

Open Information Extraction

Matthias Huck and Alexander Fraser

Center for Information and Language Processing
LMU Munich

22 January 2020

Introduction

Reminder: We know how to:

- Recognize named entities in natural language text
- Extract binary relations between named entities

We have seen an application scenario:

- Relations can be stored in a knowledge base
- And be used in question answering or spoken dialogue systems

*But so far, there are **limitations**, most notably:*

- We have dealt with narrow domains (such as geographical location, food, plant seed development)
- The sets of entity types and relations were closed and manually defined

Open IE aims at:

- Not being limited to any single domain
- Not being limited to pre-defined entity types and relations

Outline

- 1 Open IE: Motivation & Task Definition
- 2 Open Relation Extraction: TEXTRUNNER & REVERB
- 3 Open Relation Extraction: OLLIE
- 4 Open Relation Extraction: STANFORD OPENIE
- 5 Discussion: Further Challenges
- 6 Conclusion

OPEN IE:
MOTIVATION & TASK DEFINITION

Open IE: Motivation (1)

Example Queries: ²

What kills bacteria?

Who built the Pyramids?

What did Thomas Edison invent?

What contains antioxidants?

Typed Example Queries: ²

What countries are located in Africa?

What actors starred in which films?

What is the symbol of which country?

What foods are grown in which countries?

What drug ingredients has the FDA approved?

Argument 1:

Relation:

Argument 2:

Corpus:

[openie.allenai.org query, 16 Jan. 2017]

Open IE: Motivation (2)

Argument 1:

Relation: built

Argument 2: pyramids

All

 Search

62 answers from **584 sentences** (cached)

all

deceased person (7)

monarch (5)

location (3)

ethnicity (3)

misc.

Egyptians (132)

Ancient Egypt (123)

aliens (44)

the people (38)

slaves (29)

Open IE: Motivation (3)

Argument 1:	<input type="text" value="type:Country"/>	Relation:	<input type="text" value="is located in"/>
Argument 2:	<input type="text" value="Africa"/>	All	<input type="button" value="Search"/>

41 answers from **275 sentences** (cached)

Kenya (31)

Ghana (28)

Nigeria (16)

Egypt (15)

Morocco (11)

Algeria (10)

Open IE: Motivation (4)

Argument 1:

Relation:

Argument 2:

All



Search

1 answers from **27 sentences** (cached)

Elvis Presley (27)

[openie.allenai.org query, 16 Jan. 2017]

Open IE: Motivation (5)

Argument 1:	<input type="text"/>	Relation:	<input type="text" value="founded"/>
Argument 2:	<input type="text" value="Microsoft"/>	<input type="button" value="All"/>	<input type="button" value="Search"/>

7 answers from **93 sentences** (cached)

Bill Gates (64)

Paul Allen (18)

the man (3)

Allen (2)

Microsoft Bill Gates (2)

Open IE: Motivation (6)

Argument 1: Scots

Relation: eat

Argument 2:

All

 Search

2 answers from **4 sentences** (cached)

Haggis (2)

Potato (2)

[openie.allenai.org query, 16 Jan. 2017]

Open IE: Motivation (7)

Argument 1: Faust

Argument 2: devil

Relation:

All

 Search

6 answers from **55 sentences** (cached)

sold his soul to (35)

makes a pact with (11)

makes a deal with (3)

strikes a deal with (2)

had sold his soul to (2)

asks (2)

Open IE: Motivation (8)

Argument 1:

Relation:

Argument 2:

All

 Search

6 answers from **35 sentences** (cached)

Tent (20)

Dormitory (7)

classrooms (2)

car (2)

thatched-roof bandas (2)

class (2)

Open IE: Task Definition

We want to *find any relation* that is expressed in large data.

- Cannot resort to specialized domain knowledge
- Cannot think of all possible relation types beforehand
- Should rather not force all possible arguments into a rigid set of entity types

We want to *scale to billions of documents* that are heterogeneous wrt. domains, quality, credibility.

