

# Information extraction & automated knowledge graph construction

Lecture 4: From IE to AKGC

---

Lauriane Aufrant (Inria Paris)

Université Paris Cité – UFRL – M2 LI – 30/11/2023

- ▶ ~~09/11 – Overview of information extraction~~
- ▶ ~~16/11 – Entities, relations...~~
- ▶ ~~23/11 – Coreference, linking...~~
- ▶ 30/11 – From IE to automated knowledge graph construction
- ▶ 07/12 – IE annotations
- ▶ 14/12 – IE in a specialty domain

## Today: From IE to AKGC

- ▶ Comprehensiveness
- ▶ Consistency
- ▶ Task interoperability
- ▶ Combining tasks

**Comprehensiveness**

*The dean's secretary saw a guy with a green backpack jogging down the avenue to the nearest restaurant.*

## Comprehensiveness for entities

*The dean's secretary saw a guy with a green backpack jogging down the avenue to the nearest restaurant.*

## Comprehensiveness for entities

*The dean's secretary saw a guy with a green backpack jogging down the avenue to the nearest restaurant.*

- ▶ Large set of entity types
  - ↪ “a green backpack”

## Comprehensiveness for entities

*The dean's secretary saw a guy with a green backpack jogging down the avenue to the nearest restaurant.*

- ▶ Large set of entity types  
↪ “a green backpack”
- ▶ Nested entities  
↪ “The dean's secretary”



## Comprehensiveness for entities

*The dean's secretary saw a guy with a green backpack jogging down the avenue to the nearest restaurant.*

- ▶ Large set of entity types  
↪ “a green backpack”
- ▶ Nested entities  
↪ “The dean’s secretary”
- ▶ Non-named entities  
↪ All of them!

## Comprehensiveness for relations

*The dean's secretary saw a guy with a green backpack jogging down the avenue to the nearest restaurant.*

## Comprehensiveness for relations

*The dean's secretary saw a guy with a green backpack jogging down the avenue to the nearest restaurant.*

- ▶ Unforeseen relations: “jogging down”

## Comprehensiveness for relations

*The dean's secretary saw a guy with a green backpack jogging down the avenue to the nearest restaurant.*

- ▶ Unforeseen relations: “jogging down”
- ▶ “with”?

## Comprehensiveness for relations

*The dean's secretary saw a guy with a green backpack jogging down the avenue to the nearest restaurant.*

- ▶ Unforeseen relations: “jogging down”
- ▶ “with”?
- ▶ “nearest”?

## Comprehensiveness for relations

*The dean's secretary saw a guy with a green backpack jogging down the avenue to the nearest restaurant.*

- ▶ Unforeseen relations: “jogging down”
- ▶ “with”?
- ▶ “nearest”?
- ▶ “saw”: “a guy with [...] restaurant”?

# Open Information Extraction

**Bill Gates, Microsoft co-founder**, stepped down as **CEO** in January 2000. **Gates** was included in the **Forbes wealthiest list** since 1987 and **was the wealthiest** from 1995 to 2007...

relation extraction



Co-founder(Bill Gates, Microsoft)  
Director-of (MacLorraine, Ciao)  
Employee-of (MacLorraine, Ciao)  
...

It was announced that **IBM** would buy **Ciao** for an undisclosed amount. The **CEO, MacLorraine** has occupied the **corner office of the Hopkinton**, company

open IE



(Bill Gate, be, Microsoft co-founder)  
(Bill Gates, stepped down as, CEO)  
(Bill Gates, was included in, the Forbes wealthiest list)  
(Bill Gates, was, the wealthiest)  
(IBM, would buy, Ciao)  
(MacLorraine, has occupied, the corner office of the Hopkinton)  
...

The company's storage business is also threatened by new, born-on-the Web could providers like Dropbox and Box, and ...

