# Semantic Data Enrichment meets Neural-Symbolic Integration

## Vincenzo Cutrona

University of Milano - Bicocca
vincenzo.cutrona@unimib.it

City, University of London
Vincenzo.Cutrona@city.ac.uk

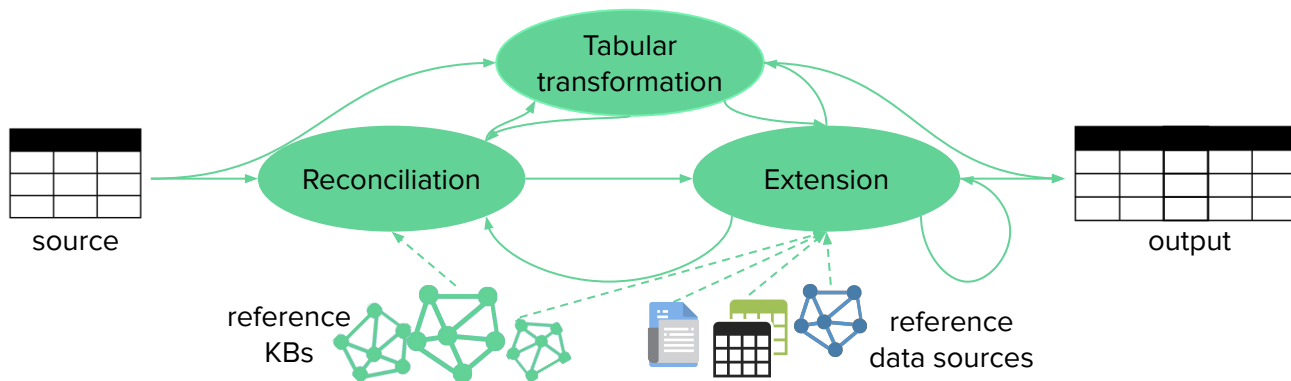# Problem Statement: Semantic Data Enrichment

Inputs:
- a **source** dataset
- a pool of **reference data sources**

Output:
- source dataset extended with more columns

**Semantic data enrichment:** a walk on the **data transformations** graph $G^T$ where at least one node is a **reconciliation**

$$t_1(s_0), ...t_n(s_{n-1}) \text{ where } s_0 = \text{source dataset and } t_n(s_{n-1}) = \text{target dataset}$$
$$t_i = \text{tabular transformation} \mid \text{reconciliation} \mid \text{extension}$$

# Problem Statement: Semantic Data Enrichment

Inputs:
- a **source** dataset
- a pool of **reference data sources**

● **large-scale** tabular data
➔ big data processing

Output:
- source dataset extended with more columns

**Semantic data enrichment:** a walk ... graph $G^T$ where at least one node is a **reconciliation**

● against **large-scale** KBs
➔ fast execution

$$t_1(s_0), ... t_n(s_{n-1}) \text{ where } s_0 = \text{source dataset and } t_n(s_{n-1}) = \text{target dataset}$$
$$t_i = \text{tabular transformation} \mid \text{reconciliation} \mid \text{extension}$$

Tabular transformation

Reconciliation

Extension

source

output

reference KBs

reference data sources

# Relevancy: Enrichment for Data Analytics

- Useful in a variety of data science applications based on the analysis
    1) **Weather-aware scheduler for digital marketing campaigns**

# Relevancy: Enrichment for Data Analytics

- Useful in a variety of data science applications based on the analysis
  2) **Event-aware scheduler for digital marketing campaigns**
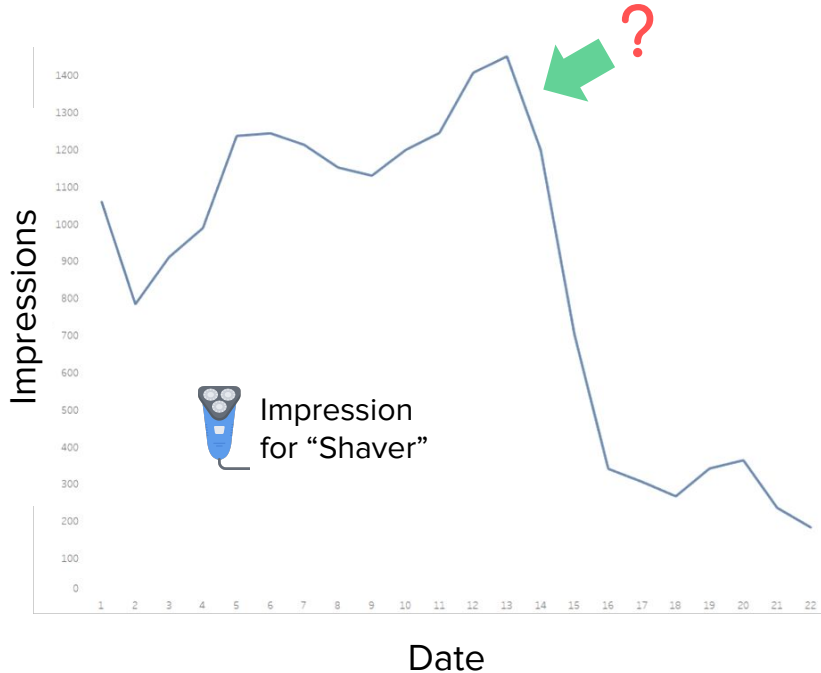
# Relevancy: Enrichment for Data Analytics

- Useful in a variety of data science applications based on the analysis
  2) **Event-aware scheduler for digital marketing campaigns**



Impression for "Shaver"

Impressions

Date

Explanation 1
Men stop shaving

(Photo: Castaway)

# Relevancy: Enrichment for Data Analytics

- Useful in a variety of data science applications based on the analysis
  2) **Event-aware scheduler for digital marketing campaigns**



Explanation 1
Men stop shaving

# Relevancy: Enrichment for Data Analytics

- Useful in a variety of data science applications based on the analysis
  2) **Event-aware scheduler for digital marketing campaigns**



Valentine's Day

Impression for "Shaver"

Impressions

Date

Explanation 1
Men stop shaving

(Photo: Castaway)