- Which relations are correct?
- Which are uninformative or incoherent?
- Which are redundant?

OPEN RELATION EXTRACTION: TEXTRUNNER & REVERB

Open Relation Extraction: Example

Hudson was born in Hampstead, which is a suburb of London.

(Hudson, was born in, Hampstead)

(Hampstead, is a suburb of, London)

[Fader *et al.*. Identifying Relations for Open Information Extraction. Proc. of EMNLP, Edinburgh, Scotland, UK, July 2011.]

Open Relation Extraction: Basic Approach

Learn a general model of how (arbitrary) relations are expressed in a particular language.

- Neither relation names nor argument types known in advance
- Bootstrap with heuristics or distant supervision
- Train a (sequence) classifier (often with unlexicalized features)

Rel. Freq.	Category	Simplified Lexico-Syntactic Pattern	Example
37.8	Verb	E_1 Verb E_2	X established Y
22.8	Noun+Prep	E_1 NP Prep E_2	X settlement with Y
16.0	Verb+Prep	E_1 Verb Prep E_2	X moved to Y
9.4	Infinitive	E_1 to Verb E_2	X plans to acquire Y
5.2	Modifier	E_1 Verb E_2 Noun	X is Y winner
1.8	Coordinate _n	E_1 (and , - :) E_2 NP	X-Y deal
1.0	Coordinate _v	E_1 (and ,) E_2 Verb	X , Y merge
0.8	Appositive	E_1 NP (: ,)? E_2	X hometown : Y

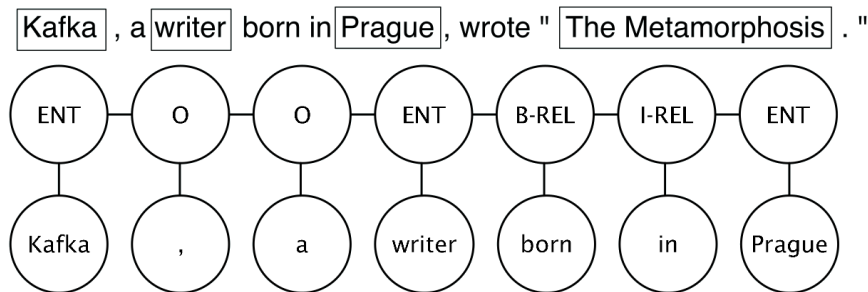
[Etzioni *et al.*. Open Information Extraction from the Web. Communications of the ACM, vol. 51 no. 12, Dec. 2008.]

[Banko and Etzioni. The Tradeoffs Between Open and Traditional Relation Extraction. Proc. of the ACL, Columbus, OH, USA, June 2008.]

Open Relation Extraction: “Three-Step Method”

- ① **Label:** Sentences are automatically labeled with extractions using heuristics or distant supervision.
- ② **Learn:** A relation phrase extractor is learned, e.g. using a sequence-labeling graphical model (CRF).
- ③ **Extract:** The system takes a sentence as input, identifies a candidate pair of NP arguments (arg1, arg2) from the sentence, and then uses the learned extractor to label each word between the two arguments as part of the relation phrase or not.

Open Relation Extraction as Sequence Labeling



[Banko and Etzioni. The Tradeoffs Between Open and Traditional Relation Extraction. Proc. of the ACL, Columbus, OH, USA, June 2008.]

Uninformative Relations

Problem 1: The sequence classifier may come up with an uninformative relation name.

Faust made a deal with the devil.

(Faust, made, deal)

(Faust, made deal with, devil)

Uninformative	Completion
is	is an album by, is the author of, is a city in
has	has a population of, has a Ph.D. in
made	made a deal with, made a promise to
took	took place in, took control over, took advantage of
gave	gave birth to, gave a talk at, gave new meaning to
got	got tickets to, got a deal on, got funding from

Incoherent Relations

Problem 2: The sequence classifier may come up with an incoherent relation name.

