	0.4	<i>b</i>	<i>v</i>	0.5	<i>w</i>	0.7
<i>c</i>	0.6	<i>b</i>	<i>w</i>	0.2	<i>b</i>	1

- The query is true iff $R(b)$ is here and one of:
 - $S(b, v)$ and $T(v)$ are here
 - $S(b, w)$ and $T(w)$ are here

→ Probability:

$$0.4 \times$$

Query evaluation on probabilistic instances

We want to evaluate the probability of a **query** on a TID instance

$$q : \exists xy R(x) \wedge S(x, y) \wedge T(y)$$

R		S			T	
a	1	a	a	1	v	0.3
b	0.4	b	v	0.5	w	0.7
c	0.6	b	w	0.2	b	1

- The query is true iff $R(b)$ is here and one of:
 - $S(b, v)$ and $T(v)$ are here
 - $S(b, w)$ and $T(w)$ are here

→ **Probability:**

$$0.4 \times (1 -$$

Query evaluation on probabilistic instances

We want to evaluate the probability of a **query** on a TID instance

$$q : \exists xy R(x) \wedge S(x, y) \wedge T(y)$$

R		S			T	
a	1	a	a	1	v	0.3
b	0.4	b	v	0.5	w	0.7
c	0.6	b	w	0.2	b	1

- The query is true iff $R(b)$ is here and one of:
 - $S(b, v)$ and $T(v)$ are here
 - $S(b, w)$ and $T(w)$ are here

→ **Probability:**

$$0.4 \times (1 - (1 - 0.5 \times 0.3))$$

Query evaluation on probabilistic instances

We want to evaluate the probability of a **query** on a TID instance

$$q : \exists xy R(x) \wedge S(x, y) \wedge T(y)$$

R		S			T	
<i>a</i>	1	<i>a</i>	<i>a</i>	1	<i>v</i>	0.3
<i>b</i>	0.4	<i>b</i>	<i>v</i>	0.5	<i>w</i>	0.7
<i>c</i>	0.6	<i>b</i>	<i>w</i>	0.2	<i>b</i>	1

- The query is true iff $R(b)$ is here and one of:
 - $S(b, v)$ and $T(v)$ are here
 - $S(b, w)$ and $T(w)$ are here

→ **Probability:**

$$0.4 \times (1 - (1 - 0.5 \times 0.3) \times (1 - 0.2 \times 0.7))$$

Query evaluation on probabilistic instances

We want to evaluate the probability of a **query** on a TID instance

$$q : \exists xy R(x) \wedge S(x, y) \wedge T(y)$$

R		S			T	
a	1	a	a	1	v	0.3
b	0.4	b	v	0.5	w	0.7
c	0.6	b	w	0.2	b	1

- The query is true iff $R(b)$ is here and one of:
 - $S(b, v)$ and $T(v)$ are here
 - $S(b, w)$ and $T(w)$ are here

→ **Probability:**

$$0.4 \times (1 - (1 - 0.5 \times 0.3) \times (1 - 0.2 \times 0.7)) = 0.1076$$

Complexity of probabilistic query evaluation (PQE)

What is the **data complexity** of probabilistic query evaluation depending on the class \mathcal{Q} of **queries** and class \mathcal{I} of **instances**?

Complexity of probabilistic query evaluation (PQE)

What is the **data complexity** of probabilistic query evaluation depending on the class \mathcal{Q} of **queries** and class \mathcal{I} of **instances**?

- **Existing dichotomy result:** (Dalvi, Suciu 2012)
 - \mathcal{Q} are (unions of) conjunctive queries, \mathcal{I} is all instances
 - There is a class $\mathcal{S} \subseteq \mathcal{Q}$ of **safe queries**

Complexity of probabilistic query evaluation (PQE)

What is the **data complexity** of probabilistic query evaluation depending on the class \mathcal{Q} of **queries** and class \mathcal{I} of **instances**?

- **Existing dichotomy result:** (Dalvi, Suciu 2012)
 - \mathcal{Q} are (unions of) conjunctive queries, \mathcal{I} is all instances
 - There is a class $\mathcal{S} \subseteq \mathcal{Q}$ of **safe queries**
 - PQE is **PTIME** for any $q \in \mathcal{S}$ on all instances

Complexity of probabilistic query evaluation (PQE)

What is the **data complexity** of probabilistic query evaluation depending on the class \mathcal{Q} of **queries** and class \mathcal{I} of **instances**?

- **Existing dichotomy result:** (Dalvi, Suciu 2012)
 - \mathcal{Q} are (unions of) conjunctive queries, \mathcal{I} is all instances
 - There is a class $\mathcal{S} \subseteq \mathcal{Q}$ of **safe queries**
 - PQE is **PTIME** for any $q \in \mathcal{S}$ on all instances
 - PQE is **#P-hard** for any $q \in \mathcal{Q} \setminus \mathcal{S}$ on all instances

Complexity of probabilistic query evaluation (PQE)

What is the **data complexity** of probabilistic query evaluation depending on the class \mathcal{Q} of **queries** and class \mathcal{I} of **instances**?