Explanation 2
Women interest for shavers starts decreasing the day before Valentine's Day
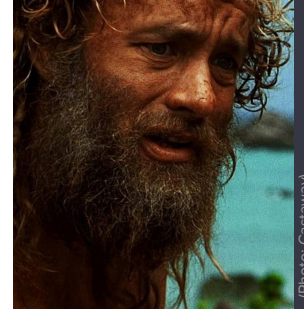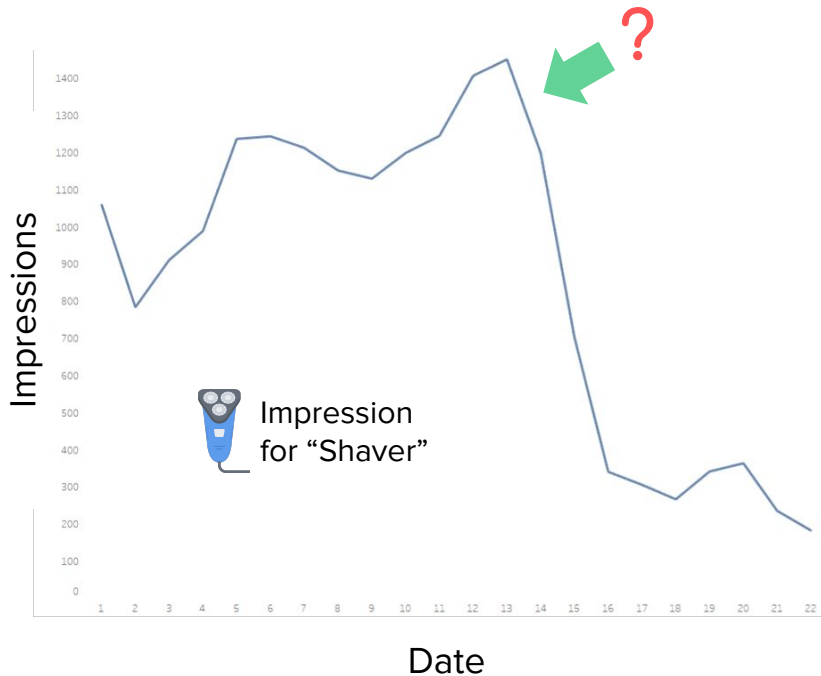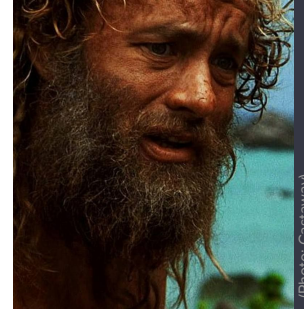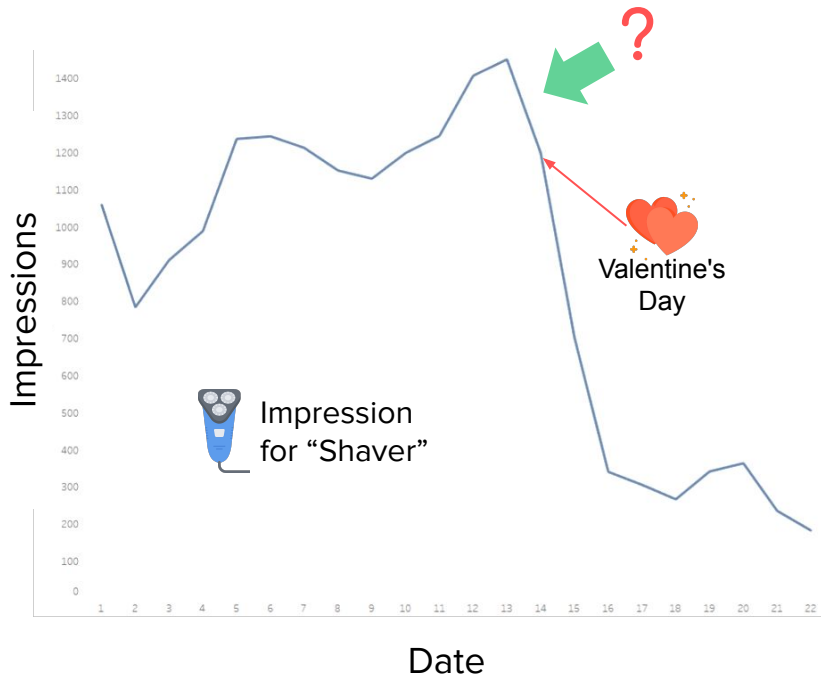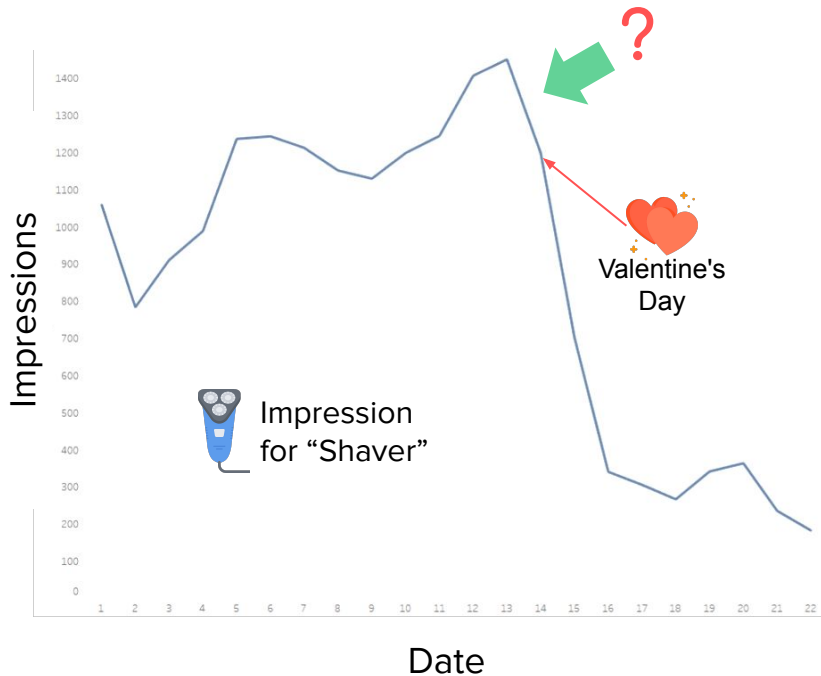
www.spreadshirt.com

# Relevancy: Enrichment for Data Analytics

- Useful in a variety of data science applications based on the analysis
  2) **Event-aware scheduler for digital marketing campaigns**



Figure 12: Correlation between German national handball team matches and online digital indicators in the "SportFitness" category.

# An Example of                    Data Enrichment

Heterogeneous data (different systems of identifiers)

**Google AdWords**

| KEYWORD | CITY | REGION | Clicks | Category | Date |
|---------|------|--------|--------|----------|------|
| 517827 | Ingolstadt | Bavaria | 50 | NewsMediaPublications | 12/03/2018 |
| 459143 | Berlin | Berlin | 42 | TravelTourism | 12/03/2018 |
| 891139 | Munich | Bavaria | 36 | HomeGarden | 11/03/2018 |

**dmoz**

Events for Recreation → Travel

**GeoNames**

Intuitively: more data attached to input data

Basically, a **LEFT OUTER JOIN** between datasets...

