



UnSAID: Uncertainty and Structure in the Access to Intensional Data

Pierre Senellart





Uncertain data is everywhere

Numerous sources of **uncertain data**:

- Measurement errors
- Data integration from contradicting sources
- Imprecise mappings between heterogeneous schemas
- Imprecise automatic processes (information extraction, natural language processing, etc.)
- Imperfect human judgment
- Lies, opinions, rumors



Uncertain data is everywhere

Numerous sources of **uncertain data**:

- Measurement errors
- Data integration from contradicting sources
- Imprecise mappings between heterogeneous schemas
- Imprecise automatic processes (**information extraction**, natural language processing, etc.)
- Imperfect human judgment
- Lies, opinions, rumors



Uncertainty in Web information extraction

instance	iteration	date learned	confidence
<u>arabic, egypt</u>	406	08-sep-2011	(Seed) 100.0
<u>chinese, republic of china</u>	439	24-oct-2011	100.0
<u>chinese, singapore</u>	421	21-sep-2011	(Seed) 100.0
<u>english, britain</u>	439	24-oct-2011	100.0
<u>english, canada</u>	439	24-oct-2011	(Seed) 100.0
<u>english, england001</u>	439	24-oct-2011	100.0
<u>arabic, morocco</u>	422	23-sep-2011	100.0
<u>cantonese, hong kong</u>	406	08-sep-2011	100.0
<u>english, uk</u>	436	19-oct-2011	100.0
<u>english, south vietnam</u>	427	27-sep-2011	99.9
<u>french, morocco</u>	422	23-sep-2011	99.9
<u>greek, turkey</u>	430	07-oct-2011	99.9

Never-ending Language Learning (NELL, CMU),

<http://rtw.ml.cmu.edu/rtw/kbbrowser/>

Uncertainty in Web information extraction

Google squared labs

comedy movies

Square it Add

Item Name	Language	Director	Release Date
<input type="checkbox"/> The Mask	English	Chuck Russell	29 July 1994
<input type="checkbox"/> Scary M	<input checked="" type="radio"/> English language for the mask www.infibeam.com - all 9 sources » Other possible values	<input checked="" type="radio"/> Chuck Russell directed by for The Mask www.infibeam.com - all 9 sources » Other possible values	
<input type="checkbox"/> Superba	<input type="radio"/> English Language Low confidence language for Mask www.freebase.com	<input type="radio"/> John R. Dilworth Low confidence director for The Mask www.freebase.com	
<input type="checkbox"/> Music	<input type="radio"/> english, french Low confidence languages for the mask www.dvdreview.com	<input type="radio"/> Fiorella Infascelli Low confidence directed by for The Mask www.freebase.com - all 2 sources »	
<input type="checkbox"/> Knocked	<input type="radio"/> Italian Language Low confidence language for The Mask www.freebase.com Search for more values »	<input type="radio"/> Charles Russell Low confidence directed by for The Mask www.freebase.com - all 2 sources » Search for more values »	

Google Squared (terminated),
screenshot from (Fink, Hogue, Olteanu, and Rath 2011)



Uncertainty in Web information extraction

Subject	Predicate	Object	Confidence
Elvis Presley	diedOnDate	1977-08-16	97.91%
Elvis Presley	isMarriedTo	Priscilla Presley	97.29%
Elvis Presley	influences	Carlo Wolff	96.25%

YAGO, <http://www.mpi-inf.mpg.de/yago-naga/yago>

(Suchanek, Kasneci, and Weikum 2007)



Structured data is everywhere

Data is **structured**, not flat:

- Variety of **representation formats** of data in the wild:
 - relational tables
 - trees, semi-structured documents
 - graphs, e.g., social networks or semantic graphs
 - data streams
 - complex views aggregating individual information
- **Heterogeneous schemas**
- Additional **structural constraints**: keys, inclusion dependencies



Intensional data is everywhere

Lots of data sources can be seen as **intensional**: accessing all the data in the source (**in extension**) is **impossible** or **very costly**, but it is possible to access the data through **views**, with some **access constraints**, associated with some **access cost**.

