

What a Beautiful Multilingual World: BabelNet 2.0 & Friends!

Roberto Navigli



Linguistic Computing Laboratory
<http://lcl.uniroma1.it>

TIA 2013, Paris, France

Tiziano Flati

David Jurgens

Andrea Moro

Daniele Vannella

Simone Ponzetto

BabelNet & friends
Roberto Navigli

2

SAYS HI

GETS SALE

WHAT IF MEMES

ARE A NEW FORM OF GLOBAL LANGUAGE

BabelNet & friends
Roberto Navigli

04/11/2013

3

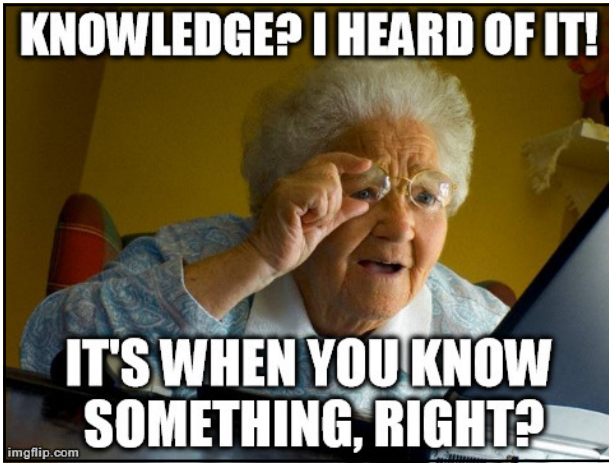
It's all about knowledge!

- Intuitively, we all know what knowledge is...
- ...and why we need it

BabelNet & friends
Roberto Navigli

04/11/2013

4



It's all about knowledge!

- But can we expect computers to *know*?
- Can't computers just use, e.g., *statistical techniques*?

KNOWLEDGE?

MY COMPUTER DOESN'T NEED IT!

BabelNet & Roberto Navigli

6

State-of-the-art Machine Translation

EN: I love chocolate, so I bought a bar in the supermarket.



IT: Amo il cioccolato, così ho comprato un bar in un supermercato.

FR: J'aime le chocolat, donc j'ai acheté un bar dans un supermarché.




BabelNet & friends Roberto Navigli

04/11/2013

7

State-of-the-art Machine Translation

- EN: These are movies in which the music genre, e.g. **rock**, is an important element but not necessarily central to the plot. Examples are Easy Rider (1969), The Graduate (1969), and Saturday Night Fever (1978).



BabelNet & friends Roberto Navigli

04/11/2013

8

State-of-the-art Machine Translation

- EN: These are movies in which the music genre, e.g. **rock**, is an important element but not necessarily central to the plot. Examples are Easy Rider (1969), The Graduate (1969), and Saturday Night Fever (1978).
- IT: Questi sono i film in cui il genere musicale, ad es **roccia**, è un elemento importante, ma non necessariamente al centro della trama.



State-of-the-art Machine Translation

- EN: Knowledge of the distribution of underground **rock** densities can assist in interpreting subsurface geologic structure and rock type.

Danger here!



State-of-the-art Machine Translation

- EN: Knowledge of the distribution of underground **rock** densities can assist in interpreting subsurface geologic structure and rock type.
- IT: La conoscenza della distribuzione di densità di **rock underground** può aiutare a interpretare in sottosuolo struttura geologica e tipo di roccia.



It's not that the "big data" approach is bad,
it's just that mere statistics is not enough



The Knowledge Acquisition Bottleneck

- Knowledge is crucial in language-related research areas
 - Word Sense Disambiguation
 - Named Entity Recognition/Linking
 - Information Extraction
 - (your favourite area here)
- However, providing knowledge is difficult and costly



On a large scale, I mean

AKA: The Hamster Wheel



04/11/2013

13

Resources to the rescue!

- Various projects undertaken to make lexical knowledge available in machine readable form
 - WordNet [Fellbaum, 1998]
 - Open Mind Word Expert [Chklovski & Mihalcea, 2002]
 - EuroWordNet [Vossen, 1998]
 - Multilingual Central Repository [Atserias et al. 2004]
 - The WordNetPlus project [Boyd-Graber et al., 2006]
 - OntoNotes [Hovy et al., 2006]
 - Wikipedia (collaborative effort)
 - Wiktionary (collaborative effort)
 - Omega Wiki (collaborative effort)
 - ...



Wisdom of the Crowd

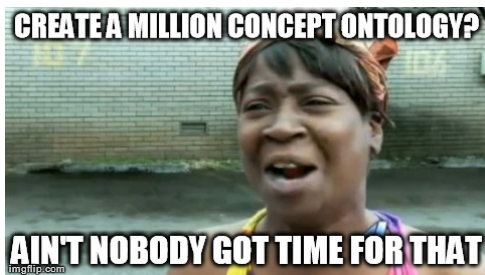
BabelNet & friends
Roberto Navigli

04/11/2013

14

But we need an ontology, not just an encyclopedia!