EVENT
**PM Modi to meet German Chancellor Merkel on April 20**

WHEN
Fri, April 20, 2018

WHERE
Berlin, Germany

ARTICLES
11

VIRALITY
5

👤 Angela Merkel   👤 Narendra Modi   ⚑ Chancellor of Germany   📍 Berlin, Germany   📍 Sweden
📍 United Kingdom   📍 New Delhi, India   ⚑ Commonwealth Heads of Government Meeting

Prime Minister Narendra Modi will meet German Chancellor Angela Merkel during a stopover in Berlin on April 20 after concluding his visits to Sweden a... (Hindustan Times)

EVENT
**Finland Named Partner Destination of ITB China 2018**

WHEN
Mon, March 12, 2018

WHERE
Berlin, Germany

ARTICLES
1

VIRALITY
15

📍 Tourism   ⚑ ITB Berlin   📍 Malaysia   📍 Berlin, Germany   📍 Germany
🏠 Bandung Institute of Technology   📍 China   📍 Tourism Malaysia   ⚑ Trade fair

Finland has been named as the Partner Destination of ITB China, scheduled to take place from 16 to 18 May 2018. Around 500 ...

**EVENTREGISTRY**

10

# An Example of        Data Enrichment



| KEYWORD | CITY | REGION | Clicks | Category | Date |
|---------|------|--------|--------|----------|------|
| 517827 | Ingolstadt | Bavaria | 50 | NewsMediaPublications | 12/03/2018 |
| 459143 | Berlin | Berlin | 42 | TravelTourism | 12/03/2018 |
| 891139 | Munich | Bavaria | 36 | HomeGarden | 11/03/2018 |

**PROBLEM**

No common identifiers ➜ no JOIN

Events for Recreation → Travel

**EVENT**
PM Modi to meet German Chancellor Merkel on April 20

| WHEN | WHERE | ARTICLES | VIRALITY |
|------|-------|----------|----------|
| Fri, April 20, 2018 | Berlin, Germany | 11 | 5 |

Angela Merkel    Narendra Modi    Chancellor of Germany    Berlin, Germany    Sweden
United Kingdom    New Delhi, India    Commonwealth Heads of Government Meeting

Prime Minister Narendra Modi will meet German Chancellor Angela Merkel during a stopover in Berlin on April 20 after concluding his visits to Sweden a... (Hindustan Times)

**EVENT**
Finland Named Partner Destination of ITB China 2018

| WHEN | WHERE | ARTICLES | VIRALITY |
|------|-------|----------|----------|
| Mon, March 12, 2018 | Berlin, Germany | 1 | 15 |

Tourism    ITB Berlin    Malaysia    Berlin, Germany    Germany
Bandung Institute of Technology    China    Tourism Malaysia    Trade fair

Finland has been named as the Partner Destination of ITB China, scheduled to take place from 16 to 18 May 2018. Around 500 exhibitors from...

# An Example of **Semantic** Data Enrichment



| KEYWORD | CITY | REGION | Clicks | Category | Date |
|---------|------|--------|--------|----------|------|
| 517827 | Ingolstadt | Bavaria | 50 | NewsMediaPublications | 12/03/2018 |
| 459143 | Berlin | Berlin | 42 | TravelTourism | 12/03/2018 |
| 891139 | Munich | Bavaria | 36 | HomeGarden | 11/03/2018 |

Semantic Web of Data to support the enrichment task



**WEB OF DATA**

# An Example of **Semantic** Data Enrichment



Semantic Web of Data to support the enrichment task

| KEYWORD | CITY | REGION | Clicks | Category | Date |
|---------|------|--------|--------|----------|------|
| 517827 | Ingolstadt | Bavaria | 50 | NewsMediaPublications | 12/03/2018 |
| 459143 | Berlin | Berlin | 42 | TravelTourism | 12/03/2018 |
| 891139 | Munich | Bavaria | 36 | HomeGarden | 11/03/2018 |

**SEMANTIC ENRICHMENT**
**=**
**RECONCILIATION**
**+**
**EXTENSION**

Target KB

**WEB OF DATA**

13

# Fundamentals

# Knowledge Graphs

# Knowledge Graph: Instances



(subject, predicate, object)

John Keats —death place→ Rome —country→ Italy

Rome —capital→ Italy

# Knowledge Graph: Types

# Knowledge Graph: FOL View vs. Graph View

Predicate(Subject,Object)



predicate

**Subject** → **Object**

<Subject,predicate,Object>

**FOL VIEW**

**GRAPH VIEW**

Country(Rome,Italy)



country

**Rome** → **Italy**

<Rome,country,Italy>

**FOL VIEW**

**GRAPH VIEW**

# Semantics in KGs

Schema (≈ ontology) defines the meaning of general terms in the KG

- Types, e.g., City(x)
- Relations, e.g., country(x,y)

Schema definition supports inference (by deduction, or by induction, etc.)

- E.g., $\forall$x,y country(x,y) $\Rightarrow$ (City(x) $\land$ Country(y))
- E.g., $\forall$x,y capital(x,y) $\Rightarrow$ country(x,y)

# Reconciliation

# Reconciliation



Cleaned dataset

3rd-party datasets

Value Reconciliation

Reconciled Dataset

Given a text that is a name or label for something, returns a ranked list of **potential entities** (based on some matching rules)

**Matching rules** assign a score to the entities (e.g., distance)

Fully-automated approaches must also **select the best candidate** to return

apple →

*Apple (0.0)* 🍎

*Apple, Inc. (0.20)* 

*The Big Apple (0.30)* 

*Apple (album) (0.24)* 💿

**< .25**

*Apple (0.0)* 🍎

*Apple, Inc. (0.20)* 

*The Big Apple (0.30)* 

*Apple (album) (0.24)* 💿

**min**

*Apple (0.0)* 🍎

*Apple, Inc. (0.20)* 

*The Big Apple (0.30)* 

*Apple (album) (0.24)* 💿

Input text          Distance score computation                    Candidates Selection                          Decision-Making

# Reconciliation: Main issues

- Precision depends on the **ambiguities**
- Impossible to explore the whole **candidate entities space**
- Human experts can not **check all results**

apple →

*Apple (0.0)* 🍎
*Apple, Inc. (0.20)* 🍎
*The Big Apple (0.30)* 🗽
*Apple (album) (0.24)* 💿

**< .25** →

*Apple (0.0)* 🍎
*Apple, Inc. (0.20)* 🍎
*The Big Apple (0.30)* 🗽
*Apple (album) (0.24)* 💿