- ▶ Often, none: fully unsupervised & manual evaluation
- ▶ ORE [Mesquita et al., 2013] (662 binary & 222 n-ary relations + automatic annotations, arguments limited to named entities, 0.6-1 fact / sentence)
- ▶ QA-SRL OIE [Stanovsky & Dagan, 2016] (10,359 tuples, automatic extraction from questions & answers, 3 facts / sentence)
- ▶ Wire57 [Léchelle et al., 2019] (347 tuples, axiomatic guidelines designed for completeness, 6 facts / sentence)



# Evaluating OpenIE

- ▶ Exact match F1
- ▶ Relaxed containment match
- ▶ Greedy matching of tuple pairs & computing token-level overlap
- ▶ Penalizing verbosity? Over-specific tuples? Incorrect number of arguments?
- ▶ F2 for higher focus on recall

# Methods & tools for OpenIE

	Extractions	Matches	Exact matches	Prec. of matches	Recall of matches	Prec.	Recall	F1
ReVerb (Fader et al., 2011)	79	54	13	.83	.77	<b>.569</b>	.121	.200
Ollie (Mausam et al., 2012)	145	74	8	.73	.81	.347	.175	.239
ClausIE (Del Corro and Gemulla, 2013)	223	121	<b>24</b>	.74	.84	.401	.298	.342
Stanford (Angeli et al., 2015)	371	99	2	.79	.65	.210	.188	.198
OpenIE 4 (Mausam, 2016)	101	74	5	.68	.84	.501	.182	.267
PropS (Stanovsky et al., 2016)	184	69	0	.59	.80	.222	.162	.187
MinIE (Gashteovski et al., 2017)	252	<b>134</b>	10	.75	.83	.400	<b>.323</b>	<b>.358</b>

# OpenIE as sequence labeling

---

## Open IE Encoding Examples

---

(a) *The president claimed that he won the majority vote.*

(The president; **claimed that he won**; the majority vote)

The<sub>A0-B</sub> president<sub>A0-I</sub> claimed<sub>P-B</sub> that<sub>P-I</sub> he<sub>P-I</sub> won<sub>P-I</sub> the<sub>A1-B</sub> majority<sub>A1-I</sub> vote<sub>A1-I</sub>

---

(b) *Barack Obama, a former U.S president, was born in Hawaii.*

(Barack Obama; **was born in**; Hawaii)

(a former U.S. president; **was born in**; Hawaii)

Barack<sub>A0-B</sub> Obama<sub>A0-I</sub> ,O a<sub>A0-B</sub> former<sub>A0-I</sub> U.S.<sub>A0-I</sub> president<sub>A0-I</sub> ,O was<sub>P-B</sub> born<sub>P-I</sub> in<sub>P-I</sub> Hawaii<sub>A1-B</sub>

---

(c) *Theresa May plans for Brexit, on which the UK has voted last June.*

(the UK; **has voted on**; Brexit; last June)

Theresa<sub>O</sub> May<sub>O</sub> plans<sub>O</sub> for<sub>O</sub> Brexit<sub>A1-B</sub> ,O on<sub>O</sub> which<sub>O</sub> the<sub>A0-B</sub> UK<sub>A0-I</sub> has<sub>P-B</sub> voted<sub>P-I</sub> on<sub>P-I</sub> last<sub>A2-B</sub> June<sub>A2-I</sub>

---

## OpenIE as sequence generation

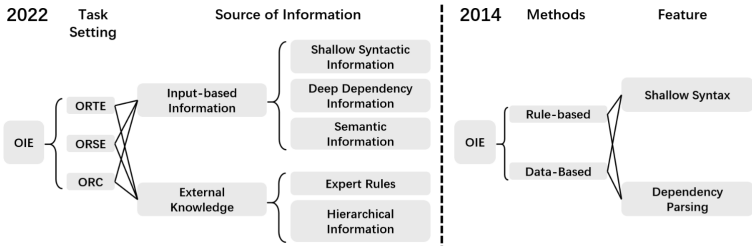
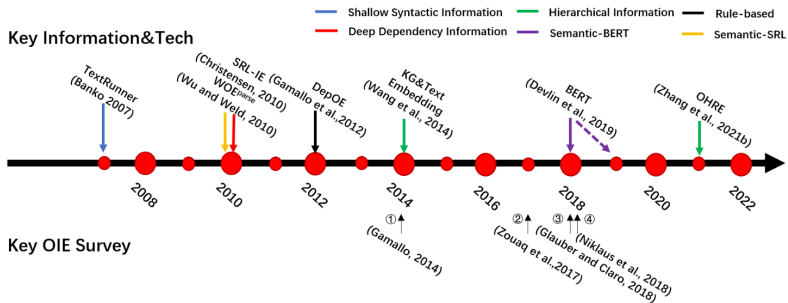
“deep learning is a subfield of machine learning”