- And, ideally, we need it to be large-scale, wide-coverage

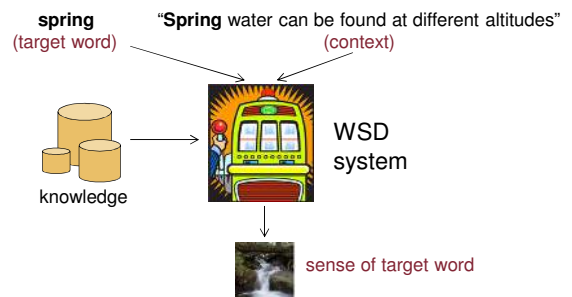


BabelNet & friends
Roberto Navigli

04/11/2013

15

Word Sense Disambiguation in a Nutshell



BabelNet & friends
Roberto Navigli

04/11/2013

16

The Richer, The Better



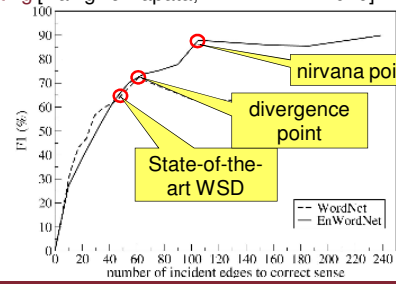
BabelNet & friends
Roberto Navigli

04/11/2013

17

The Richer, The Better

- Highly-interconnected semantic networks have a **great impact** on knowledge-based WSD even in a **fine-grained setting** [Navigli & Lapata, IEEE TPAMI 2010]



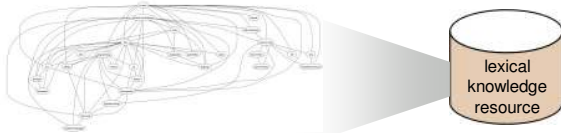
BabelNet & friends
Roberto Navigli

04/11/2013

18

State of the Art "in a nutshell"

- Supervised approaches
 - Require **large amounts** of training data
 - Do **not generalize** across domains and languages
- Knowledge-based approaches have a **higher potential**
 - Lexical knowledge resources **only partly available**



BabelNet & friends
Roberto Navigli

04/11/2013

19

State of the Art "in a nutshell"

- Supervised approaches
 - Require **large amounts** of training data
 - Do **not generalize** across domains and languages
- Knowledge-based approaches have a **higher potential**
 - Lexical knowledge resources **only partly available**
 - Only for **few languages** (e.g. not all 23 EU official languages)
 - Heterogenous** and with **low coverage**



BabelNet & friends
Roberto Navigli

04/11/2013

20



This is where the ERC (and my project) comes into play



erc **MULTILINGUAL JOINT WORD SENSE DISAMBIGUATION**
Multilingual Joint word sense Disambiguation

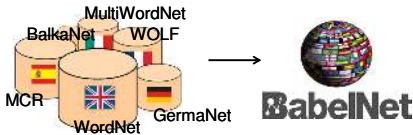
A 5-year ERC Starting Grant (2011-2016)
on Multilingual Word Sense Disambiguation

BabelNet & friends
Roberto Navigli

04/11/2013 23

Multilingual Joint Word Sense Disambiguation
(MultiJEDI)

Key Objective 1: create knowledge for all languages



BalkaNet WOLF
MCR WordNet GermaNet **BabelNet**

BabelNet & friends
Roberto Navigli

04/11/2013 24

Multilingual Joint Word Sense Disambiguation (MultiJEDI)

Key Objective 2: use all languages to disambiguate one

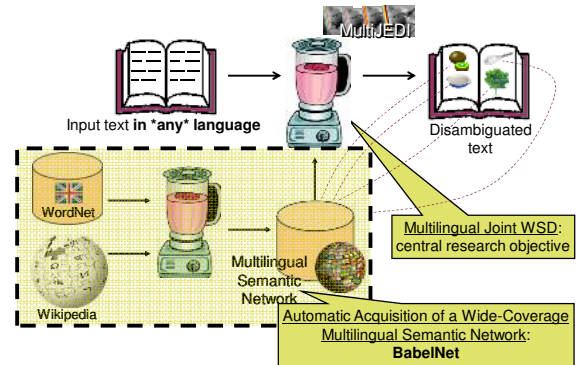


BabelNet & friends
Roberto Navigli

04/11/2013

25

The Vision



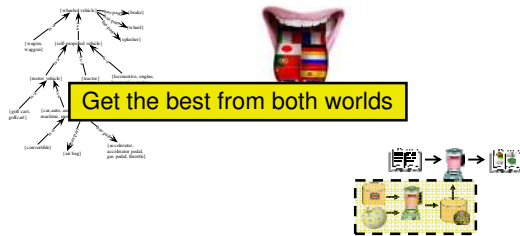
BabelNet & friends
Roberto Navigli

04/11/2013

25

Objective 1: Creating a Multilingual Semantic Network

- Start from two large **complementary** resources:
 - WordNet: full-fledged taxonomy
 - Wikipedia: multilingual and continuously updated