Sentence	Incoherent Relation
The guide contains dead links and omits sites.	contains omits
The Mark 14 was central to the torpedo scandal of the fleet.	was central torpedo
They recalled that Nungesser began his career as a precinct leader.	recalled began

[Fader *et al.*. Identifying Relations for Open Information Extraction. Proc. of EMNLP, Edinburgh, Scotland, UK, July 2011.]

POS-based Constraints to Avoid Incoherence & Uninformativeness

Extendicare agreed to buy Arbor Health Care for about US \$432 million in cash and assumed debt.

(Arbor Health Care, for assumed, debt)

- POS-based **regular expressions** help avoid extraction of uninformative or incoherent relation phrases
- Manually written; e.g. the relation phrase must match:

$V \mid V P \mid V W^* P$
$V = \text{verb particle? adv?}$
$W = (\text{noun} \mid \text{adj} \mid \text{adv} \mid \text{pron} \mid \text{det})$
$P = (\text{prep} \mid \text{particle} \mid \text{inf. marker})$

- Choose longest possible match
- Require the relation phrase to appear between its arguments

Overspecific Relations & How to Avoid Them

Problem 3: Some relations are specific to an argument pair, or have only a few possible instances.

The Obama administration is offering only modest greenhouse gas reduction targets at the conference.

(Obama administration, is offering only modest greenhouse gas reduction targets at, conference)

- Intuition: a valid relation phrase should take many distinct arguments in a large corpus
- **Lexical constraint:** relation phrases are required to match at least k distinct argument pairs in the data (e.g., $k = 20$)

Relation Phrase Normalization

Shakespeare (*has written* | *wrote* | *was writing*)
Hamlet.

Allow for minor variations in relation phrases.

- Remove inflection
- Remove auxiliary verbs, adjectives, adverbs

[Fader *et al.*. Identifying Relations for Open Information Extraction. Proc. of EMNLP, Edinburgh, Scotland, UK, July 2011.]

Confidence Function

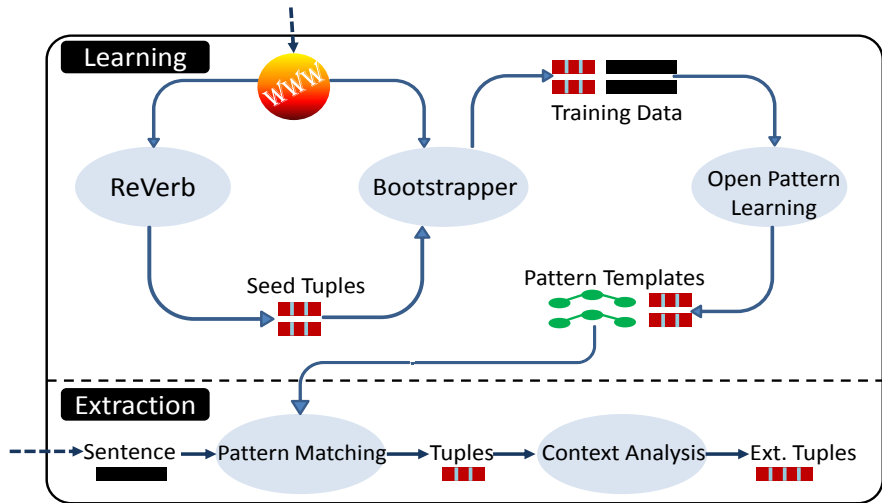
- Train classifier to assign a confidence score to each extraction
- Trade recall for precision by tuning a confidence threshold

Weight	Feature
1.16	(x, r, y) covers all words in s
0.50	The last preposition in r is <i>for</i>
0.49	The last preposition in r is <i>on</i>
0.46	The last preposition in r is <i>of</i>
0.43	$\text{len}(s) \leq 10$ words
0.43	There is a <i>WH</i> -word to the left of r
0.42	r matches $VW^* P$
0.39	The last preposition in r is <i>to</i>
0.25	The last preposition in r is <i>in</i>
0.23	$10 \text{ words} < \text{len}(s) \leq 20 \text{ words}$
0.21	s begins with x
0.16	y is a proper noun
0.01	x is a proper noun
-0.30	There is an NP to the left of x in s
-0.43	$20 \text{ words} < \text{len}(s)$
-0.61	r matches V
-0.65	There is a preposition to the left of x in s
-0.81	There is an NP to the right of y in s
-0.93	Coord. conjunction to the left of r in s