- **Existing dichotomy result:** (Dalvi, Suciu 2012)
 - \mathcal{Q} are (unions of) conjunctive queries, \mathcal{I} is all instances
 - There is a class $\mathcal{S} \subseteq \mathcal{Q}$ of **safe queries**
 - PQE is **PTIME** for any $q \in \mathcal{S}$ on all instances
 - PQE is **#P-hard** for any $q \in \mathcal{Q} \setminus \mathcal{S}$ on all instances
 - $q : \exists xy R(x) \wedge S(x, y) \wedge T(y)$ is **unsafe!**

Complexity of probabilistic query evaluation (PQE)

What is the **data complexity** of probabilistic query evaluation depending on the class \mathcal{Q} of **queries** and class \mathcal{I} of **instances**?

- **Existing dichotomy result:** (Dalvi, Suciu 2012)
 - \mathcal{Q} are (unions of) conjunctive queries, \mathcal{I} is all instances
 - There is a class $\mathcal{S} \subseteq \mathcal{Q}$ of **safe queries**
 - PQE is **PTIME** for any $q \in \mathcal{S}$ on all instances
 - PQE is **#P-hard** for any $q \in \mathcal{Q} \setminus \mathcal{S}$ on all instances
 - $q : \exists xy R(x) \wedge S(x, y) \wedge T(y)$ is **unsafe!**

Is there a **smaller class** \mathcal{I} such that PQE is tractable for a **larger** \mathcal{Q} ?

Trees and treelike instances

- **Goal:** find an **instance class** \mathcal{I} where PQE is tractable

Trees and treelike instances

- **Goal:** find an **instance class** \mathcal{I} where PQE is tractable
- **Idea:** let \mathcal{I} be **treelike instances** (constant bound on **treewidth**)

Trees and treelike instances

- **Goal:** find an **instance class** \mathcal{I} where PQE is tractable
- **Idea:** let \mathcal{I} be **treelike instances** (constant bound on **treewidth**)
 - **Trees** have treewidth **1**
 - **Cycles** have treewidth **2**
 - **k -cliques** and **$(k - 1)$ -grids** have treewidth **$k - 1$**

Trees and treelike instances

- **Goal:** find an **instance class** \mathcal{I} where PQE is tractable
 - **Idea:** let \mathcal{I} be **treelike instances** (constant bound on **treewidth**)
 - **Trees** have treewidth **1**
 - **Cycles** have treewidth **2**
 - **k -cliques** and **$(k - 1)$ -grids** have treewidth **$k - 1$**
- Known results (Courcelle 1990):
- \mathcal{I} : **treelike** instances; \mathcal{Q} : **monadic second-order** queries
 - **non-probabilistic** QE is in **linear time**

Trees and treelike instances

- **Goal:** find an **instance class** \mathcal{I} where PQE is tractable
 - **Idea:** let \mathcal{I} be **treelike instances** (constant bound on **treewidth**)
 - **Trees** have treewidth **1**
 - **Cycles** have treewidth **2**
 - **k -cliques** and **$(k - 1)$ -grids** have treewidth **$k - 1$**
- Known results (Courcelle 1990):
- \mathcal{I} : **treelike instances**; \mathcal{Q} : **monadic second-order queries**
 - **non-probabilistic QE** is in **linear time**
- Does this extend to **probabilistic QE**?

Our main result

An **instance-based** dichotomy result:

Upper bound.

For \mathcal{I} the **treelike** instances and \mathcal{Q} the **MSO queries**

→ PQE is in **linear time** modulo arithmetic costs

Our main result

An **instance-based** dichotomy result:

Upper bound.

For \mathcal{I} the **treelike** instances and \mathcal{Q} the **MSO queries**

→ PQE is in **linear time** modulo arithmetic costs

- Also for expressive **provenance representations**
- Also with bounded-treewidth **correlations**

Our main result

An **instance-based** dichotomy result:

Upper bound.

For \mathcal{I} the **treelike** instances and \mathcal{Q} the **MSO queries**

- PQE is in **linear time** modulo arithmetic costs
- Also for expressive **provenance representations**
- Also with bounded-treewidth **correlations**

Lower bound.