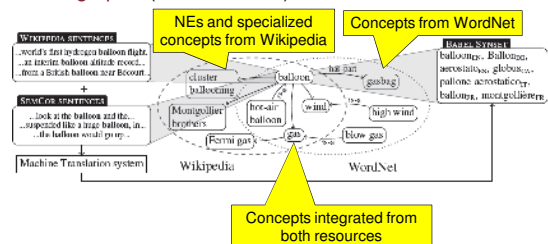
BabelNet & friends
Roberto Navigli

04/11/2013

27

BabelNet [Navigli and Ponzetto, AIJ 2012]

- A wide-coverage multilingual semantic network including both **encyclopedic** (from Wikipedia) and **lexicographic** (from WordNet) entries



BabelNet & friends
Roberto Navigli

04/11/2013

28

BabelNet integrates the best of both worlds

WordNet

Wikipedia

balloon

- n** balloon (large tough corrugated bag filled with gas or heated air)
- n** balloon (small thin inflatable rubber bag with narrow neck)

Speech balloon

BALLOON

04/11/2013 29

Taming the long tail...

Number of results

← more generic → more specific →

BabelNet & friends

Roberto Navigli

04/11/2013 30

WordNet [Miller et al., 1990; Fellbaum, 1998]

semantic relation

concepts

BabelNet & friends

Roberto Navigli

04/11/2013 31

Wikipedia [The Web Community, 2001-today]

(unspecified) semantic relation

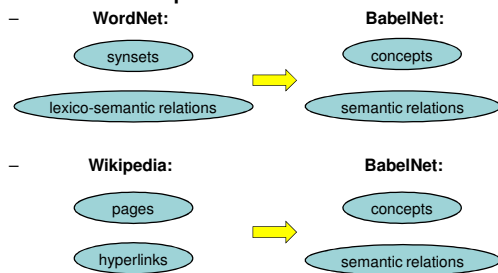
concepts

BabelNet & friends

Roberto Navigli

BabelNet: concepts and semantic relations (1)

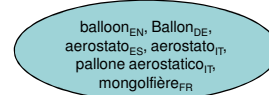
- Concepts and relations in BabelNet are harvested from **WordNet** and **Wikipedia**:



BabelNet: concepts and semantic relations (2)

- We encode knowledge as a **labeled directed graph**:

- Each vertex is a **Babel synset**

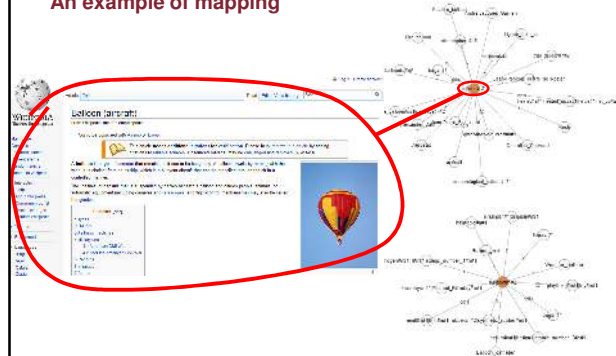


- Each edge is a **semantic relation** between synsets:
 - **is-a** (balloon is-a aircraft)
 - **part-of** (gasbag part-of balloon)
 - **instance-of** (Einstein instance-of physicist)
 - ...
 - **unspecified/relatedness** (balloon related-to flight)

BabelNet: objectives

1. Provide a **unified resource**
 - By establishing an **automated mapping** between Wikipedia pages and WordNet senses
2. Enable **multilinguality**
 - By collecting the **lexicalizations** of concepts in different languages using:
 - a) Wikipedia interlanguage links
 - b) Statistical Machine Translation

An example of mapping



Creation of the Wikipedia disambiguation contexts

Balloon (aircraft)

This article needs additional citations for verification. Please help improve this article by adding citations to reliable sources. Unsourced material may be challenged and removed. [Learn how and when to remove this template message](#).

A **balloon** is a type of aerostat that remains aloft due to buoyancy. A balloon floats by rising with the wind. It is inflated with an inert gas, such as helium, or hot air, and can be steered through the air via a controlled means.

The "load" or payload that is suspended by cables from a basket and three people, basket, or basket or equipment (including camera and telescope, and flight control mechanisms) may also be called a **payload**.

Categories: [Balloons \(aircraft\)](#) | [Aerostats](#) | [Aircraft technology](#) | [Hydrogen technologies](#) | [Aeronautics](#)

$ctx(\text{Balloon (aircraft)}) = \{ \}$

BabelNet & friends
Roberto Navigli
04/11/2013
37

Creation of the Wikipedia disambiguation contexts

Balloon (aircraft)

This article needs additional citations for verification. Please help improve this article by adding citations to reliable sources. Unsourced material may be challenged and removed. [Learn how and when to remove this template message](#).

A **balloon** is a type of aerostat that remains aloft due to buoyancy. A balloon floats by rising with the wind. It is inflated with an inert gas, such as helium, or hot air, and can be steered through the air via a controlled means.

The "load" or payload that is suspended by cables from a basket and three people, basket, or basket or equipment (including camera and telescope, and flight control mechanisms) may also be called a **payload**.