- **Indexes** over regular data sources
- **Deep Web** sources: Web forms, Web services
- The Web or social networks as partial graphs that can be expanded by **crawling**
- Outcome of **complex automated processes**: information extraction, natural language analysis, machine learning, ontology matching
- **Crowd data**: (very) partial views of the world
- **Logical consequences** of facts, costly to compute



Interactions between uncertainty, structure, intensionality

- If the data has complex structure, uncertain models should represent **possible worlds over these structures** (e.g., probability distributions over graph completions of a known subgraph in Web crawling).
- If the data is intensional, we can use uncertainty to represent **prior distributions** about what may happen if we access the data. Sometimes good enough to reach a decision without having to make the access!
- If the data is a RDF^F graph accessed by semantic Web services, each intensional data access will **not give a single data point**, but a **complex** subgraph.



State of the art and opportunities

Probabilistic databases cover limited structure variations, do not consider intensionality

Active and reinforcement learning deals with uncertainty and intensionality, but assumes trivial structures and simple goals

Crowdsourcing, focused crawling, deep Web crawling focus on specific applications of the uncertainty/structure/intensionality problem

Answering queries using views assumes simplistic cost models

Opportunities for Web data management systems that take **all dimensions into account**



Introducing UnSAID

- Uncertainty and Structure in the Access to Intensional Data
- Jointly deal with Uncertainty, Structure, and the fact that access to data is **limited** and has a **cost**, to solve a user's **knowledge need**
- **Lazy evaluation** whenever possible
- Evolving probabilistic, structured view of the **current knowledge of the world**
- Solve at each step the problem: **What is the next best access to do** given my current knowledge of the world and the knowledge need
- **Knowledge acquisition plan** (recursive, dynamic, adaptive) that minimizes access cost, and provides probabilistic guarantees