For **any** unbounded-tw family \mathcal{I} and \mathcal{Q} the **FO queries**

- PQE is **#P-hard under RP reductions** assuming:
 - Signature **arity is 2** (graphs)
 - High-tw instances in \mathcal{I} are **easily constructible**

Technical tool: lineages

The **lineage** of a query q on an instance I :

- Boolean function ϕ whose **variables** are the facts of I
- A subinstance of I satisfies q **iff** ϕ is true for that valuation

Technical tool: lineages

The **lineage** of a query q on an instance I :

- Boolean function ϕ whose **variables** are the facts of I
- A subinstance of I satisfies q **iff** ϕ is true for that valuation

Example query: $\exists xyz R(x, y) \wedge R(y, z)$

Technical tool: lineages

The **lineage** of a query q on an instance I :

- Boolean function ϕ whose **variables** are the facts of I
- A subinstance of I satisfies q iff ϕ is true for that valuation

Example query: $\exists xyz R(x, y) \wedge R(y, z)$

<hr/>		
<i>R</i>		
<hr/>		
<i>a</i>	<i>b</i>	<i>f</i> ₁
<i>b</i>	<i>c</i>	<i>f</i> ₂
<i>d</i>	<i>e</i>	<i>f</i> ₃
<i>e</i>	<i>d</i>	<i>f</i> ₄
<i>f</i>	<i>f</i>	<i>f</i> ₅

Technical tool: lineages

The **lineage** of a query q on an instance I :

- Boolean function ϕ whose **variables** are the facts of I
- A subinstance of I satisfies q iff ϕ is true for that valuation

Example query: $\exists xyz R(x, y) \wedge R(y, z)$

<hr/>		
R		
<hr/>		
a	b	f_1
b	c	f_2
d	e	f_3
e	d	f_4
f	f	f_5
<hr/>		

→ **Lineage:** $(f_1 \wedge f_2)$

Technical tool: lineages

The **lineage** of a query q on an instance I :

- Boolean function ϕ whose **variables** are the facts of I
- A subinstance of I satisfies q iff ϕ is true for that valuation

Example query: $\exists xyz R(x, y) \wedge R(y, z)$

R		
a	b	f_1
b	c	f_2
d	e	f_3
e	d	f_4
f	f	f_5

→ Lineage: $(f_1 \wedge f_2) \vee (f_3 \wedge f_4)$

Technical tool: lineages

The **lineage** of a query q on an instance I :

- Boolean function ϕ whose **variables** are the facts of I
- A subinstance of I satisfies q iff ϕ is true for that valuation

Example query: $\exists xyz R(x, y) \wedge R(y, z)$

<hr/>		
R		
<hr/>		
a	b	f_1
b	c	f_2
d	e	f_3
e	d	f_4
f	f	f_5
<hr/>		

→ **Lineage:** $(f_1 \wedge f_2) \vee (f_3 \wedge f_4) \vee f_5$

Technical tool: lineages

The **lineage** of a query q on an instance I :

- Boolean function ϕ whose **variables** are the facts of I
- A subinstance of I satisfies q iff ϕ is true for that valuation

Example query: $\exists xyz R(x, y) \wedge R(y, z)$

R		
a	b	f_1
b	c	f_2
d	e	f_3
e	d	f_4
f	f	f_5

→ **Lineage:** $(f_1 \wedge f_2) \vee (f_3 \wedge f_4) \vee f_5$

Technical tool: lineages

The **lineage** of a query q on an instance I :

- Boolean function ϕ whose **variables** are the facts of I
- A subinstance of I satisfies q **iff** ϕ is true for that valuation

Example query: $\exists xyz R(x, y) \wedge R(y, z)$

R		
a	b	f_1
b	c	f_2
d	e	f_3
e	d	f_4
f	f	f_5

→ **Lineage:** $(f_1 \wedge f_2) \vee (f_3 \wedge f_4) \vee f_5$

→ For all $\nu : I \rightarrow \{0, 1\}$ we have $\nu(\phi) = 1$ **iff** $\{F \in I \mid \nu(F) = 1\} \models q$

Using lineages

- To solve the PQE problem on **treelike instances** for **MSO**

Using lineages

- To solve the PQE problem on **treelike instances** for **MSO**
 - First solve the problem on **trees**

Using lineages

- To solve the PQE problem on **treelike instances** for **MSO**
 - First solve the problem on **trees**
 - Then use the results of (Courcelle 1990)

Using lineages

- To solve the PQE problem on **treelike instances** for **MSO**
 - First solve the problem on **trees**
 - Then use the results of (Courcelle 1990)
- Use **lineage** for PQE:

Using lineages

- To solve the PQE problem on **treelike instances** for **MSO**
 - First solve the problem on **trees**
 - Then use the results of (Courcelle 1990)
- Use **lineage** for PQE:
 - Compute a lineage representation **efficiently**