Categories: [Balloons \(aircraft\)](#) | [Aerostats](#) | [Aircraft technology](#) | [Hydrogen technologies](#) | [Aeronautics](#)

$ctx(\text{Balloon (aircraft)}) = \{ \text{aircraft} \}$

BabelNet & friends
Roberto Navigli
04/11/2013
38

Creation of the Wikipedia disambiguation contexts

Balloon (aircraft)

This article needs additional citations for verification. Please help improve this article by adding citations to reliable sources. Unsourced material may be challenged and removed. [Learn how and when to remove this template message](#).

A **balloon** is a type of aerostat that remains aloft due to buoyancy. A balloon floats by rising with the wind. It is inflated with an inert gas, such as helium, or hot air, and can be steered through the air via a controlled means.

The "load" or payload that is suspended by cables from a basket and three people, basket, or basket or equipment (including camera and telescope, and flight control mechanisms) may also be called a **payload**.

Categories: [Balloons \(aircraft\)](#) | [Aerostats](#) | [Aircraft technology](#) | [Hydrogen technologies](#) | [Aeronautics](#)

$ctx(\text{Balloon (aircraft)}) = \{ \text{aircraft, aerostat, buoyancy, airship, ..., gondola} \}$

BabelNet & friends
Roberto Navigli
04/11/2013
39

Creation of the Wikipedia disambiguation contexts

Balloon (aircraft)

This article needs additional citations for verification. Please help improve this article by adding citations to reliable sources. Unsourced material may be challenged and removed. [Learn how and when to remove this template message](#).

A **balloon** is a type of aerostat that remains aloft due to buoyancy. A balloon floats by rising with the wind. It is inflated with an inert gas, such as helium, or hot air, and can be steered through the air via a controlled means.

The "load" or payload that is suspended by cables from a basket and three people, basket, or basket or equipment (including camera and telescope, and flight control mechanisms) may also be called a **payload**.

Categories: [Balloons \(aircraft\)](#) | [Aerostats](#) | [Aircraft technology](#) | [Hydrogen technologies](#) | [Aeronautics](#)

$ctx(\text{Balloon (aircraft)}) = \{ \text{aircraft, aerostat, buoyancy, airship, ..., gondola, ballooning, hydrogen, aeronautics} \}$

BabelNet & friends
Roberto Navigli
04/11/2013
40

Building BabelNet: Mapping Wikipedia to WordNet

$$\mu(w) = \underset{s \in \text{Senses}_{\text{WN}}(w)}{\text{argmax}} p(s|w) = \underset{s}{\text{argmax}} \frac{p(s, w)}{p(w)}$$

$$= \underset{s}{\text{argmax}} p(s, w)$$

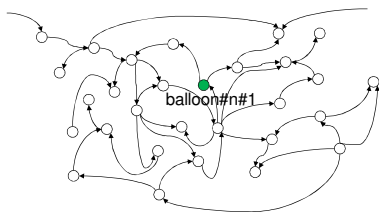
- Given a Wikipage w and its disambiguation context $\text{ctx}(w)$:
 - For each WordNet sense s of w , calculate $\text{score}(s, w)$ as follows:

$$\text{score}(s, w) = \sum_{c \in \text{Ctx}(w)} \sum_{s' \in \text{Senses}_{\text{WN}}(c)} \sum_{p \in \text{paths}_{\text{WN}}(s, s')} e^{-(\text{length}(p)-1)}$$

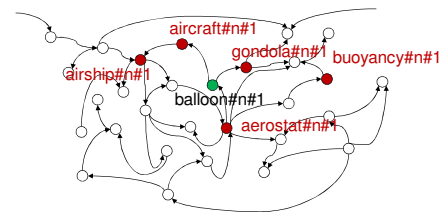


The Wikipedia page context in the WordNet graph

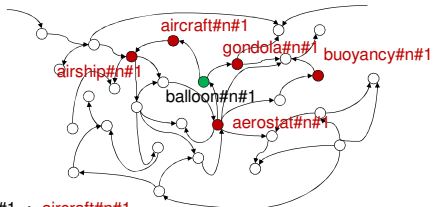
$\text{ctx}(\text{Balloon (aircraft)}) = \{ \text{aircraft, aerostat, buoyancy, airship, ..., gondola} \}$



The Wikipedia page context in the WordNet graph

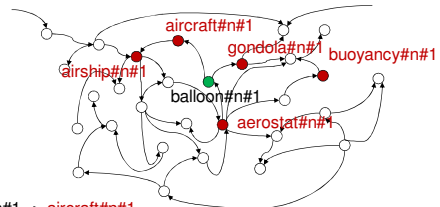


The Wikipedia page context in the WordNet graph



balloon#n#1 -> aircraft#n#1
 balloon#n#1 -> aircraft#n#1 -> airship#n#1
 balloon#n#1 -> gondola#n#1
 balloon#n#1 -> gondola#n#1 -> flight#n#1 -> buoyancy#n#1
 balloon#n#1 -> aerostat#n#1

The Wikipedia page context in the WordNet graph



balloon#n#1 -> aircraft#n#1
 balloon#n#1 -> aircraft#n#1 -> airship#n#1 \rightarrow 0.35
 balloon#n#1 -> gondola#n#1
 balloon#n#1 -> gondola#n#1 -> flight#n#1 -> buoyancy#n#1
 balloon#n#1 -> aerostat#n#1

Building BabelNet: Translating Babel synsets

1. Exploiting Wikipedia interlanguage links

Building BabelNet: Translating Babel synsets

2. Filling the lexical translation gaps using a Machine Translation system to translate the English lexicalizations of a concept

On August 27, 1783 in Paris, Franklin witnessed the world's first hydrogen **(Balloon (aircraft)|balloon)** flight.