formulation



Knowledge
need





formulation

Knowledge
need

Structured
source profiles

modeling





formulation

Knowledge
need

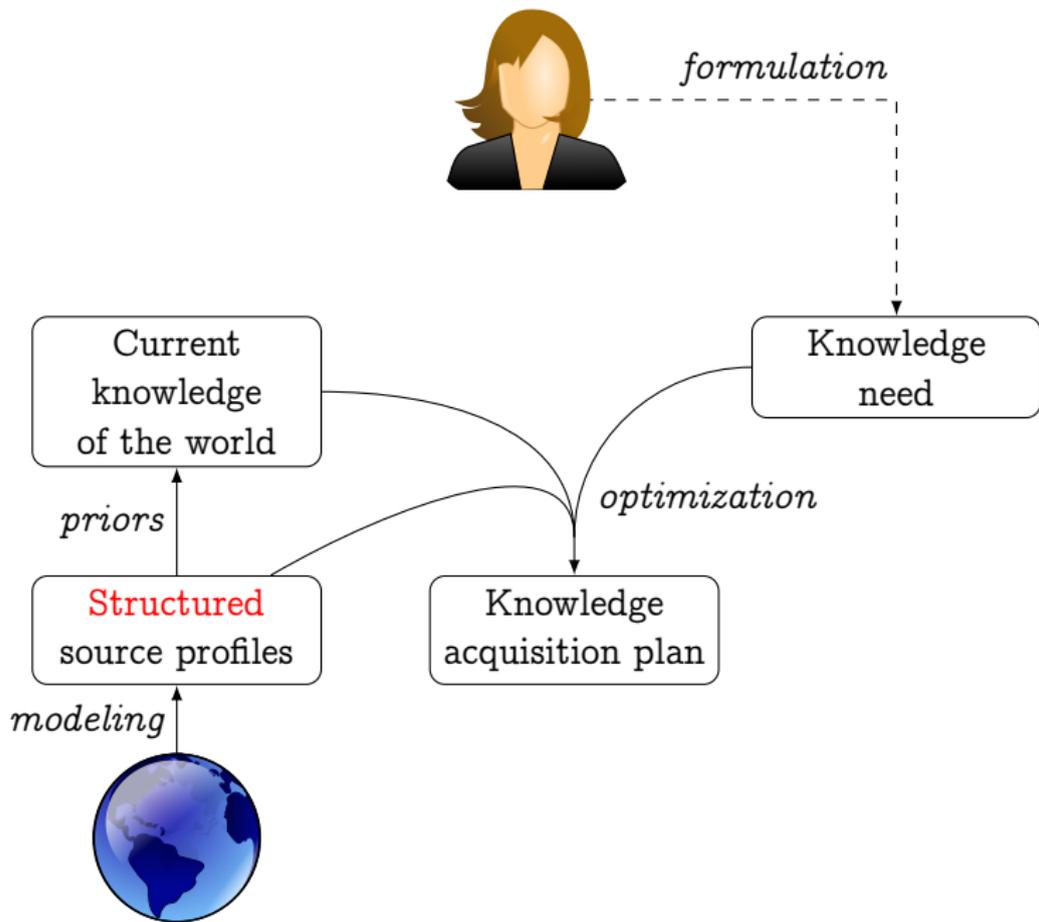
Current
knowledge
of the world

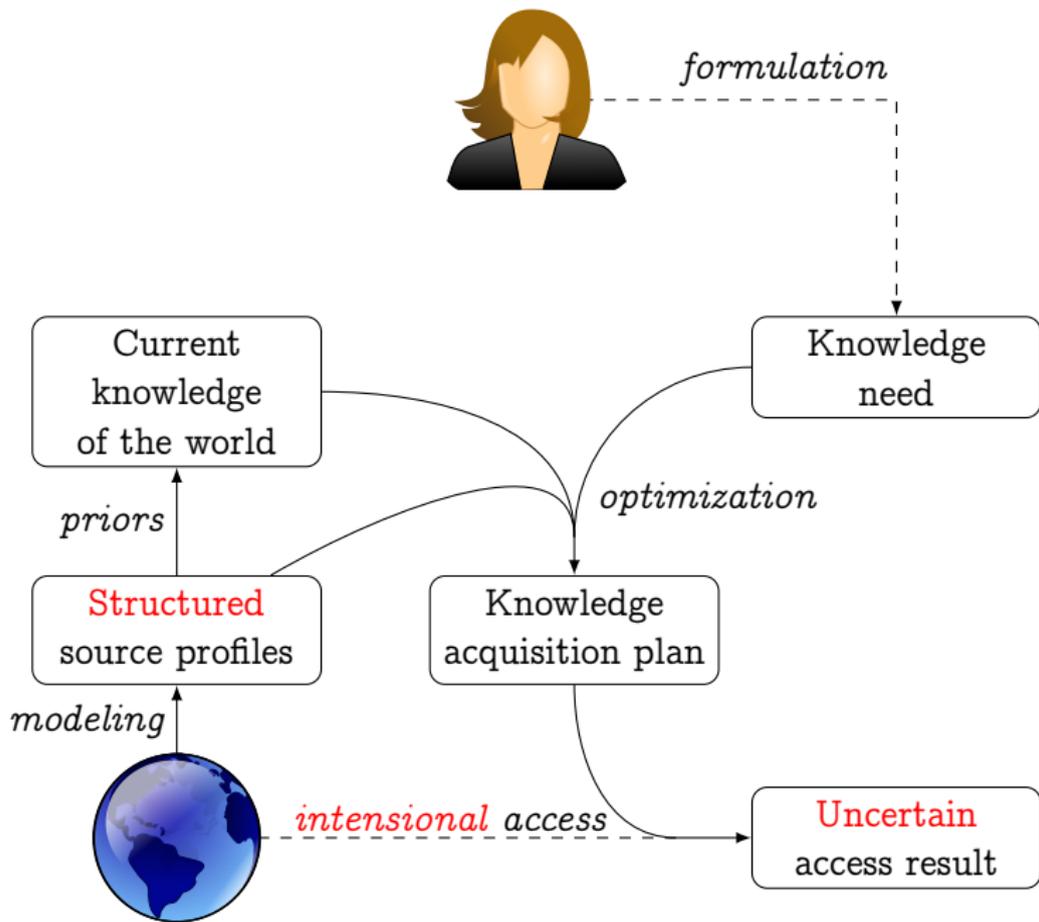
priors

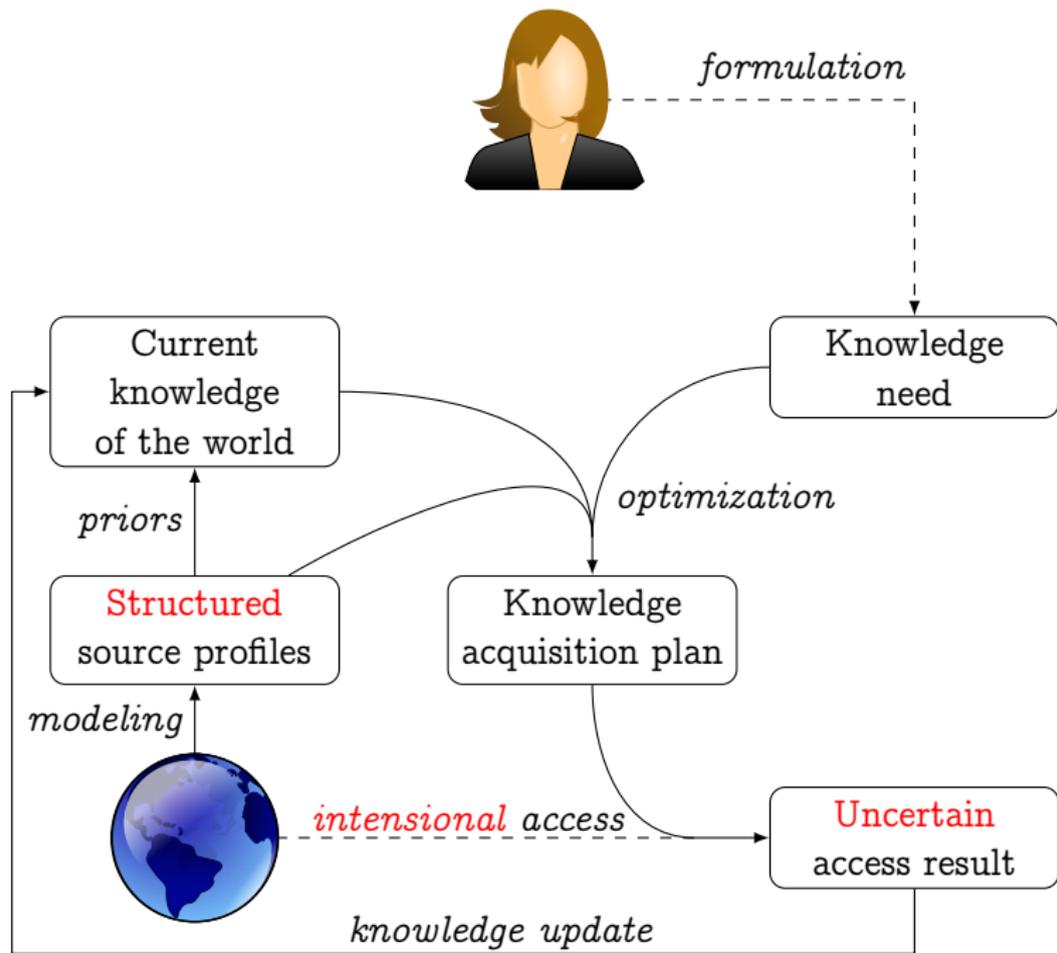
Structured
source profiles

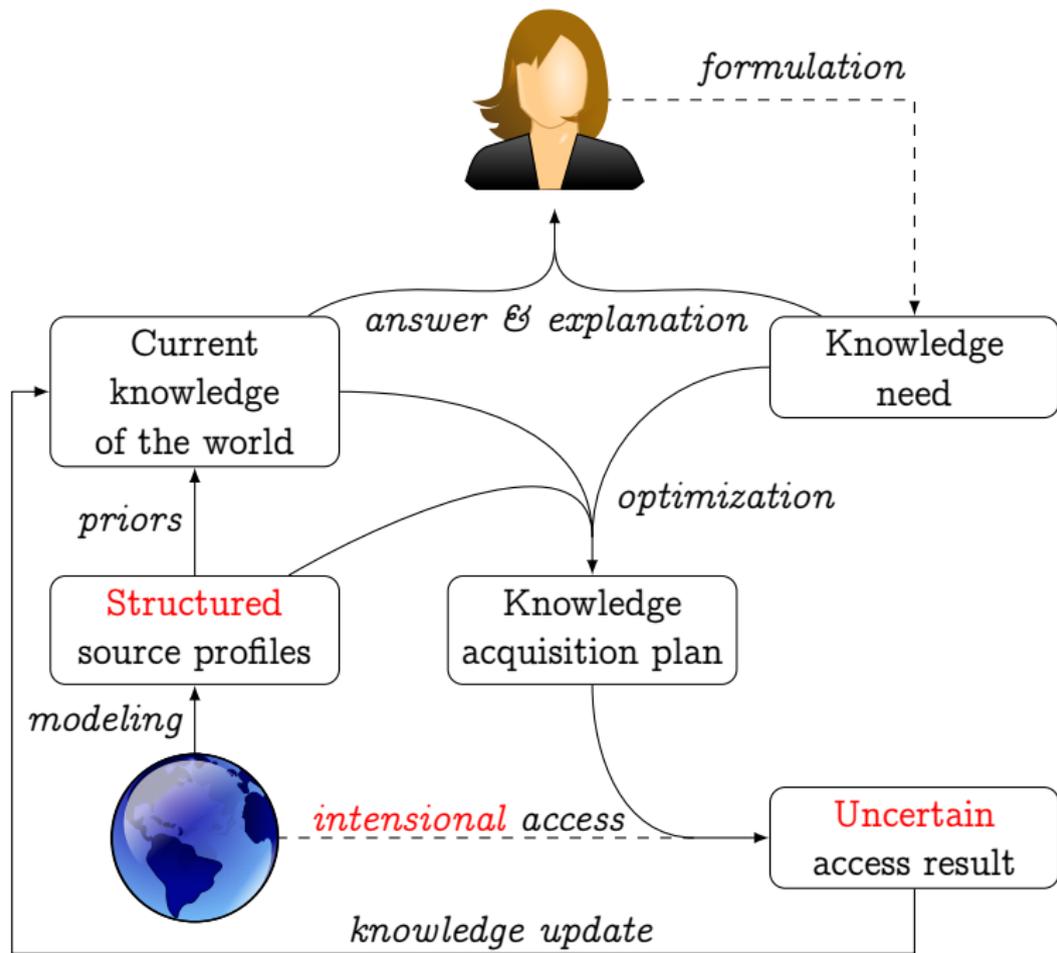
modeling