OPEN RELATION EXTRACTION: OLLIE

OLLIE (Open Language Learning for Information Extraction)

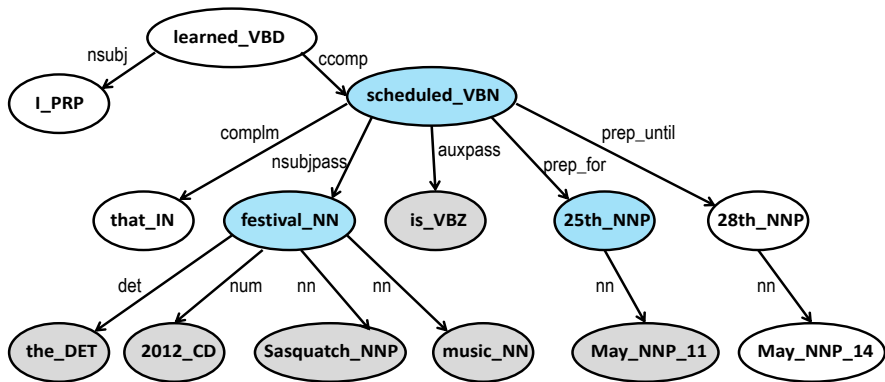
Bootstrapping with high precision seed tuples from existing system (REVERB, cf. previous part)



Employing Dependency Parses

I learned that the 2012 Sasquatch music festival is scheduled for May 25th until May 28th.

(the 2012 Sasquatch Music Festival, is scheduled for, May 25th)



Open Pattern Templates

Open pattern templates encode the ways in which a relation may be expressed in a sentence.

- Based on a dependency parse path
- with lexical constraint
- and POS constraint

Extraction Template	Open Pattern
1. (arg1; be {rel} {prep}; arg2)	{arg1} ↑nsubjpass↑ {rel:postag=VBN} ↓{prep_*}↓ {arg2}
2. (arg1; {rel}; arg2)	{arg1} ↑nsubj↑ {rel:postag=VBD} ↓dobj↓ {arg2}
3. (arg1; be {rel} by; arg2)	{arg1} ↑nsubjpass↑ {rel:postag=VBN} ↓agent↓ {arg2}
4. (arg1; be {rel} of; arg2)	{rel:postag=NN; type=Person } ↑nn↑ {arg1} ↓nn↓ {arg2}
5. (arg1; be {rel} {prep}; arg2)	{arg1} ↑nsubjpass↑ {slot:postag=VBN; lex ∈ announce name choose... } ↓dobj↓ {rel:postag=NN} ↓{prep_*}↓ {arg2}

OLLIE: Advantages

Previously (in REVERB), we required the relation phrase to appear between its arguments:

```
Elvis married Priscilla.
```

Open pattern templates may help with:

```
Elvis and Priscilla are married.
```

Other systems are designed to have verb-mediated relation phrases:

```
Bill Gates founded Microsoft.
```

OLLIE can deal with noun-mediated relations:

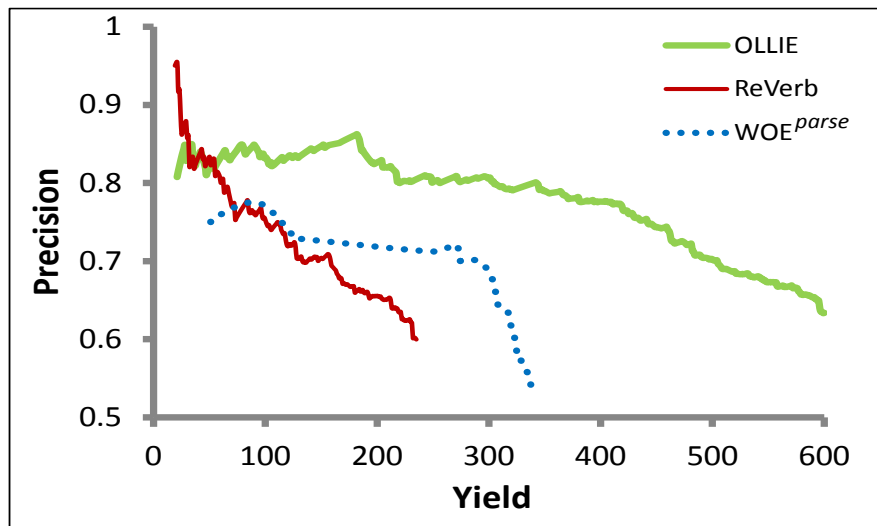
```
Bill Gates is founder of Microsoft.
```

Many relationships are most naturally expressed via noun phrases:

```
is capital of, is president of, is professor at, ...
```

Dependency parse is useful; parsers not deemed too slow any more.

OLLIE: Evaluation



OPEN RELATION EXTRACTION: STANFORD OPENIE

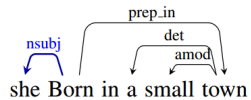
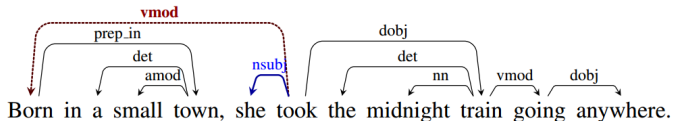
STANFORD OPENIE

Heavily based on dependency parses.

- ① Each dependency-parsed sentence is first split into a set of entailed clauses
- ② Clauses are then maximally shortened, producing a set of entailed shorter sentence fragments
- ③ The fragments are segmented into relation triples, and output by the system

[Angeli *et al.*. Leveraging Linguistic Structure For Open Domain Information Extraction. Proc. of the ACL, Beijing, China, July 2015.]

STANFORD OPENIE: Illustration



(input)

(extracted clause)

she took the midnight train going anywhere

she took the midnight train

Born in a small town, she took the midnight train

she took midnight train

Born in a town, she took the midnight train

...

she Born in small town

she Born in a town

she Born in town

(she; took; midnight train)

(she; born in; small town)

(she; born in; town)

[<http://nlp.stanford.edu/software/openie.html>, 17 Jan. 2017]

[Angeli *et al.*. Leveraging Linguistic Structure For Open Domain Information Extraction. Proc. of the ACL, Beijing, China, July 2015.]

Clause Splitting as a Classification Task

- Inspect the dependency structure
- Decide whether to split on a dependency arc
- Classifier using a set of dependency-based features
- Distant supervision for training:
sequence which recovers a known relation is correct

Feature Class	Feature Templates
Edge taken	$\{l, \text{short_name}(l)\}$
Last edge taken	$\{\text{incoming_edge}(p)\}$
Neighbors of parent	$\{\text{nbr}(p), (p, \text{nbr}(p))\}$
Grandchild edges	$\{\text{out_edge}(c),$ $(e, \text{out_edge}(c))\}$
Grandchild count	$\{\text{count}(\text{nbr}(e_{\text{child}}))$ $(e, \text{count}(\text{nbr}(e_{\text{child}})))\}$
Has subject/object	$\forall e \in \{e, e_{\text{child}}\} \forall l \in \{\text{subj}, \text{obj}\}$ $\mathbb{1}(l \in \text{nbr}(e))$
POS tag signature	$\{\text{pos}(p), \text{pos}(c),$ $(\text{pos}(p), \text{pos}(c))\}$
Features at root	$\{\mathbb{1}(p = \text{root}), \text{POS}(p)\}$