Le 27 Août, 1783 à Paris, Franklin vu le premier vol en **ballon** d'hydrogène.

Building BabelNet: Translating Babel synsets

2. Filling the lexical translation gaps using a Machine Translation system to translate the English lexicalizations of a concept

- For each word sense *s*, we translate:
 - sentences from **SemCor** (a corpus annotated with WordNet senses) which contain *s*
 - sentences from **Wikipedia** linked to the Wikipage of *s*
- The most frequent translation of *s* is selected for each target language



The most frequent translation of a word in a given meaning

left context	term	right context
	wikification	may refer to: the...
geoinformatics services' and '	wikification	of GIS by the masses'
the process may be called	wikification	(as in ...
which is then called "	wikification	and to the related problem
reason needs copyediting,	wikification	, reduction of POV, work on references
huge amount of cleanup,	wikification	, etc. Version of 12 Nov

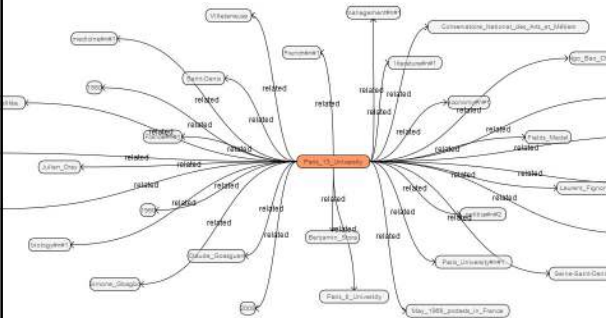
The most frequent translation of a word in a given meaning

left context	term	right context
	wikificazione	potrebbe riferirsi a: il...
servizi geoinformatici' e '	wikification	di GIS dalle masse'
il processo chiamato	wikificazione	(come in ...
che è quindi chiamato	wikificazione	e al problema correlato...
ragione richiede copyediting,	wikification	, riduzione di POV, lavoro su reference
grandi quantità di pulizia,	wikificazione	, ecc. Versione del 12 Novembre

The most frequent translation of a word in a given meaning

left context	term	right context
	wikificazione	potrebbe riferirsi a: il...
servizi geoinformatici' e '	wikification	di GIS dalle masse'
il processo chiamato	wikificazione	(come in ...
che è quindi chiamato	wikificazione	e al problema correlato...
ragione richiede copyediting,	wikification	, riduzione di POV, lavoro su reference
grandi quantità di pulizia,	wikificazione	, ecc. Versione del 12 Novembre

BabelNet knows Paris 13!



The BabelNet API

```

BabelNet or = BabelNet.getInstance();
System.out.println("SYNSETS WITH English words: 'bank'");
List<BabelSynset> synsets = tr.getSynsets(Language.EN, "bank");
for (BabelSynset synset : synsets)
{
    System.out.println("ID: " + synset.getID() + " | SOURCE: " + synset.getSource() +
        " | MESH: " + synset.getMeshOffset() + " | " +
        " | IRI: " + synset.getIRI() + " | " +
        " | URI: " + synset.getURI() + " | " +
        " | NAME: " + synset.getCanonicalName() + " | " +
        " | " + synset.getCanonicalName());
    for (BabelSynset relsynset : synset.getRelatedSynsets())
    {
        System.out.println("----");
        System.out.println("Relation: " + relsynset.getRelationType().getIRI());
        for (BabelSynset relsynset : relsynset.getRelatedSynsets())
        {
            System.out.println("ID: " + relsynset.getID() +
                " | " + relsynset.getURI() +
                " | " + relsynset.getCanonicalName());
        }
    }
}
System.out.println();
    
```

Retrieve all synsets with the English lemma "bank"

Print information about each synset

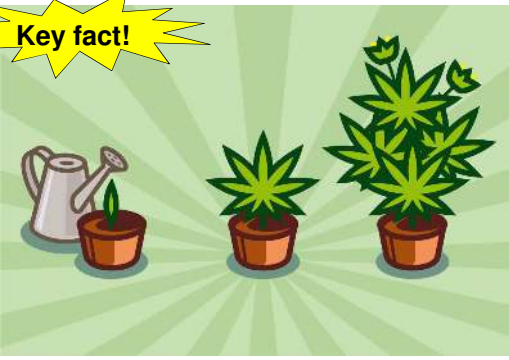
Get the (relation, synsets) map of the synset neighbours

Get the synsets related by a given relation type

Print the information of each related synset

BabelNet goes at a faster pace than I can cope with

Key fact!



Anatomy of BabelNet 2.0

Previous version had 6!