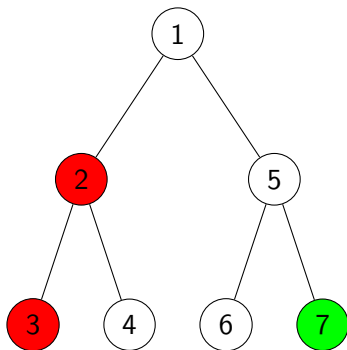
Using lineages

- To solve the PQE problem on **treelike instances** for **MSO**
 - First solve the problem on **trees**
 - Then use the results of (Courcelle 1990)
- Use **lineage** for PQE:
 - Compute a lineage representation **efficiently**
 - Probability of the **lineage** = probability of the **query**

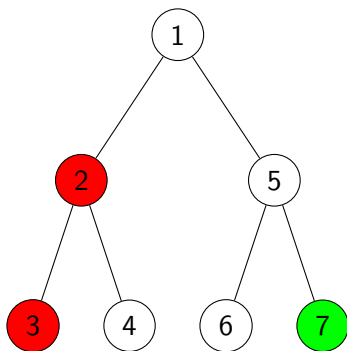
Using lineages

- To solve the PQE problem on **treelike instances** for **MSO**
 - First solve the problem on **trees**
 - Then use the results of (Courcelle 1990)
- Use **lineage** for PQE:
 - Compute a lineage representation **efficiently**
 - Probability of the **lineage** = probability of the **query**
 - Compute the lineage probability **efficiently**
(show it is not **#P-hard** as in the general case)

Uncertain trees

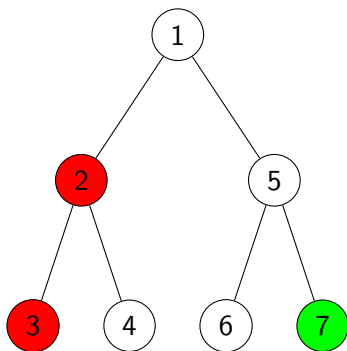


Uncertain trees



A **valuation** of a tree decides whether to **keep** or **discard** node labels.

Uncertain trees



A **valuation** of a tree decides whether to **keep** or **discard** node labels.

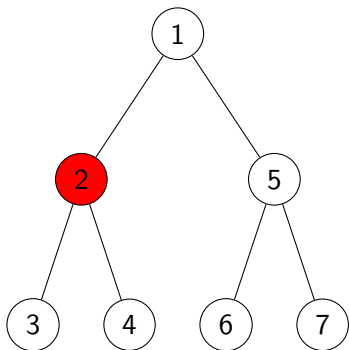
Example query:

“Is there both a red and a green node?”

Valuation: $\{2, 3, 7\}$

The query is **true**

Uncertain trees



A **valuation** of a tree decides whether to **keep** or **discard** node labels.

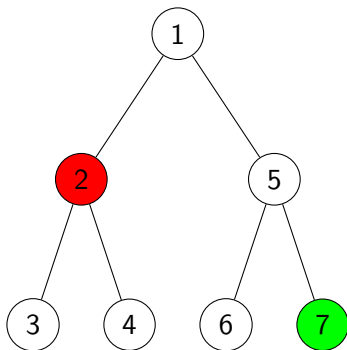
Example query:

“Is there both a red and a green node?”

Valuation: {2}

The query is **false**

Uncertain trees



A **valuation** of a tree decides whether to **keep** or **discard** node labels.

Example query:

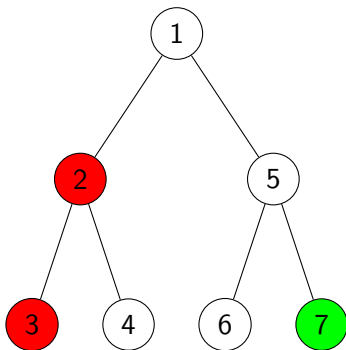
“Is there both a red and a green node?”

Valuation: $\{2, 7\}$

The query is **true**

Lineage circuits on trees

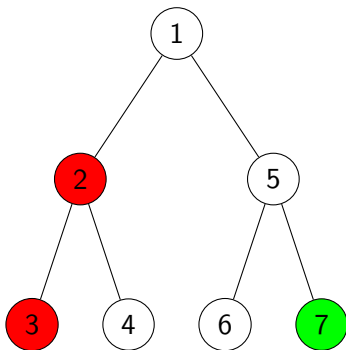
- q : Is there both a **red** and a **green** node?
- Which **valuations** satisfy q ?



Lineage circuits on trees

q : Is there both a red and a green node?