“<arg1> deep learning </arg1> <rel> is a subfield  
of </rel> <arg2> machine learning </arg2>”

↔ Pure generation or with copy mechanisms

# Overview of OpenIE methods



- ▶ “Paul’s wife turned 35.”
- ▶ “They met him at the French-German border.”
- ▶ “The author of this paper is William Léchelle.”  
↔ “This paper has been written by William Léchelle.”

- ▶ “I read an interesting paper about OpenIE. The main author is William L chelle.”

# Implicit relations

- ▶ “Paris, France”
- ▶ “The Turing paper” (*written by*)  
“The Nature paper” (*published in*)



- ▶ “A second earthquake occurred in Tokyo that day.”

# Limits of light inference

**Sentence CH 7** – “His parents are Ashkenazi Jews who had to flee from Hungary during World War II.”

## Annotations

- (His/(Chilly Gonzales’s) parents ; are ; Ashkenazi Jews)
- (His/(Chilly Gonzales’s) parents ; are ; Jews)
- (His/(Chilly Gonzales’s) parents ; had to flee from ; Hungary ; during World War II)
- (His/(Chilly Gonzales’s) parents ; [fled] from ; Hungary ; during World War II)
- ((Chilly Gonzales) ; [has] ; parents)

**Sentence FI 2** – “A police statement did not name the man in the boot, but in effect indicated the traveler was State Secretary Samuli Virtanen, who is also the deputy to Foreign Minister Timo Soini.”

## Annotations

- (A police/(Finnish police) statement ; did not name ; (the man in the boot)/(Samuli Virtanen))
- ((the man in the boot)/(Samuli Virtanen) ; was ; Samuli Virtanen) [attributed]
- ((the traveler)/(Samuli Virtanen) ; was ; Samuli Virtanen) [attributed]
- (Samuli Virtanen ; [is] ; State Secretary)
- (Samuli Virtanen ; is ; the deputy to Foreign Minister Timo Soini)
- (Samuli Virtanen ; is ; [a] deputy)
- (Timo Soini ; [is] ; Foreign Minister)
- (Timo Soini ; [has] ; [a] deputy)

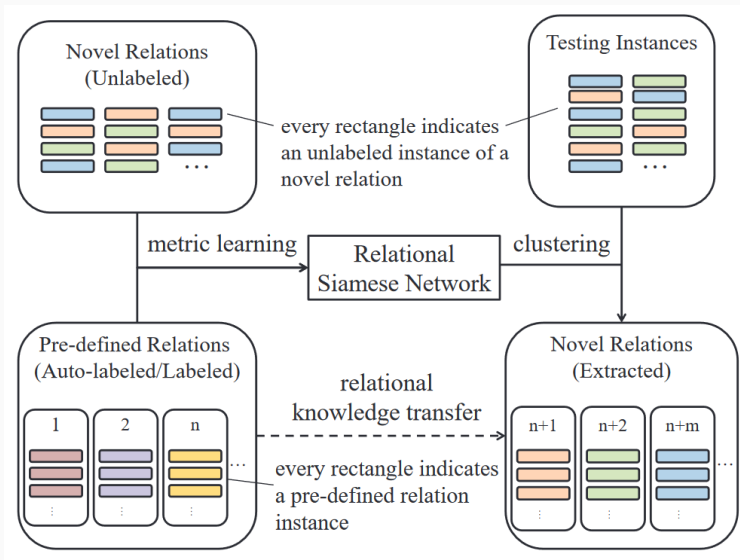
- ▶ Extracting properties of entities: age, colour, size...  
↳ Still triples!
- ▶ Known lists of (some) properties for (some) entity types: and the rest?  
↪ Same as relations, but more often implicit