- 50 languages covered (including Latin!)
- List at <http://babelnet.org/stats.jsp>

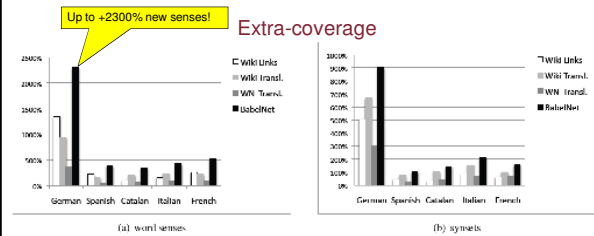


WordNet-Wikipedia mapping accuracy

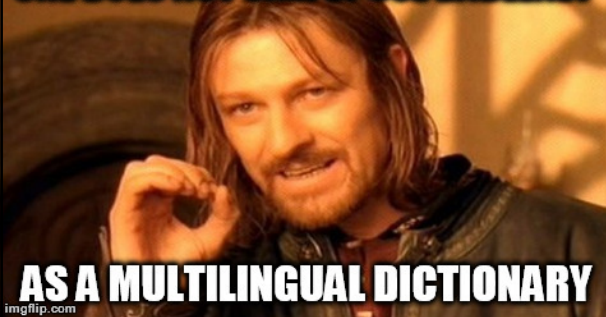
- Overall quality of the mapping: ~84%
 - On a random sample of 1k Wikipages
 - Note: this concerns only those 50k synsets in the intersection
- Quality of the mapping of frequent words: ~91%



Evaluation of BabelNet against gold standard resources



ONE DOES NOT SIMPLY USE BABELNET



Coarse-grained Word Sense Disambiguation with BabelNet

Resource	Algorithm	Nouns only	All words
		P/R/F ₁	P/R/F ₁
WordNet	Degree	80.1	79.7
	PLength	80.2	79.8
	SProbability	79.8	79.3
	PageRank	79.9	79.3
BabelNet	Degree	84.7	82.3
	PLength	85.4	82.7
	SProbability	84.6	82.1
	PageRank	82.3	80.1
	SUSSX-FR	81.1	77.0
	TreeMatch	N/A	73.6
	NUS-PT	82.3	82.5
SSI	84.1	83.2	
MI'S BL	77.4	78.9	
Random BL	63.5	62.7	

State of the art results!

Current state of the art

Home Task Description Participate! Data Results

Multilingual Word Sense Disambiguation

Multilingual Word Sense Disambiguation

This is the webpage for the SemEval-2013 task on multilingual Word Sense Disambiguation. Here you can find [information about the task](#), find out [how to participate](#), and [download the data](#).

News:

- The competition starts on **March 5** (when the test set will be available) and ends on **March 15**!
- March 3, 2013: [BabelNet 1.1.1](#) released for gold standard task sense inventory ([download](#)).
- July 2012: Trial data are out! You can find them [here](#).
- November 2012: Availability of a training dataset. Please note that **no training data will be released for this task**, since we do not assume the availability of any pre-existing labeled data for multilingual disambiguation - cf. also the [task description](#). However, task participants are free to explore any WSD framework they want. This includes both unsupervised (e.g.

Contact Info

Organizers
 Roberto Navigli
 Sapientza University of Rome, Italy
 David A. Jurgens
 Sapientza University of Rome, Italy

Other Info

Announcements

- March 3, 2013: [Callabout 1.1.1](#) Released for Task Sense Inventory.
- July 2012: Trial data are out! You can find them [here](#).

BabelNet & friends
 Roberto Navigli
 04/11/2013
 73

Key fact! Annotating with BabelNet: all in one!

- Annotating with **BabelNet** implies annotating with **WordNet** and **Wikipedia**
- (now also **OmegaWiki** and **Open Multilingual WordNet**!)

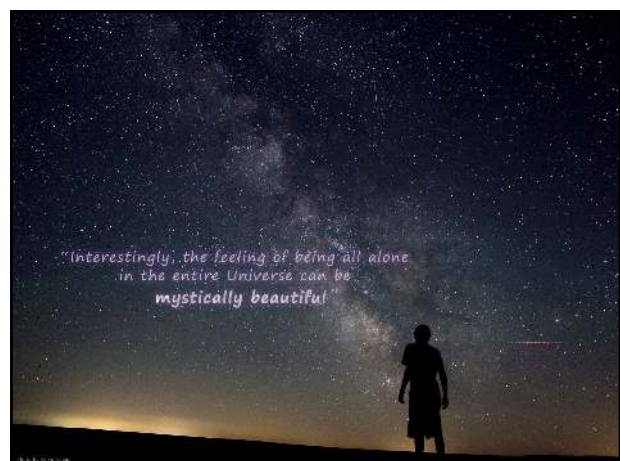


BabelNet & friends
 Roberto Navigli
 04/11/2013
 74

Dataset for Multilingual Word Sense Disambiguation

Language	Instances	Single-words	Multiword expressions	Named Entities
BabelNet				
English	1931	1604	127	200
French	1656	1389	89	176
German	1467	1267	21	176
Italian	1706	1454	211	41
Spanish	1481	1103	129	249
Wikipedia				
English	1242	945	102	195
French	1039	790	72	175
German	1156	957	21	176
Italian	1977	869	85	41
Spanish	1103	758	107	248
WordNet				
English	1644	1502	85	57

BabelNet & friends
 Roberto Navigli
 04/11/2013



We are not alone in the (resource) universe!