**min** →

*Apple (0.0)* 🍎
*Apple, Inc. (0.20)* 🍎
*The Big Apple (0.30)* 🗽
*Apple (album) (0.24)* 💿

Input text       Distance score computation                Candidates Selection                Decision-Making

# Reconciliation in Tables

| KEYWORD | REGION | Clicks | Category | Date |
|---------|--------|--------|----------|------|
| 194906 | Thuringia | 64 | BusinessManagement | 11/03/2018 |
| 517827 | Bavaria | 50 | NewsMediaPublications | 12/03/2018 |
| 459143 | Berlin | 42 | TravelTourism | 12/03/2018 |
| 891139 | Bavaria | 36 | Vehicles | 11/03/2018 |
| 459143 | Bavaria | 30 | HomeGarden | 10/03/2018 |

Reconcile region names versus Geonames identifiers (11.7M entities)

Reconcile category names versus DMOZ taxonomy (1M entities)

| KEYWORD | REGION | Geonames ID | Clicks | Category | DMOZ ID | Date |
|---------|--------|-------------|--------|----------|---------|------|
| 194906 | Thuringia | 2822542 | 64 | BusinessManagement | dmoz/Business/Management | 11/03/2018 |
| 517827 | Bavaria | 2951839 | 50 | NewsMediaPublications | dmoz/News | 12/03/2018 |
| 459143 | Berlin | 2950157 | 42 | TravelTourism | dmoz/Recreation/Travel | 12/03/2018 |
| 891139 | Bavaria | 2951839 | 36 | Vehicles | dmoz/Shopping/Vehicles | 11/03/2018 |
| 459143 | Bavaria | 2951839 | 30 | HomeGarden | dmoz/Home/Gardening | 10/03/2018 |

# Reconciliation in Tables - Cell-by-Cell

*Altenburg (0.0) (city)* 🇺🇸

*Altenburg (0.0) (mountain)* 🏔️

*Altenburg (0.0) (city)* 🇩🇪

*Altenburg (0.0) (city)* 🇺🇸

| CITY | REGION |
|------|--------|
| Altenburg | Thuringia |
| Ingolstadt | Bavaria |
| Berlin | Berlin |

# Reconciliation in Tables - Cell-by-Cell

| CITY | REGION |
|------|--------|
| Altenburg | Thuringia |
| Ingolstadt | Bavaria |
| Berlin | Berlin |

Ingolstadt (0.0) (city)

Altenburg (0.0) (city)

Ingolstadt (0.0) (city)

# Reconciliation in Tables - Cell-by-Cell

| CITY | REGION |
|------|--------|
| Altenburg | Thuringia |
| Ingolstadt | Bavaria |
| Berlin | Berlin |

*Altenburg (0.0) (city)*

*Ingolstadt (0.0) (city)*

*Berlin (0.0) (region)*

*Berlin (0.0) (region)*

*Berlin (0.0) (city)*

# Reconciliation in Tables - Cell-by-Cell

| CITY | REGION |
|------|--------|
| Altenburg | Thuringia |
| Ingolstadt | Bavaria |
| Berlin | Berlin |

*Altenburg (0.0) (city)*

*Ingolstadt (0.0) (city)*

*Berlin (0.0) (region)*

**different types (city/region)
AND
different c (dataset
contains German cities only)**

- We did not exploit the tabular structure!
- Cells in the same column talk about the same things
  - Not always true! Sometimes data are very noisy…

# Reconciliation in Tables - Column-by-Column



- By looking at the columns, we are focusing on **CATEGORIES**
- We have to identify which is the **category that has at least one candidate in each subgroup**
- How many categories exist? cities, cities in Europe, cities in Italy …
  - ~$2^{(m \cdot n)}$, where m = #attributes and n = #possible values for each attribute

# Reconciliation in Tables - Row-by-Row

# Reconciliation in Tables - Row-by-Row

property available
in the KG

Bavaria (region)

inRegion

| CITY | REGION |
|------|--------|
| Altenburg | Thuringia |
| Ingolstadt | **Bavaria** |
| Berlin | Berlin |

Altenburg (0.0) (mountain)

Ingolstadt (0.0) (city)

Ingolstadt (0.0) (city)

# Reconciliation in Tables - Row-by-Row

# Reconciliation in Tables - Row-by-Row

| CITY | REGION |
|------|--------|
| Altenburg | Thuringia |
| Ingolstadt | Bavaria |
| Berlin | Berlin |

*Altenburg (0.0) (mountain)* 🔺

*Ingolstadt (0.0) (city)* 🇩🇪

*Berlin (0.0) (city)* 🇩🇪

**different types (city/region)
BUT
all entities have the right value
for the *inRegion* property**

- By looking at the rows, we are focusing on **PROPERTIES**
- We have to identify which are the **most discriminative properties to consider**
- How many properties to compare for each row?
  - ~$(m \cdot n)$, where m = #attributes and n = #candidates

# Logic Tensor Networks

# Terminological Recap

A **constant** is an element of a domain (set) taken in consideration

*S : {Rome, Paris, ...}*
*T : {Italy, France, ...}*

A **function** is a relation f: S → T between sets that associates to every element of a first set exactly one element of the second set.

*Capital*: T → S
Capital(Italy) = Rome

A **predicate** is a Boolean-valued function P: S → {1 (= True), 0 (=False)}.

*city: S → {1, 0}*          *country: S x T → {1, 0}*
*city(Rome) = 1*          *country(Rome, Italy) = 1*

# Terminological Recap (cont)

An **axiom:** a statement in a logical language:

$$R(a, b)$$

A **grounded axiom** contains grounded constants:

$$country(Rome, Italy)$$

A **quantified axiom** is an axiom that contains quantified variables:

$$\forall \; x,y \; capital(x, y)$$

A **formula** is a combination of grounded and quantified axioms:

$$\forall \; x,y \; country(Rome, Italy) \; \& \; country(Paris, France) \; \& \; capital(x, y)$$

# Logic Tensor Networks

Logic Tensor Networks [Serafini+,2016] (LTNs) => neuro-symbolic [Garcez+,2008;Garcez+,2012] combines neural network and symbolic AI.

**LTNs** = **Neural Networks** + **First Order Fuzzy Logic**

**Key Aspects**:

- LTNs **ground fuzzy logic in a vector space: continuous values in [0,1]**
- LTNs assign truth values to formulas using neural networks
- LTNs can learn from both data and rules
- LTNs can be used to do inferences over rules after training

**Key Idea**: LTNs provide a method to learn reasoning over vector spaces

# Logic Tensor Networks

Logic Tensor Networks [Serafini+,2016] (LTNs) => neuro-symbolic [Garcez+,2008;Garcez+,2012] combines neural network and symbolic AI.