Atomic Patterns over Short Entailed Sentences

	Input	Extraction
Verb-mediated:	<i>cats play with yarn</i>	(cats; play with; yarn)
	<i>fish like to swim</i>	(fish; like to; swim)
	<i>cats have tails</i>	(cats; have; tails)
	<i>cats are cute</i>	(cats; are; cute)
	<i>Tom and Jerry are fighting</i>	(Tom; fighting; Jerry)
	<i>There are cats with tails</i>	(cats; have; tails)

	Input	Extraction
Noun-mediated:	<i>Durin, son of Thorin</i>	(Durin; is son of; Thorin)
	<i>Thorin's son, Durin</i>	(Thorin; 's son; Durin)
	<i>IBM CEO Rometty</i>	(Rometty; is CEO of; IBM)
	<i>President Obama</i>	(Obama; is; President)
	<i>Fischer of Austria</i>	(Fischer; is of; Austria)
	<i>IBM's research group</i>	(IBM; 's; research group)
	<i>US president Obama</i>	(Obama; president of; US)
	<i>Our president, Obama,</i>	(Our president; be; Obama)

Validating Deletions with Natural Logic

Scopes of operators *all*, *no*, *many*, ...

- all rabbits eat fresh vegetables
yields (rabbits, eat, vegetables)
- all young rabbits drink milk
does not yield (rabbits, drink, milk)

Non-subjective adjectives

- a *fake gun* is not a gun

Prepositional attachment

- Alice played baseball on Sunday
entails Alice played on Sunday
- Obama signed the bill on Sunday
should not entail Obama signed on Sunday

STANFORD OPENIE: Example Extractions

Born in Honolulu, Hawaii, Obama is a US Citizen.

Our System

(Obama; is; US citizen)

(Obama; born in;
Honolulu, Hawaii)

Ollie

(Obama; is; a US citizen)

(Obama; be born in; Honolulu)
(Honolulu; be born in; Hawaii)
(Obama; is citizen of; US)

Friends give true praise.

Enemies give fake praise.

Our System

(friends; give; true praise)

(friends; give; praise)

(enemies; give; fake praise)

Ollie

(friends; give; true praise)

(enemies; give; fake praise)

Heinz Fischer of Austria visits the US

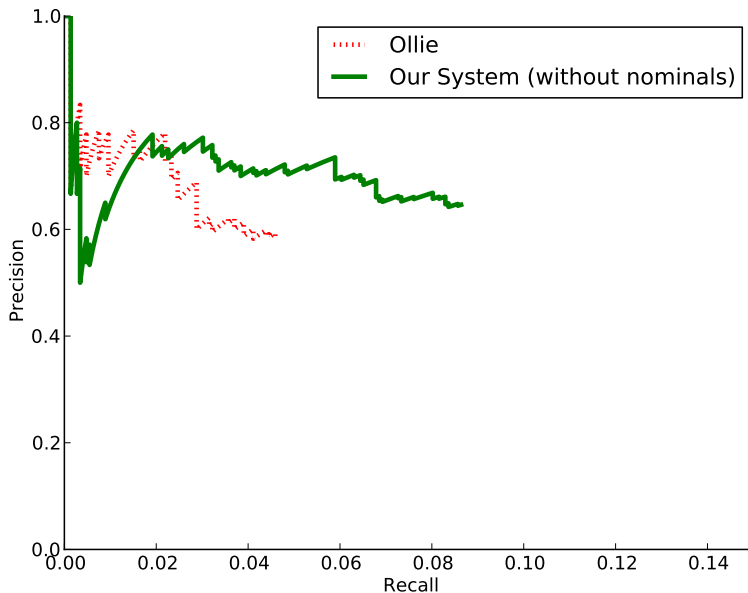
Our System

(Heinz Fischer; visits; US)

Ollie

(Heinz Fischer **of Austria;**
visits; **the US**)

STANFORD OPENIE: Evaluation



DISCUSSION:
FURTHER CHALLENGES

Synonym Resolution

The same entity may be referred to by a variety of names.

- *Michael Jackson; Jacko; The King of Pop; . . .*

The same fact may be expressed in a variety of ways.