What this talk is about

- **General overview** of some of my current (and recent) research, through one-slide presentation of individual works
- Hopefully, emerging **consistent themes**
- **Connections** with the UnSAID problem



Plan

Introduction

Instances of UnSAID

Uncertainty and Structure

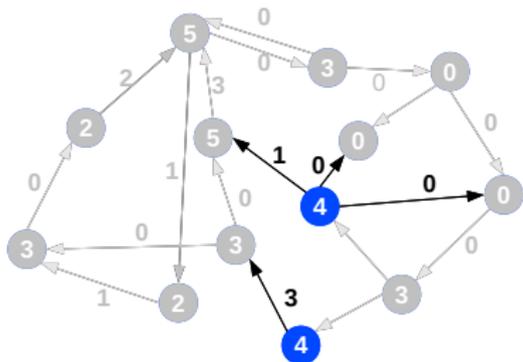
UnSAID Applications

Conclusion



Adaptive focused crawling

(Gouriten, Maniu, and Senellart 2014)



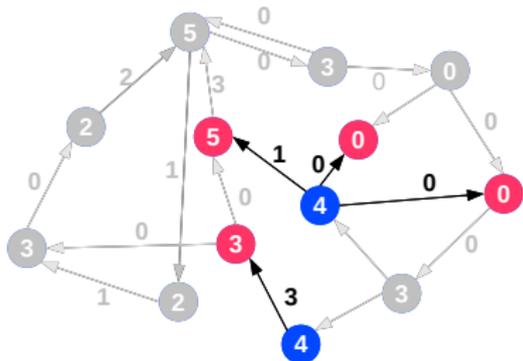
- **Problem:** Efficiently crawl nodes in a graph such that **total score is high**
- **Challenge:** The score of a node is **unknown till it is crawled**
- **Methodology:** Use various predictors of node scores, and **adaptively select the best one so far** with multi-armed bandits