**LTNs** = **Neural Networks** + **First Order Fuzzy Logic**

**Key Aspects**:

- LTNs ground fuzzy logic in a vector space: continuous values in [0,1]
- LTNs **assign truth values to formulas using neural networks**
- LTNs can learn from both data and rules
- LTNs can be used to do inferences over rules after training

**Key Idea**: LTNs provide a method to learn reasoning over vector spaces

# Logic Tensor Networks

Logic Tensor Networks [Serafini+,2016] (LTNs) => neuro-symbolic [Garcez+,2008;Garcez+,2012] combines neural network and symbolic AI.

**LTNs** = **Neural Networks** + **First Order Fuzzy Logic**

**Key Aspects**:

- LTNs ground fuzzy logic in a vector space: continuous values in [0,1]
- LTNs assign truth values to formulas using neural networks
- LTNs **can learn from both data and rules**
- LTNs can be used to do inferences over rules after training

**Key Idea**: LTNs provide a method to learn reasoning over vector spaces

# Logic Tensor Networks

Logic Tensor Networks [Serafini+,2016]  (LTNs) => neuro-symbolic [Garcez+,2008;Garcez+,2012]  combines neural network and symbolic AI.

**LTNs** = **Neural Networks** + **First Order Fuzzy Logic**

**Key Aspects**:

- LTNs ground fuzzy logic in a vector space: continuous values in [0,1]
- LTNs assign truth values to formulas using neural networks
- LTNs can learn from both data and rules
- LTNs **can be used to do inferences over rules after training**

**Key Idea**: LTNs provide a method to learn reasoning over vector spaces

# Logic Tensor Networks: General Idea

**parent**(Susan, Ann)

# Logic Tensor Networks: General Idea

**parent**(Susan, Ann)

Constants are points in **R**$^k$

# Logic Tensor Networks: General Idea

**parent**(Susan, Ann)

Neural Network

Each predicate in LTN is a NN (the training phase is on many networks as predicates)

# Logic Tensor Networks: General Idea

**parent**(Susan, Ann)

**Neural Network** ( [ ] , [ ] )

# Logic Tensor Networks: General Idea

**parent**(Susan, Ann)

NN + sigmoid

Neural Network $\left( \begin{array}{c} \phantom{x} \end{array}, \begin{array}{c} \phantom{x} \end{array} \right) = [0,1]$

# Logic Tensor Networks: General Idea

**parent**(Susan, Ann) & **parent(**Mike, Robert)

# Logic Tensor Networks: General Idea

**parent**(Susan, Ann) & **parent(**Mike, Robert)

Neural Network $\left(\begin{array}{c}\\\end{array}, \begin{array}{c}\\\end{array}\right) = x$

Neural Network $\left(\begin{array}{c}\\\end{array}, \begin{array}{c}\\\end{array}\right) = y$

# Logic Tensor Networks: General Idea

**parent**(Susan, Ann) & **parent(**Mike, Robert)

Neural Network $( \ , \ ) = x$

Neural Network $( \ , \ ) = y$

min(x,y)

t-norm

# Logic Tensor Networks: General Idea

**parent**(Susan, Ann) & **parent(**Mike, Robert)

Neural Network ( [ ] , [ ] ) = x      Neural Network ( [ ] , [ ] ) = y

**min(x,y)**

t-norm

How do we learn these representations?
**Backpropagation**

# Example

**KB:**

- ¬parent(mark, john)
- parent(john,mark)
- ancestor(mark, lucas)
- parent(john, susan) | parent(john, dania)

**Parent**

**Ancestor**

J  M  L  S  D

# Example: Forward Pass

- **¬parent(mark, john)**
- parent(john,mark)
- ancestor(mark, lucas)
- parent(john, susan) | parent(john, dania)

$$1 - \text{Parent}(\underset{M}{\square}, \underset{J}{\square}) = 0.8$$

**We want to maximize this**

Parent

Ancestor

J  M  L  S  D

# Example: Back Pass

- **¬parent(mark, john)**
- parent(john,mark)
- ancestor(mark, lucas)
- parent(john, susan) | parent(john, dania)

$$1 - \text{Parent}\left(\boxed{\phantom{M}}_M, \boxed{\phantom{J}}_J\right) = 0.8$$

**Update using backpropagation**



51

# Example: Forward and Back Pass

- ¬parent(mark, john)
- **parent(john,mark)**
- ancestor(mark, lucas)
- parent(john, susan) | parent(john, dania)

$$\text{Parent}\left(\,\Box_J\,,\,\Box_M\,\right) = 0.8$$

**We want to maximize this and thus we update the respective values**

Parent

Ancestor

J M L S D

52

# Example: Forward and Back Pass

Ancestor $\left(\begin{array}{c}\\\end{array}, \begin{array}{c}\\\end{array}\right) = 0.7$

M   L

**We want to maximize this and thus we update the respective values**

Parent

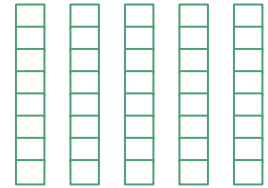Ancestor

J   M   L   S   D    53

# Example: Forward and Back Pass

KB:

- ¬parent(mark, john)
- parent(john,mark)
- ancestor(mark, lucas)
- **parent(john, susan) |
  parent(john, dania)**

$$\text{Parent} \left( \Box , \Box \right) = 0.2 \mid \text{Parent} \left( \Box , \Box \right) = 0.9$$

J   S       J   D

$$\max(0.2, 0.9) = 0.9$$

**We want to maximize this and thus we update the respective values**

Parent

Ancestor

J   M   L   S   D    54

# Logic Tensor Networks: Learning

The network is trained on a **best satisfiability task**:

- Learn the representations
  - **vectors** for the constants
  - **parameters** for the predicates

  in such a way that **the axioms are satisfied in the best possible way**.

Given **parent(Ann, Susan)** we expect the network to **learn representations** for **Ann**, **Susan** and **parent** in such a way that the predicted value is close to 1

# Implementing Logic in Tensor Networks

The grounding of *m*-ary predicate P, **G(P)**, is defined as a function from $R^{mn}$ to [0,1]

$$\mathcal{G}(P) = \sigma \left( u_P^T \tanh \left( \mathbf{v}^T W_P^{[1:k]} \mathbf{v} + V_P \mathbf{v} + B_P \right) \right)$$  [Serafini+,2016]



(Image adapted from [Socher+,2013])

# Implementing Logic in Tensor Networks: an example



Tensor net for P(x, y) ➡ A(y), with G(x) = v and G(y) =u and k= 2 (from [Serafini+,2016])

# Logic Tensor Networks: Data and Rules

LTNs can learn from both data and rules.