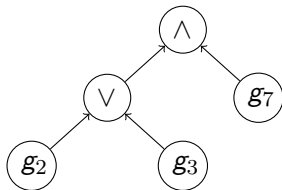
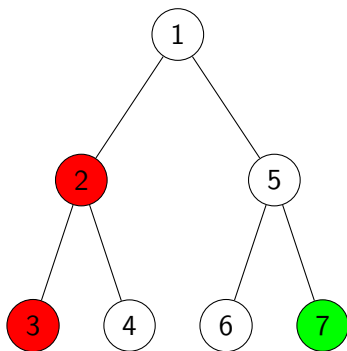
- Which valuations satisfy q ?
 - Lineage circuit of a query q on an uncertain tree T
 - Boolean circuit C
 - with input gates g_2, g_3, g_7
- $\nu(T)$ satisfies q iff $\nu(C)$ is true



Lineage circuits on trees

q : Is there both a red and a green node?

- Which **valuations** satisfy q ?
 - **Lineage circuit** of a query q on an uncertain tree T
 - **Boolean circuit** C
 - with **input gates** g_2, g_3, g_7
- $\nu(T)$ satisfies q iff $\nu(C)$ is **true**



Our main results

Theorem

*For any query q given as a bottom-up **tree automaton** A , for any input **tree** T , we can build a **lineage circuit** of A on T in **linear time** in A and T .*

Our main results

Theorem

*For any query q given as a bottom-up **tree automaton** A , for any input **tree** T , we can build a **lineage circuit** of A on T in **linear time** in A and T .*

MSO on treelike instances \Rightarrow MSO on trees (Courcelle 1990).

Our main results

Theorem

For any query q given as a bottom-up *tree automaton* A ,
for any input *tree* T , we can build a *lineage circuit* of A on T
in *linear time* in A and T .

MSO on treelike instances \Rightarrow MSO on trees (Courcelle 1990).

Theorem

For any fixed *MSO query* q and $k \in \mathbb{N}$,
for any input *instance* I of *treewidth* $\leq k$,
we can build in *linear time* in I a lineage circuit of q on I .

Our main results

Theorem

For any query q given as a bottom-up *tree automaton* A ,
for any input *tree* T , we can build a *lineage circuit* of A on T
in *linear time* in A and T .

MSO on treelike instances \Rightarrow MSO on trees (Courcelle 1990).

Theorem

For any fixed *MSO query* q and $k \in \mathbb{N}$,
for any input *instance* I of *treewidth* $\leq k$,
we can build in *linear time* in I a lineage circuit of q on I .

The lineage circuits are themselves *treelike*, hence:

Our main results

Theorem

For any query q given as a bottom-up *tree automaton* A ,
for any input *tree* T , we can build a *lineage circuit* of A on T
in *linear time* in A and T .

MSO on treelike instances \Rightarrow MSO on trees (Courcelle 1990).

Theorem

For any fixed *MSO query* q and $k \in \mathbb{N}$,
for any input *instance* I of *treewidth* $\leq k$,
we can build in *linear time* in I a lineage circuit of q on I .

The lineage circuits are themselves *treelike*, hence:

Corollary

Probabilistic query evaluation of MSO queries on treelike instances is in linear time up to arithmetic operations.

Extension 1: general semirings

- **Semiring** of positive Boolean functions ($\text{PosBool}[X], \vee, \wedge, f, t$)

Extension 1: general semirings

- **Semiring** of positive Boolean functions ($\text{PosBool}[X], \vee, \wedge, f, t$)
- **Provenance semirings:** (Green, Karvounarakis, Tannen 2007)
 - Provenance for **arbitrary (commutative) semirings**
 - For queries in the **positive relational algebra** and Datalog

Extension 1: general semirings

- **Semiring** of positive Boolean functions ($\text{PosBool}[X], \vee, \wedge, f, t$)
- **Provenance semirings:** (Green, Karvounarakis, Tannen 2007)
 - Provenance for **arbitrary (commutative) semirings**
 - For queries in the **positive relational algebra** and Datalog

Our construction can be extended to $\mathbb{N}[X]$ -provenance for conjunctive queries and **unions of conjunctive queries** (UCQ):

Extension 1: general semirings

- **Semiring** of positive Boolean functions ($\text{PosBool}[X], \vee, \wedge, f, t$)
- **Provenance semirings**: (Green, Karvounarakis, Tannen 2007)
 - Provenance for **arbitrary (commutative) semirings**
 - For queries in the **positive relational algebra** and Datalog

Our construction can be extended to $\mathbb{N}[X]$ -provenance for conjunctive queries and **unions of conjunctive queries** (UCQ):

Theorem

For any fixed **UCQ** q and $k \in \mathbb{N}$,
for any input **instance** I of **treewidth** $\leq k$,
we can build in **linear time** a $\mathbb{N}[X]$ -provenance circuit of q on I .