- *IBM built Watson*
IBM created Watson
IBM invented Watson
...
• *Dookie is a record by Green Day*
Dookie is an album by Green Day
...

RESOLVER identifies synonymous relations and objects

[Yates and Etzioni. Unsupervised Resolution of Objects and Relations on the Web. Proc. of the NAACL, Rochester, NY, USA, April 2007.]

[Etzioni *et al.*. Open Information Extraction from the Web. Communications of the ACM, vol. 51 no. 12, Dec. 2008.]

Disambiguation

The same string may refer to different entities (especially across different domains).

- *Watson*, the founder of IBM; *Watson*, the computer system
- *mouse*, the animal; *mouse*, the input device
- *1984*, the year; *1984*, the book
- *Paris*, France; *Paris*, Texas

Vagaries of Natural Language

- pronoun resolution
- metaphor
- anaphora
- complex or ungrammatical sentences
- irony, sarcasm
- ...

Incorrect Information

Nowadays referred to as “fake news”.

- *Elvis killed JFK*

Rate the reliability of an extracted relation.

- The relation extractor may have made an error:
cf. the previously discussed confidence function
- Occurrence frequencies over the whole corpus
can give an indication
- Credibility of the source of a document

YAGO-NAGA ranks facts f via:

$$\text{confidence}(f) = \max \{ \text{accuracy}(f, s) \times \text{trust}(s) \mid s \in \text{witnesses}(f) \}$$

[Kasneji *et al.*. The YAGO-NAGA Approach to Knowledge Discovery. ACM SIGMOD Record Volume 37 Issue 4, Dec. 2008.
<https://suchanek.name/work/publications/sigmodrec2008.pdf>]

[Etzioni *et al.*. Open Information Extraction from the Web. Communications of the ACM, vol. 51 no. 12, Dec. 2008.]

Temporal and Spatial Aspects

Time.

- The capital city of the Federal Republic of Germany? Berlin today. Bonn from 1949 to 1990.
- Plato has not met with Tsipras

Space.

- An elephant does not fit into a coffee mug
- Trees don't travel
- Somebody who pays in GBP is probably located in Britain
- Plato has never seen a kangaroo

Fact Consistency Checks

Avoid contradictory facts within the knowledge base.

- Elvis died in 460 AD cannot refer to *Elvis Presley* if we already know that Elvis Presley was born in 1935
- $\text{born}(X,Y) \wedge \text{died}(X,Z) \Rightarrow Y < Z$
- $\text{appears}(A,P,B) \wedge R(A,B) \Rightarrow \text{expresses}(P,R)$
 $\text{appears}(A,P,B) \wedge \text{expresses}(P,R) \Rightarrow R(A,B)$
- $\text{means}(\text{"Elvis"}, \text{Elvis_Presley}, 0.8)$
 $\text{means}(\text{"Elvis"}, \text{Elvis_Costello}, 0.2)$

Implemented in the SOFIE IE system, which aims to extend the YAGO knowledge base

[Suchanek. Information Extraction for Ontology Learning. Book chapter in Völker and Lehman: Perspectives on Ontology Learning, 2014. <https://suchanek.name/work/publications/ontologybookchapter.pdf>]

CONCLUSION

Summary: Open IE

- Discovering relations without a closed set of pre-defined relation types
- Open-domain
- Learning from the whole Web
- Distant supervision / bootstrapping to get started
- Attention to detail required to avoid pitfalls
- The system should benefit from the sheer size of the data
- It should learn more by itself when being run perpetually, and become more reliable

THE END!

Thank you for your attention

Matthias Huck
mhuck@cis.lmu.de

Argument 1:	<input type="text" value="thanks"/>	Relation:	<input type="text"/>
Argument 2:	<input type="text" value="questions"/>	All	<input type="button" value="Search"/>

27 answers from **129 sentences** (cached)

feel free to ask (25)

here are (17)

have (16)

here are my answers to (7)

please feel free to ask (6)

Give your answer to (5)

go ahead with (5)