Adaptive focused crawling

(Gouriten, Maniu, and Senellart 2014)



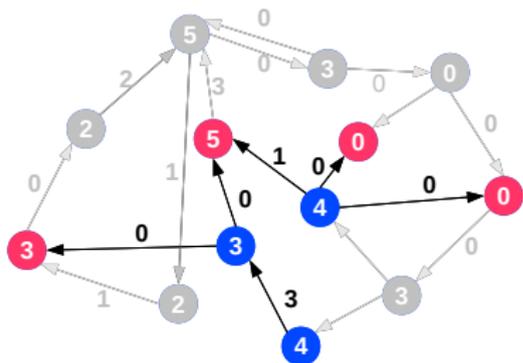
- **Problem:** Efficiently crawl nodes in a graph such that **total score is high**
- **Challenge:** The score of a node is **unknown till it is crawled**
- **Methodology:** Use various predictors of node scores, and **adaptively select the best one so far** with multi-armed bandits





Adaptive focused crawling

(Gouriten, Maniu, and Senellart 2014)

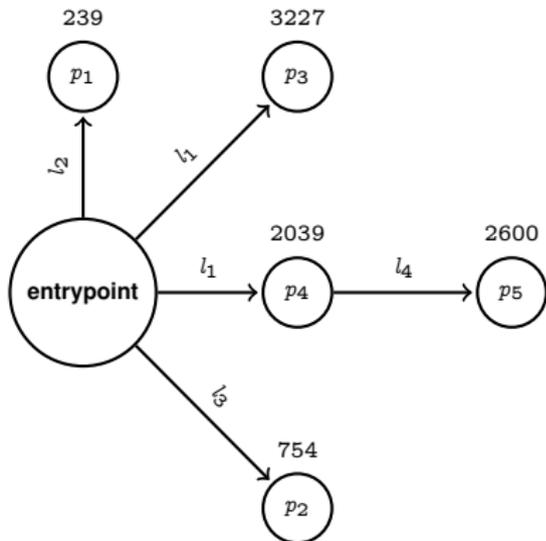


- **Problem:** Efficiently crawl nodes in a graph such that **total score is high**
- **Challenge:** The score of a node is **unknown till it is crawled**
- **Methodology:** Use various predictors of node scores, and **adaptively select the best one so far** with multi-armed bandits





Adaptive Web application crawling

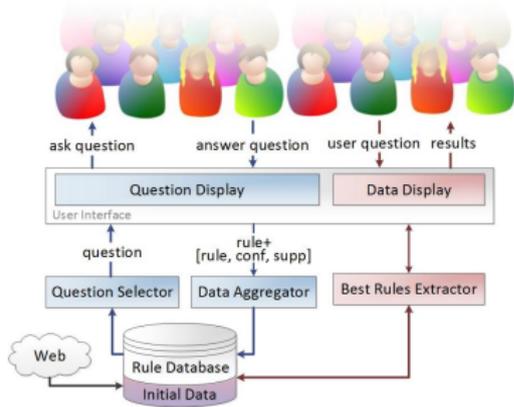


- **Problem:** Optimize the **amount of distinct content** retrieved from a Web site w.r.t. the **number of HTTP requests**
- **Challenge:** No way to know a priori **where the content lies** on the Web site
- **Methodology:** **Sample** a small part of the Web site and discover **optimal crawling patterns** from it





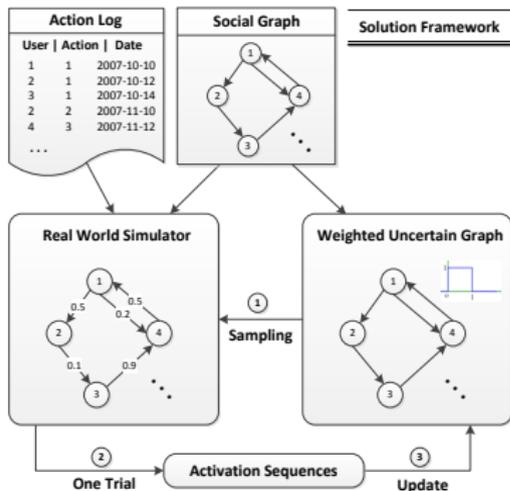
Optimizing crowd queries under order



- **Problem:** Given a query, what is the next best question to ask the crowd when crowd answers are **constrained by a partial order**
- **Challenge:** Order constraints make questions **not independent** of each other
- **Methodology:** Construct a **polytope of admissible regions** and uniformly sample from it to **determine the impact of a data item**