Extension 2: correlations

- Our **probabilistic instances** assume **independence** on all facts

Extension 2: correlations

- Our **probabilistic instances** assume **independence** on all facts
- More expressive: **Block-Independent Disjoint** instances:

Extension 2: correlations

- Our **probabilistic instances** assume **independence** on all facts
- More expressive: **Block-Independent Disjoint** instances:

<u>name</u>	favorite	p
john	kougelhopf	0.8
john	bretzel	0.2
jane	kougelhopf	0.1
jane	bretzel	0.9

Extension 2: correlations

- Our **probabilistic instances** assume **independence** on all facts
- More expressive: **Block-Independent Disjoint** instances:

<u>name</u>	favorite	p
john	kougelhopf	0.8
john	bretzel	0.2
jane	kougelhopf	0.1
jane	bretzel	0.9

Theorem

*Probabilistic query evaluation of MSO queries on treelike **BID** is in linear time up to arithmetic operations.*

Extension 2: correlations

- Our **probabilistic instances** assume **independence** on all facts
- More expressive: **Block-Independent Disjoint** instances:

<u>name</u>	favorite	p
john	kougelhopf	0.8
john	bretzel	0.2
jane	kougelhopf	0.1
jane	bretzel	0.9

Theorem

*Probabilistic query evaluation of MSO queries on treelike **BID** is in linear time up to arithmetic operations.*

Generalises to **pc-tables** with **treelike** correlations

Lower bound goal

- Class \mathcal{I} of **unbounded-treewidth instances**, query q in class \mathcal{Q}
- Show that **probabilistic query evaluation** of q on \mathcal{I} is **hard**

Lower bound goal

- Class \mathcal{I} of **unbounded-treewidth instances**, query q in class \mathcal{Q}
 - Show that **probabilistic query evaluation** of q on \mathcal{I} is **hard**
- Restrict to **arity-2** (= labeled graphs) for technical reasons

Lower bound goal

- Class \mathcal{I} of **unbounded-treewidth instances**, query q in class \mathcal{Q}
- Show that **probabilistic query evaluation** of q on \mathcal{I} is **hard**
- Restrict to **arity-2** (= labeled graphs) for technical reasons
- Impose that \mathcal{I} is **tw-constructible**:

Lower bound goal

- Class \mathcal{I} of **unbounded-treewidth instances**, query q in class \mathcal{Q}
- Show that **probabilistic query evaluation** of q on \mathcal{I} is **hard**
- Restrict to **arity-2** (= labeled graphs) for technical reasons
- Impose that \mathcal{I} is **tw-constructible**:
 - Given $k \in \mathbb{N}$, we can construct in **time $\text{Poly}(k)$**
an instance of \mathcal{I} of **treewidth $\geq k$**

Lower bound goal

- Class \mathcal{I} of **unbounded-treewidth instances**, query q in class \mathcal{Q}
- Show that **probabilistic query evaluation** of q on \mathcal{I} is **hard**
- Restrict to **arity-2** (= labeled graphs) for technical reasons
- Impose that \mathcal{I} is **tw-constructible**:
 - Given $k \in \mathbb{N}$, we can construct in **time $\text{Poly}(k)$** an instance of \mathcal{I} of **treewidth $\geq k$**

Theorem

There is a **first-order** query q such that for any unbounded-tw, tw-constructible, arity-2 **instance family \mathcal{I}** , probabilistic query eval for q on \mathcal{I} is **$\#P$ -hard** under RP reductions.

Lower bound goal

- Class \mathcal{I} of **unbounded-treewidth instances**, query q in class \mathcal{Q}
- Show that **probabilistic query evaluation** of q on \mathcal{I} is **hard**
- Restrict to **arity-2** (= labeled graphs) for technical reasons
- Impose that \mathcal{I} is **tw-constructible**:
 - Given $k \in \mathbb{N}$, we can construct in **time $\text{Poly}(k)$** an instance of \mathcal{I} of **treewidth $\geq k$**

Theorem

*There is a **first-order** query q such that for any unbounded-tw, tw-constructible, arity-2 **instance family** \mathcal{I} , probabilistic query eval for q on \mathcal{I} is **$\#P$ -hard** under RP reductions.*

Proven by extracting arbitrary graphs as **minors** of high-treewidth families using (Chekuri, Chuzhoy 2014)

Table of contents

- 1 Databases
- 2 Uncertainty
- 3 Overview of my PhD Research
- 4 Probabilities and Provenance on Trees and Treelike Instances
- 5 Conclusion**

Conclusion

Main contributions to the study of **uncertain data management**:

- New **decidable** languages to reason on **incomplete data**

Conclusion

Main contributions to the study of **uncertain data management**:

- New **decidable** languages to reason on **incomplete data**
- New techniques and results for **finite reasoning**

Conclusion

Main contributions to the study of **uncertain data management**:

- New **decidable** languages to reason on **incomplete data**
- New techniques and results for **finite reasoning**
- Representations and complexity for uncertain **ordered data**

Conclusion

Main contributions to the study of **uncertain data management**:

- New **decidable** languages to reason on **incomplete data**
- New techniques and results for **finite reasoning**
- Representations and complexity for uncertain **ordered data**
- **Instance-based dichotomy** for probabilistic instances:

Conclusion

Main contributions to the study of **uncertain data management**:

- New **decidable** languages to reason on **incomplete data**
- New techniques and results for **finite reasoning**
- Representations and complexity for uncertain **ordered data**
- **Instance-based dichotomy** for probabilistic instances:
 - **Tractable** data complexity for MSO on treelike families
(based on treelike lineage circuits via tree automata)

Conclusion

Main contributions to the study of **uncertain data management**:

- New **decidable** languages to reason on **incomplete data**
- New techniques and results for **finite reasoning**
- Representations and complexity for uncertain **ordered data**
- **Instance-based dichotomy** for probabilistic instances:
 - **Tractable** data complexity for MSO on treelike families
(*based on treelike lineage circuits via tree automata*)
 - Extends to general **provenance semirings** for UCQs

Conclusion

Main contributions to the study of **uncertain data management**:

- New **decidable** languages to reason on **incomplete data**
- New techniques and results for **finite reasoning**
- Representations and complexity for uncertain **ordered data**
- **Instance-based dichotomy** for probabilistic instances:
 - **Tractable** data complexity for MSO on treelike families
(*based on treelike lineage circuits via tree automata*)
 - Extends to general **provenance semirings** for UCQs
 - Extends to probabilistic **correlations**

Conclusion

Main contributions to the study of **uncertain data management**:

- New **decidable** languages to reason on **incomplete data**
- New techniques and results for **finite reasoning**
- Representations and complexity for uncertain **ordered data**
- **Instance-based dichotomy** for probabilistic instances:
 - **Tractable** data complexity for MSO on treelike families
(*based on treelike lineage circuits via tree automata*)
 - Extends to general **provenance semirings** for UCQs
 - Extends to probabilistic **correlations**
 - **Lower bound** for FO on **any** non-treelike family
(*assuming arity-two and treewidth-constructibility*)

Ongoing and future work

- **Probabilistic** query answering
 - Tractability in **combined complexity** for some queries
 - **Hybrid** tractability criteria based on instance and query
 - **Practical implementation** with partial decompositions

Ongoing and future work

- Probabilistic query answering
 - Tractability in **combined complexity** for some queries
 - **Hybrid** tractability criteria based on instance and query
 - **Practical implementation** with partial decompositions
- Open-world query answering
 - Managing **order relations** and transitive relations
 - Extending **provenance** techniques to open-world reasoning

Ongoing and future work





- **Probabilistic** query answering
 - Tractability in **combined complexity** for some queries
 - **Hybrid** tractability criteria based on instance and query
 - **Practical implementation** with partial decompositions
- **Open-world** query answering
 - Managing **order relations** and transitive relations
 - Extending **provenance** techniques to open-world reasoning

Thanks for your attention!

Image sources

- Slides 2 and 14:
<https://openclipart.org/download/163711/database-server.svg>
- Slide 3: SMSSecure <https://smssecure.org/> and AOSP
<https://source.android.com/>
- Slide 7: <https://openclipart.org/download/36529/interrogation.svg>
- Slide 8: <http://rtw.ml.cmu.edu/>,
<https://openclipart.org/download/25537/HMTL.svg>, and
<https://twitter.com/cmunell>
- Slide 9: <https://en.wikipedia.org/wiki/Template:Disputed>
- Slide 10: Zhang 2015, p. 9, Dong, Halevy, Yu 2009, p. 4,
[https://atlas.web.cern.ch/Atlas/GROUPS/PHYSICS/CONFNOTES/
ATLAS-CONF-2015-041/fig_06b.png](https://atlas.web.cern.ch/Atlas/GROUPS/PHYSICS/CONFNOTES/ATLAS-CONF-2015-041/fig_06b.png),
<https://code.google.com/p/transducersaurus/wiki/CascadeTutorial>,
<https://www.cs.washington.edu/robotics/mcl/>
- Slide 16:
<https://diaryofawhinyguy.files.wordpress.com/2013/01/rage-guy.png>
- Slide 17: [http://mylolface.com/assets/faces/
happy-everything-went-better-than-expected.jpg](http://mylolface.com/assets/faces/happy-everything-went-better-than-expected.jpg)





References I

-  **Amarilli, Antoine (2014)**. “The Possibility Problem for Probabilistic XML”. In: *Proc. AMW*. URL: http://ceur-ws.org/Vol-1189/paper_2.pdf.
-  **Amarilli, Antoine (2015a)**. “Possibility for Probabilistic XML”. In: *Ingénierie des Systèmes d'Information*. URL: <http://arxiv.org/abs/1404.3131>.
-  **Amarilli, Antoine (2015b)**. “Structurally Tractable Uncertain Data”. In: *Proc. PhD Symposium of SIGMOD/PODS*. URL: <http://arxiv.org/abs/1507.04955>.
-  **Amarilli, Antoine, Cyril Allauzen, Mehryar Mohri (2015)**. *Minimum Bayesian Risk Methods for Automatic Speech Recognition*. United States Patent 9123333. URL: <https://a3nm.net/publications/amarilli2014minimum.pdf>.