- **DBPedia** [Bizer et al. 2009] - a resource obtained from structured information in Wikipedia
 - «Describes 3.77M things»
 - Core of the Linked Open Data Cloud
- **YAGO** [Suchanek et al. 2007]
 - «Contains 10M entities and 120M facts about these entities»
 - Links Wikipedia categories to WordNet synsets
- **MENTA** [de Melo and Weikum, 2010]
 - A «multilingual taxonomy with 5.4M entities»
- **WikiNet** [Nastase and Strube, 2013]
 - Semantic network connecting Wikipedia entities
 - «3M concepts and 38+M relations»
- **Freebase** (<http://freebase.com>): collaborative effort
 - Structured data; started from Wikipedia, MusicBrainz, ChefMoz, etc.

BabelNet & friends
Roberto Navigli

04/11/2013

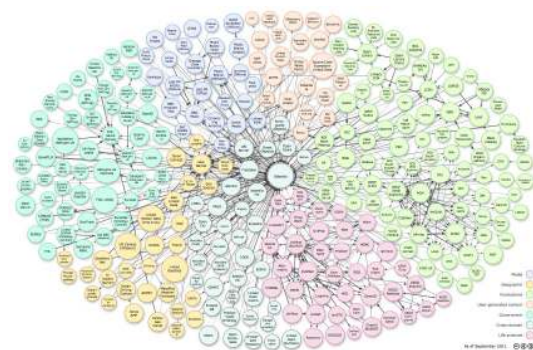
78

So where is the novelty?

- We provide a **unified, integrated** inventory and network for **both word senses and named entities**



Now in the Linked Open Data cloud...

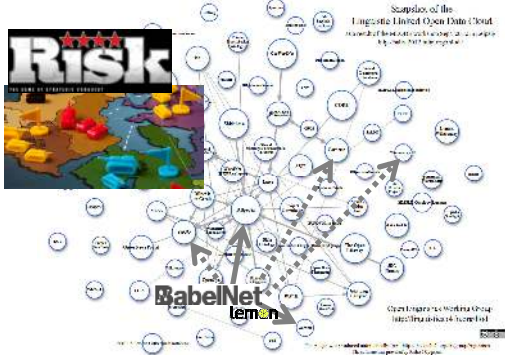


BabelNet goes to the (Multilingual) Semantic Web
Roberto Navigli

04/11/2013

81

Actually, in the *Linguistic* Linked Open Data cloud...



SPred: Semantic Predicates on a large scale [Flati & Navigli, ACL 2013]

- Choose a set of semantic classes
- Given a lexical predicate

$$w_1 w_2 \dots w_i^* w_{i+1} \dots w_n$$

identify the semantic class distribution for *

- Classify new filling arguments

The idea in a nutshell: start from text (1/6)

cup of *

...pounded millet mixed with two cups of butter, every day...
 ...senior citizens would linger over cups of coffee and exchange news...
 ...began over a cup of yogurt and a broken printer that ...
 ...in the quarterfinal cup of Yugoslavia as a renowned team...
 ...which is equivalent to 3 cups of regular coffee, he will be able to...
 He drank several more cups of wine and started behaving wildly...
 ...or fruit, as well as small cups of kosher wine or other beverages...
 ...two Euroleague titles, two Cups of Italy, one Korać Cup and one...

The idea in a nutshell: focus on lexical predicates (2/6)

cup of *

...pounded millet mixed with two **cups of** butter, every day...
 ...senior citizens would linger over **cups of** coffee and exchange news...
 ...began over a **cup of** yogurt and a broken printer that ...
 ...in the quarterfinal **cup of** Yugoslavia as a renowned team...
 ...which is equivalent to 3 **cups of** regular coffee, he will be able to...
 He drank several more **cups of** wine and started behaving wildly...
 ...or fruit, as well as small **cups of** kosher wine or other beverages...
 ...two Euroleague titles, two **Cups of** Italy, one Korać Cup and one...

The idea in a nutshell: focus on filling arguments (3/6)

cup of *

...pounded millet mixed with two cups of butter, every day...
...senior citizens would linger over cups of coffee and exchange news...
...began over a cup of yogurt and a broken printer that ...
...in the quarterfinal cup of Yugoslavia as a renowned team...
...which is equivalent to 3 cups of coffee, he will be able to...
He drank several more cups of wine and started behaving wildly...
...or fruit, as well as small cups of kosher wine or other beverages...
...two Euroleague titles, two Cups of Italy, one Korac Cup and one...