Online influence maximization



- **Problem:** Run **influence campaigns** in social networks, optimizing the amount of influenced nodes
- **Challenge:** Influence probabilities are **unknown**
- **Methodology:** Build a model of influence probabilities and focus on influential nodes, with an **exploration/exploitation trade-off**





Query answering under uncertain rules

$\text{Pope}(X) \Rightarrow \text{BuriedIn}(X, \text{Rome})$ (98%)
 $\text{LocatedIn}(X, \text{Lombardy}) \Rightarrow$
 $\text{BelongsTo}(X, \text{AustrianEmpire})$ (45%)

$\text{Pope}(\text{PiusXI})$
 $\text{BornIn}(\text{PiusXI}, \text{Desio})$
 $\text{LocatedIn}(\text{Desio}, \text{Lombardy})$

$\exists X, \text{BornIn}(X, Y) \wedge$
 $\text{BelongsTo}(Y, \text{AustrianEmpire}) \wedge$
 $\text{BuriedIn}(X, \text{Rome})?$

- **Problem:** Determine efficiently the probability of a query being true, given some **data and uncertain rules** over this data
- **Challenge:** Produced facts may be **correlated**, the same facts can be generated in **different ways**, probability computation is **hard in general**...
- **Methodology:** Find **restrictions** on the rules (guarded?) and the data (bounded tree-width?) that make the problem **tractable**





Plan

Introduction

Instances of UnSAID

Uncertainty and Structure

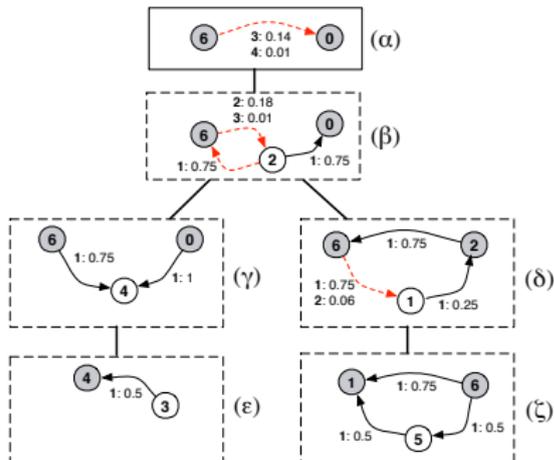
UnSAID Applications

Conclusion



Efficient querying of uncertain graphs

(Maniu, Cheng, and Senellart 2014)



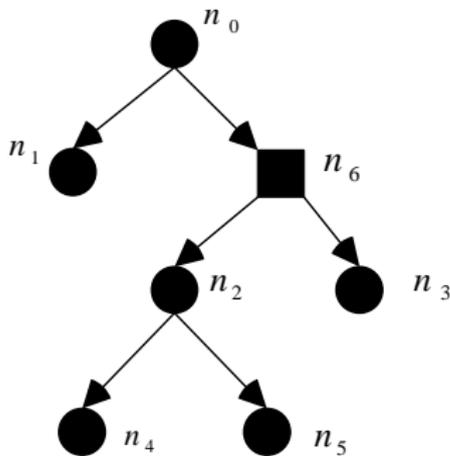
- **Problem:** Optimize query evaluation on probabilistic graphs
- **Challenge:** Probabilistic query evaluation is hard, and standard indexing techniques for large graphs do not work
- **Methodology:** Build a tree decomposition that preserves probabilities and run the query on this tree decomposition





Uniform sampling of XML documents

(Tang, Amarilli, Senellart, and Bressan 2014)