Quantifiers are defined **over a domain sample.**

**parent(Mark,Susan)**

**parent(Ron,Susan)**

**∀ x,y parent(x,y) →**
**ancestor(x,y)**

Optimize the
**representations** of
the parameters to
support the axioms

Quantifiers interpreted using an **aggregation function** (e.g., average):

**∀ x P(x)** = average value of P(x) in LTNs.

# Logic Tensor Networks: After Training Inference

The trained network defines a new **compositional language** built on constants, functions and predicates, which can be combined arbitrary.

The trained network can be used for discovering novel inferences.

Suppose we train using a dataset of *parents* and *ancestors* relationships.

# Logic Tensor Networks: After Training Inference

The trained network defines a new **compositional language** built on constants, functions and predicates, which can be combined arbitrary.

The trained network can be used for discovering novel inferences.

Suppose we train using a dataset of **parents** and **ancestors** relationships.

**After training** we can query LTNs on:

**∀ x,y ancestor(x,y) ➜ parent(x,y)** has truth value close to 0

# Reconcile tables with LTNs

# LTN-based Reconciliation

**①** Embed the KB in a vector space KGE

- Each entity in the graph is mapped to a $n$-dimensional point in $R^n$
  - e.g., by Graph Embedding [Wang+,2017]

DBpedia

Rome$_{Italy}$
Italy
Paris$_{France}$
France
…

$v(Rome_{Italy}) = .5\ .1\ .4\ …\ .85$
$v(Italy) = .5\ .1\ .4\ …\ .85$
$v(Paris_{France}) = .15\ .31\ .44\ …\ .5$
$v(France) = .3\ .11\ .14\ …\ .95$
…

KGE

# LTN-based Reconciliation

Get axioms from the KB ontology

$\forall$x City(x) $\rightarrow$ $\exists$y: country(x, y)     (A city must be in a country)
$\forall$x Country(x) $\rightarrow$ $\exists$y: capital(y, x)   (A country must have a capital)
$\forall$x,y capital(x, y) $\rightarrow$ country(x, y)     (A capital must be a city of its country)
$\forall$x ¬country(x,x)                  (The countryOf property is non-reflexive)
...

country($Rome_{Italy}$, Italy) = 1     (Rome is located in Italy)
country($Paris_{France}$, Italy) = 0    (Paris is not located in Italy)
City($Rome_{Italy}$) = 1           (Rome is a city)
City(Italy) = 0                 (Italy is not a city)
...

DBpedia

Axioms

63

# LTN-based Reconciliation

**③** Train the LTN with axioms and KGE, and obtained the trained model (which represents a new language!)

Axioms

KGE

Constants replaced by their own vector in KGE
E.g., $Rome_{Italy} \Rightarrow v(Rome_{Italy})$
E.g, $country(Rome_{Italy}, Italy) = 1 \Rightarrow country(v(Rome_{Italy}), v(Italy)) = 1$

LTN

# LTN-based Reconciliation

**4** Schema-level table annotation

- With the language defined by the LTN we can made infinite annotations by combining symbols

country

City    Country

| Paris | France |
| Paris | Texas |

| Paris | France |
| Paris | Texas |

**AR = City(x) $\wedge$ Country(y) $\wedge$ country(x, y)**
(x is a city, y is a country, and y is the country of x)

**User
Dataset**

**Schema
Annotation**

**Annotation
to Rule**

# LTN-based Reconciliation

**5** Iterate over table rows and test AR for all candidates (pairwise)

**Paris**

Paris$_{mythology}$

Paris$_{France}$

Paris$_{Texas}$

**France**

Île-de-France

France

Tour_de_France

| Paris | France |

**ROW 1**

**Candidates Generation**

**Pairwise Test**

**Scoring**

# LTN-based Reconciliation

**5** Iterate over table rows and test AR for all candidates (pairwise)



| Paris | France |
|-------|--------|

**ROW 1**

**Paris**
Paris$_{mythology}$
Paris$_{France}$
Paris$_{Texas}$

**France**
Île-de-France
France
Tour_de_France

1 Paris$_{France}$
France

2 Paris$_{France}$
Île-de-France

3 Paris$_{France}$
Tour_de_France

4 Paris$_{mythology}$
France

5 Paris$_{mythology}$
Île-de-France

6 Paris$_{mythology}$
Tour_de_France

7 Paris$_{Texas}$
France

8 Paris$_{Texas}$
Île-de-France

9 Paris$_{Texas}$
Tour_de_France

**Candidates Generation**

**Pairwise Test**

**Scoring**

# LTN-based Reconciliation



**5** Iterate over table rows and test AR for all candidates (pairwise)

the higher the better

**Candidates Generation**

| Paris | France |

ROW 1

**Paris**
Paris$_{mythology}$
Paris$_{France}$
Paris$_{Texas}$

**France**
Île-de-France
France
Tour_de_France

**Pairwise Test**

1. Paris$_{France}$ France
2. Paris$_{France}$ Île-de-France
3. Paris$_{France}$ Tour_de_France
4. Paris$_{mythology}$ France
5. Paris$_{mythology}$ Île-de-France
6. Paris$_{mythology}$ Tour_de_France
7. Paris$_{Texas}$ France
8. Paris$_{Texas}$ Île-de-France
9. Paris$_{Texas}$ Tour_de_France