References II

-  Amarilli, Antoine, Yael Amsterdamer, Tova Milo (2014a). “On the Complexity of Mining Itemsets from the Crowd Using Taxonomies”. In: *Proc. ICDT*. URL: <http://arxiv.org/abs/1312.3248>.
-  Amarilli, Antoine, Yael Amsterdamer, Tova Milo (2014b). “Uncertainty in Crowd Data Sourcing Under Structural Constraints”. In: *Proc. UnCrowd*. URL: <http://arxiv.org/abs/1403.0783>.
-  Amarilli, Antoine, Michael Benedikt (2015a). “Combining Existential Rules and Description Logics”. In: *Proc. IJCAI*. URL: <http://arxiv.org/abs/1505.00326>.
-  Amarilli, Antoine, Michael Benedikt (2015b). “Finite Open-World Query Answering with Number Restrictions”. In: *Proc. LICS*. URL: <http://arxiv.org/abs/1505.04216>.


References III

-  Amarilli, Antoine, Pierre Bourhis, Pierre Senellart (2015). “Provenance Circuits for Trees and Treelike Instances”. In: *Proc. ICALP*. URL: <http://arxiv.org/abs/1511.08723>.
-  Amarilli, Antoine, Pierre Bourhis, Pierre Senellart (2016). “Tractable Lineages on Treelike Instances: Limits and Extensions”. In: *Proc. PODS*. To appear. URL: <https://a3nm.net/publications/amarilli2016tractable.pdf>.
-  Amarilli, Antoine, Silviu Maniu, Pierre Senellart (2015). “Intensional Data on the Web”. In: *SIGWEB Newsletter*. URL: <https://a3nm.net/publications/amarilli2015intensional.pdf>.
-  Amarilli, Antoine et al. (2014). “Recent Topics of Research around the YAGO Knowledge Base”. In: *Proc. APWEB*. URL: <https://zenodo.org/record/34912>.

References IV

-  Amarilli, Antoine et al. (2016). “Possible and Certain Answers for Queries over Order-Incomplete Data”. Preprint: <https://a3nm.net/publications/amarilli2016possible.pdf>.
-  Amarilli, Antoine et al. (2016). “Top- k Queries on Unknown Values under Order Constraints”. Preprint: <https://a3nm.net/publications/amarilli2016top.pdf>.
-  Chekuri, Chandra, Julia Chuzhoy (2014). “Polynomial Bounds for the Grid-Minor Theorem”. In: *Proc. STOC*.
-  Cosmadakis, Stavros S., Paris C. Kanellakis, Moshe Y. Vardi (1990). “Polynomial-Time Implication Problems for Unary Inclusion Dependencies”. In: *J. ACM*.
-  Courcelle, Bruno (1990). “The Monadic Second-Order Logic of Graphs. I. Recognizable Sets of Finite Graphs”. In: *Inf. Comput.*
-  Dalvi, Nilesh, Dan Suciu (2012). “The Dichotomy of Probabilistic Inference for Unions of Conjunctive Queries”. In: *J. ACM*.

References V

-  Dong, Xin Luna, Alon Halevy, Cong Yu (2009). “Data integration with uncertainty”. In: *The VLDB Journal—The International Journal on Very Large Data Bases*.
-  Green, Todd J., Grigoris Karvounarakis, Val Tannen (2007). “Provenance Semirings”. In: *Proc. PODS*.
-  Talaika, Aliaksandr et al. (2015). “IBEX: Harvesting Entities from the Web Using Unique Identifiers”. In: *Proc. WebDB*. URL: <http://arxiv.org/abs/1505.00841>.
-  Tang, Ruiming et al. (2014a). “A Framework for Sampling-Based XML Data Pricing”. In: *Transactions on Large-Scale Data and Knowledge-Centered Systems*. URL: <https://a3nm.net/publications/tang2014framework.pdf>.
-  Tang, Ruiming et al. (2014b). “Get a Sample for a Discount”. In: *Proc. DEXA*. URL: <https://a3nm.net/publications/tang2014get.pdf>.

References VI



Zhang, Ce (2015). “DeepDive: A Data Management System for Automatic Knowledge Base Construction”. <https://cs.stanford.edu/people/czhang/zhang.thesis.pdf>. PhD thesis. University of Wisconsin–Madison.