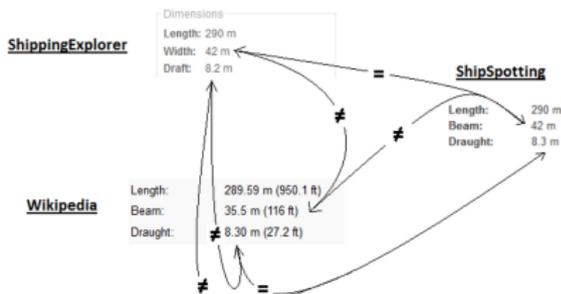


- **Problem:** Sample a subtree of fixed size/characteristics **uniformly at random** from a tree, e.g., for data pricing reasons
- **Challenge:** Naive top-down sampling does not work, will result in **biased** sampling depending on the tree structure
- **Methodology:** **Bottom-up annotation** of the tree recording distribution information, followed by **top-down sampling**

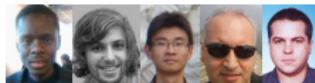




Truth discovery on heterogeneous data



- **Problem:** Determine **true values** from integrated semi-structured Web sources
- **Challenge:** Semi-structured data on the Web is **contradictory** and **copied** from source to source
- **Methodology:** Estimate the truth and determine copy patterns between sources, not only of base facts but of **subtrees** of the data





Plan

Introduction

Instances of UnSAID

Uncertainty and Structure

UnSAID Applications

Conclusion



Smarter urban mobility



- **Problem:** Smart and adaptive recommendations for mobility in cities (transit, bike rental, car, etc.)
- **Challenge:** Should take into account personal information (calendar, etc.), past trajectories, public information about transit and traffic
- **Methodology:** Map GPS tracks to routes of public transport, learn route patterns, infer destination while in transit and provide push suggestions





Introduction

Instances of UnSAID

Uncertainty and Structure

UnSAID Applications

Conclusion



What's next?

- So far, we have tackled **individual** aspects or specializations of the UnSAID problem
- Now we need to consider the general problem, and propose **general solutions**
- There is a strong potential for uncovering the **unsaid information from the Web**
- Strong **connections** with a number of research areas: active learning, reinforcement learning, adaptive query evaluation, etc. Inspiration to get from these areas.
- **Everyone is welcome** to join the effort!



What's next?

- So far, we have tackled **individual** aspects or specializations of the UnSAID problem
- Now we need to consider the general problem, and propose **general solutions**
- There is a strong potential for uncovering the **unsaid information from the Web**
- Strong **connections** with a number of research areas: active learning, reinforcement learning, adaptive query evaluation, etc. Inspiration to get from these areas.
- **Everyone is welcome** to join the effort!

Merci.

Amarilli, A., M. L. Ba, D. Deutch, and P. Senellart (Dec. 2013). *Provenance for Nondeterministic Order-Aware Queries*. Preprint available at

<http://pierre.senellart.com/publications/amarilli2014provenance.pdf>.

Ba, M. L., S. Montenez, T. Abdessalem, and P. Senellart (Apr. 2014). *Extracting, Integrating, and Visualizing Uncertain Web Information about Moving Objects*. Preprint available at

<http://pierre.senellart.com/publications/ba2014extracting.pdf>.

Fink, R., A. Hogue, D. Olteanu, and S. Rath (2011). “SPROUT²: a squared query engine for uncertain web data”. In: *SIGMOD*.

Gouriten, G., S. Maniu, and P. Senellart (Mar. 2014). *Scalable, Generic, and Adaptive Systems for Focused Crawling*. Preprint available at

<http://pierre.senellart.com/publications/gouriten2014scalable.pdf>.

Maniu, S., R. Cheng, and P. Senellart (Mar. 2014). *ProbTree: A Query-Efficient Representation of Probabilistic Graphs*. Preprint available at

<http://pierre.senellart.com/publications/maniu2014probtree.pdf>.

Suchanek, F. M., G. Kasneci, and G. Weikum (2007). “YAGO: A Core of Semantic Knowledge. Unifying WordNet and Wikipedia”. In: *WWW*, pp. 697–706. ISBN: 978-1-59593-654-7.

Tang, R., A. Amarilli, P. Senellart, and S. Bressan (Mar. 2014). *Get a Sample for a Discount: Sampling-Based XML Data Pricing*. Preprint available at

<http://pierre.senellart.com/publications/tang2014get.pdf>.