LTN

**Scoring**

1. City(Paris$_{France}$) ∧ Country(France) ∧ country(Paris$_{France}$, France) — .99
2. City(Paris$_{France}$) ∧ Country(Île-de-France) ∧ country(Paris$_{France}$, Île-de-France) — .64
3. City(Paris$_{France}$) ∧ Country(Tour_de_France) ∧ country(Paris$_{France}$, Tour_de_France) — .62
4. City(Paris$_{mythology}$) ∧ Country(France) ∧ country(Paris$_{mythology}$, France) — .58
5. City(Paris$_{mythology}$) ∧ Country(Île-de-France) ∧ country(Paris$_{mythology}$, Île-de-France) — .20
6. City(Paris$_{mythology}$) ∧ Country(Tour_de_France) ∧ country(Paris$_{mythology}$, Tour_de_France) — .15
7. City(Paris$_{Texas}$) ∧ Country(France) ∧ country(Paris$_{Texas}$, France) — .80
8. City(Paris$_{Texas}$) ∧ Country(Île-de-France) ∧ country(Paris$_{Texas}$, Île-de-France) — .65
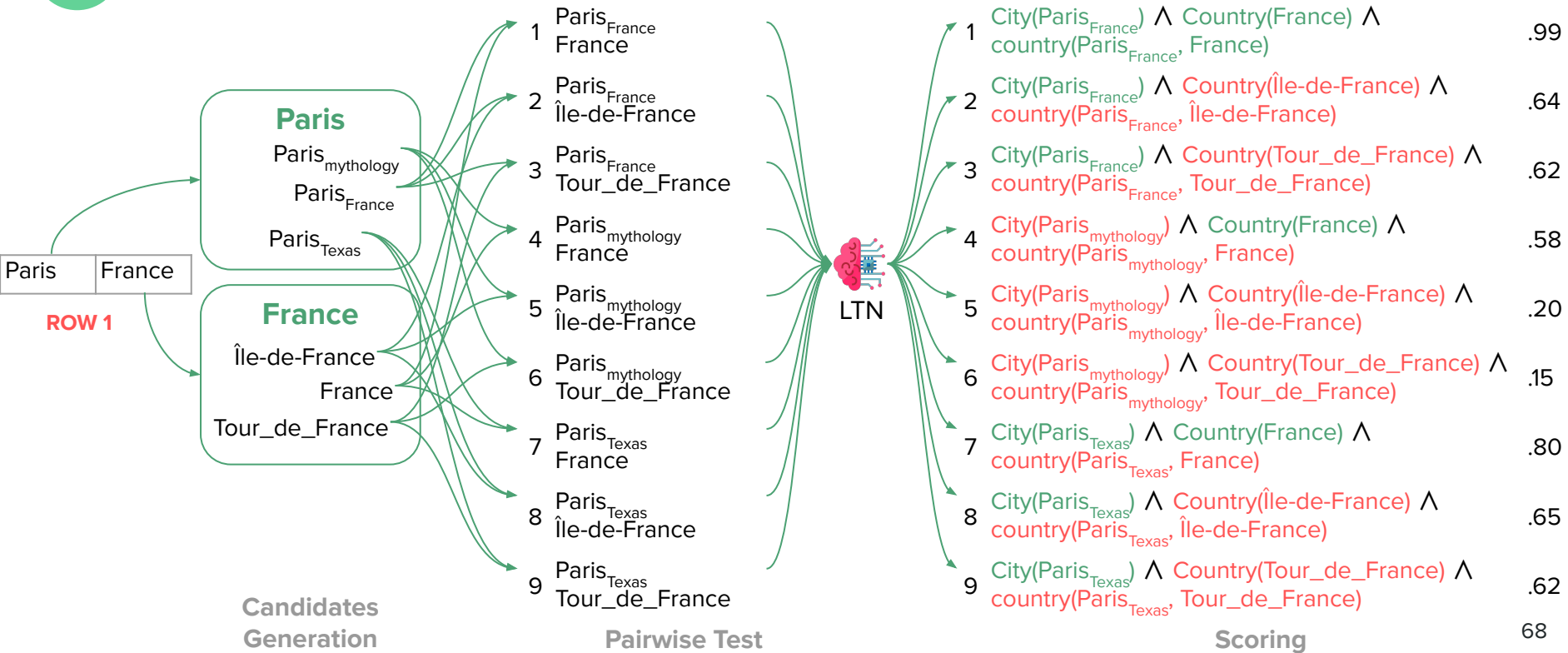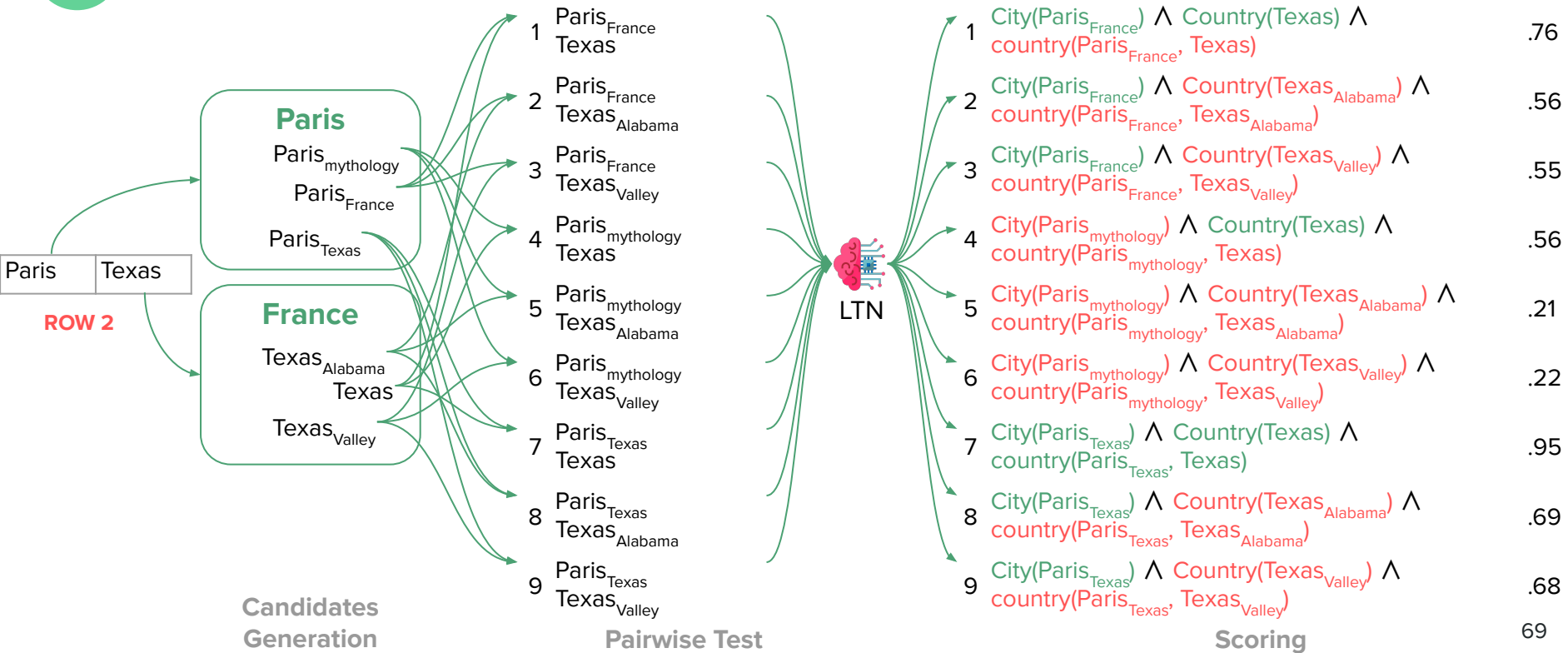9. City(Paris$_{Texas}$) ∧ Country(Tour_de_France) ∧ country(Paris$_{Texas}$, Tour_de_France) — .62