**Consistency**

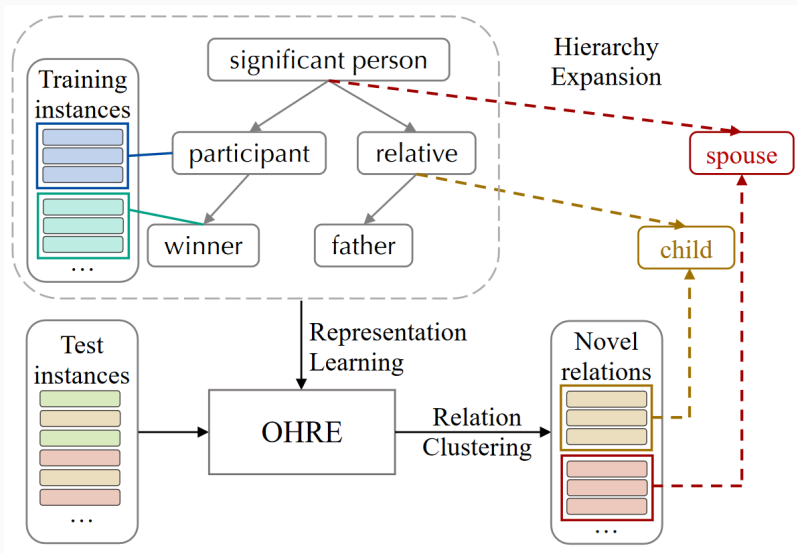
*is married to = has spouse  $\neq$  has wife*

# Relation discovery

*is married to = has spouse  $\neq$  has wife*

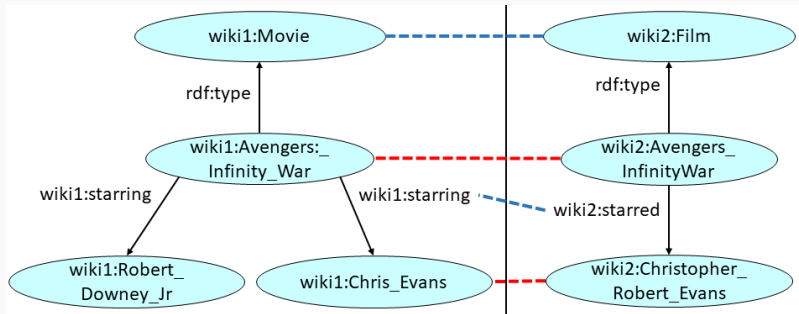


# Relation discovery with a base ontology



# Node deduplication

In case of missed coreference, concurrent entity discovery...



~> Connections with ontology matching (ontology alignment)



# Task interoperability

## Different sets of entities across tasks

- ▶ Most available models for NER: only named, even flat entities
- ▶ Usually in relation extraction: arbitrary noun phrases

# Mention detection and entity recognition

# Mention detection and entity recognition

- ▶ Mention coverage in coreference resolution: non-named, possessives, clauses, unrestricted entity types...  
↪ What about NER, EL, RE?

# Mention detection and entity recognition

- ▶ Mention coverage in coreference resolution: non-named, possessives, clauses, unrestricted entity types...  
↪ What about NER, EL, RE?
- ▶ But: no type, potentially no singleton...

# Mention detection and entity recognition

- ▶ Mention coverage in coreference resolution: non-named, possessives, clauses, unrestricted entity types...

↪ What about NER, EL, RE?

- ▶ But: no type, potentially no singleton...

↪ Heuristics to reconcile?

## Combining tasks

# Named entity recognition & relation extraction

Which shall go first?



Which shall go first?