# LTN-based Reconciliation



Iterate over table rows and test AR for all candidates (pairwise)

the higher the better

# LTN-based Reconciliation

**6** Select the best candidates for each row

the higher
the better

City(**Paris$_{France}$**) ∧ Country(**France**) ∧
country(**Paris$_{France}$, France**)

**.99** → Paris,_France, France

City(Paris$_{Texas}$) ∧ Country(France) ∧
country(Paris$_{Texas}$, France)

.80

City(Paris$_{Texas}$) ∧ Country(Île-de-France) ∧
country(Paris$_{Texas}$, Île-de-France)

.65

| Paris | France |
|-------|--------|

**ROW 1**

City(Paris$_{France}$) ∧ Country(Texas) ∧
country(Paris$_{France}$, Texas)

.76

City(**Paris$_{Texas}$**) ∧ Country(**Texas**) ∧
country(Paris$_{Texas}$, Texas)

**.95** → Paris,_Texas, Texas

City(Paris$_{Texas}$) ∧ Country(Texas$_{Alabama}$) ∧
country(Paris$_{Texas}$, Texas$_{Alabama}$)

.69

| Paris | Texas |
|-------|-------|

**ROW 2**

70

# Experimental results: Datasets

**Dataset:**

- 8 african countries
- No more than 50 cities for each country

**Embedding:**

- Pretrained DBpedia embeddings from RDF2VEC (200 dimensions)
  - Only cities and countries
- Embeddings downsized to 40 dimensions (using PCA)
  - Cosine similarity between vectors is preserved (similar vectors are still similar in the new space)

# Experimental results: Training capital() and locatedIn()

**Universally quantified axiom:**

**Legend:**
*a-b: all cities*
*c-d: all countries*

- ∀ ?a,?c,?d: locatedIn(?a,?c) -> (¬ equals(?c,?d) & ¬ locatedIn(?a,?d))
- ∀ ?a,?b,?c: capital(?a,?c) -> (¬ equals(?a,?b) & ¬ capital(?b,?c))
- ∀ ?a,?c: capital(?a,?c) -> locatedIn(?a,?c)
- ∀ ?a,?c: ¬ locatedIn(?a,?c) -> ¬ capital(?a,?c)

|  | **TRAINING** | **TEST** |
|---|---|---|

**locatedIn():**
threshold = 0.80

**Precision: 0.95**
**Recall: 0.97**

**capital():**
threshold = 0.95
(In almost all cases the right pair is the one with the highest score)

**Precision: 0.62**
**Recall: 1.00**

**TEST**
16 cities never trained (about 2 cities for each country)
- **16 / 16** cities properly assessed (**locatedIn()**)
- **16 / 16** cities properly assessed (**capital()**)

# Experimental results: Training unary predicates

**Universally quantified axiom:**

**Legend:**
*a-b: all cities*
*c-d: all countries*

- ∀ ?a,?c,?d: locatedIn(?a,?c) -> (¬ equals(?c,?d) & ¬ locatedIn(?a,?d))
- ∀ ?a,?c: locatedIn(?a, ?c) -> City(?a) & Country(?c)
- ∀ ?a: ¬ Country(?a)
- ∀ ?c: ¬ City(?c)

### TRAINING

**locatedIn():**

threshold = 0.80

**Precision: 1.00**
**Recall: 0.99**

**377/378 cities satisfy City()**
**8/8 countries satisfy Country()**

### TEST

16 cities never trained (about 2 cities for each country)

- **16 / 16** cities properly assessed (**locatedIn()**)
- **all cities** have **City()** value > **0.5** and **Country()** value < **0.5**
- **all countries** have **Country()** value > **0.5** and **City()** value < **0.5**

# Experimental results: Training all predicates

**Universally quantified axiom:**

- ∀ ?a,?c,?d: locatedIn(?a,?c) -> (¬ equals(?c,?d) & ¬ locatedIn(?a,?d))
- ∀ ?a,?b,?c: capital(?a,?c) -> (¬ equals(?a,?b) & ¬ capital(?b,?c))
- ∀ ?a,?c: capital(?a,?c) -> locatedIn(?a,?c)
- ∀ ?a,?c: ¬ locatedIn(?a,?c) -> ¬ capital(?a,?c)
- ∀ ?y: Capital(?y)
- ∀ ?x: ¬ Capital(?x)
- ∀ ?a: City(?a)
- ∀ ?a: ¬ Country(?a)
- ∀ ?c: Country(?c)
- ∀ ?c: ¬ City(?c)
- ∀ ?c: ¬ Capital(?c)

**Legend:**
*a-b: all cities*
*c-d: all countries*
*y: all capitals*
*x: all non capitals*

# Experimental results: Training all predicates (cont)

**TRAINING**

**locatedIn():**     **Precision:  0.94**
threshold = 0.70      **Recall:  0.95**

**capital():**        **Precision:  0.60**
threshold = 0.90      **Recall:  1.00**

**377/378 cities that satisfy City()**
**8/8 countries that satisfy Country()**
**8/8 cities that satisfy Capital()**

**TEST**

16 cities never trained (about 2 cities for each country)

- **16 / 16** cities properly assessed (**locatedIn()**)
- **16 / 16** cities properly assessed (**capital()**)
- **all cities** have:
    - **City()** value > **0.5**
    - **Country()** value < **0.5**
    - **Capital()** value < **0.5**
- **all countries** have:
    - **Country()** value > **0.5**
    - **City()** value < **0.5**
    - **Capital()** value < **0.5**

# Thanks!

# References

d'Avila Garcez, A., Lamb, L. C., & Gabbay, D. M. (2008). Neural-symbolic cognitive reasoning. Springer Science & Business Media.

d'Avila Garcez, A., Broda, K. B., & Gabbay, D. M. (2012). Neural-symbolic learning systems: foundations and applications. Springer Science & Business Media.

Serafini, L., & d'Avila Garcez A. (2016). Logic Tensor Networks: Deep Learning and Logical Reasoning from Data and Knowledge. 11th International Workshop on Neural-Symbolic Learning and Reasoning (NeSy16) arXiv:1606.04422.

Socher, R., Chen, D., Manning, C. D., & Ng, A. Y. (2013). Reasoning with neural tensor networks for knowledge base completion. Proceedings of the 26th International Conference on Neural Information Processing Systems - Volume 1 (NIPS'13). pp. 926–934.

Wang, Q., Mao, Z., Wang, B., & Guo, L. (2017). Knowledge Graph Embedding: A Survey of Approaches and Applications. IEEE Transactions on Knowledge and Data Engineering, vol. 29, no. 12, pp. 2724-2743.