- ▶ Relation extraction as a classification of entity pairs

Which shall go first?

- ▶ Relation extraction as a classification of entity pairs
- ▶ **Bituyp Jenozs Zoble** was born in Beijing.  
    ↪ Probably Person (or maybe Organization)

# Entity linking & relation extraction

- ▶ Born on 21/12/1977 in Amiens, **Emmanuel** graduated from the *École nationale d'administration* in 2004.

- ▶ Born on 21/12/1977 in Amiens, **Emmanuel** graduated from the *École nationale d'administration* in 2004.
- ▶ Born on 21/12/1977 in Amiens, **Emmanuel** succeeded to François in 2017.

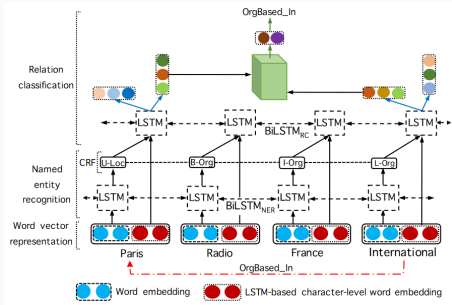
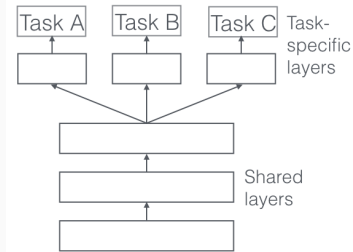
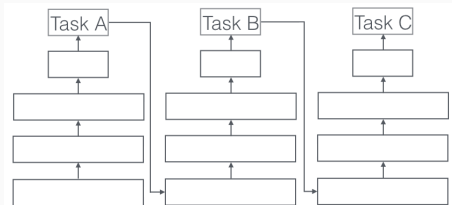


- ▶ **Emmanuel** was born on 21/12/1977 in Amiens. **He** graduated from the *École nationale d'administration* in 2004.

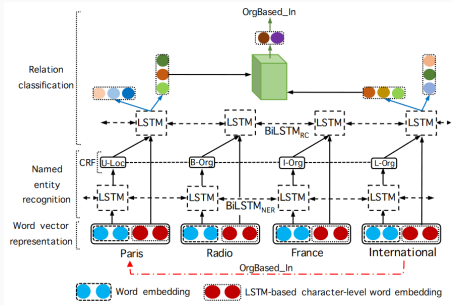
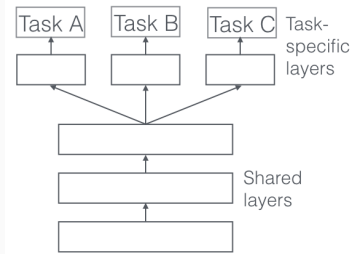
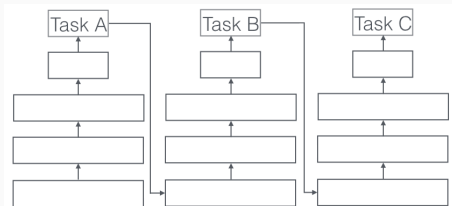
- ▶ **Emmanuel** was born on 21/12/1977 in Amiens. **He** graduated from the *École nationale d'administration* in 2004.
- ▶ Born on 21/12/1977 in Amiens, **Emmanuel** graduated from the *École nationale d'administration* in 2004. That school was abolished in 2021. It was replaced by the INSP. **Macron** succeeded to François Hollande in 2017.



# Interleaving tasks: sequential, joint, multi-task learning



# Interleaving tasks: sequential, joint, multi-task learning



And many others:

- ▶ Auxiliary inputs
- ▶ Auxiliary outputs
- ▶ ... stacking?  
iterations?

# Relation extraction & open IE

Add open IE triples to seeds and retrain RE? Bootstrapping?

Add open IE triples to seeds and retrain RE? Bootstrapping?

- ▶ Often lots of noisy triples: how to filter them?
- ▶ How to match open relations with predefined ones?

*See you next week!*

`first.last@inria.fr`