The idea in a nutshell: gather all the filling arguments (4/6)

cup of *



The idea in a nutshell: group similar arguments (5/6)

cup of *




The idea in a nutshell: associate semantic labels (6/6)

cup of *




cup of *




wine¹_n

Wine
Sack
White wine
Red wine
Claret
Kosher wine
Madeira wine
Wine in China
...




coffee¹_n

Coffee
Turkish coffee
Drip coffee
Espresso
Cappuccino
Caffè latte
Decaffeinated coffee
...



herb²_n

Earl Grey tea
Green tea
Indian tea
Black tea
Tea
...



water¹_n

Water
Seawater
...

Classes sorted by relevance!

BabelNet & friends 04/11/2013 90

**Another Plan for Knowledge Acquisition:
Ontology Learning**

- Knowledge acquisition without relying on hand-crafted knowledge
- Deals with **technical domains**
- Current approaches have **limits**:
 - Kozareva & Hovy [2010]: lexico-syntactic patterns + pruning based on the longest path
 - Yang & Callan [2009]: based on incremental clusters of terms
 - 82% F1 for small-scale WordNet is-a sub-hierarchies (39 terms on average)
 - 61% F1 on part-of sub-hierarchies
 - Snow et al. [2006]: hyponym acquisition based on a probabilistic model
 - 58% P, 21% R
- These approaches have **not** been shown to be able to extract large **specialized** domain ontologies

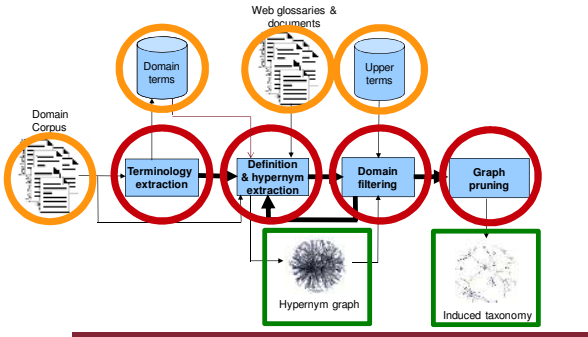
BabelNet & friends 04/11/2013 91

OntoLearn Reloaded
[Navigli, Faralli & Velardi, IJCAI 2011;
Computational Linguistics 2013]

Unlike other approaches, we learn both concepts
and relations **entirely from scratch**
for **any** domain of interest

BabelNet & OntoLearn Reloaded 04/11/2013 92

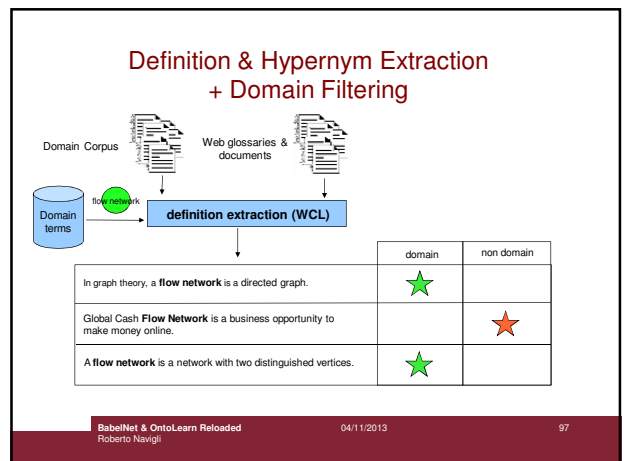
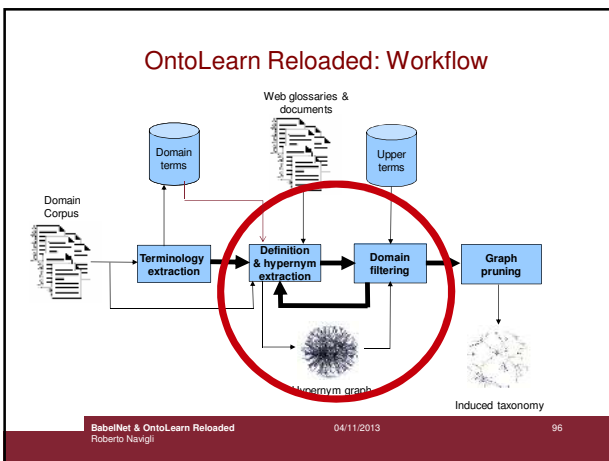
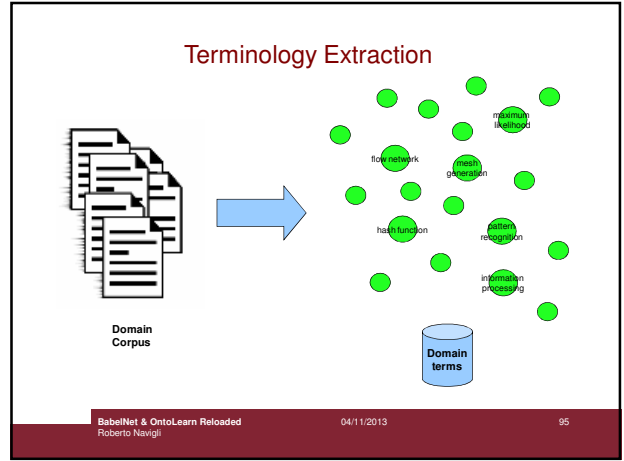
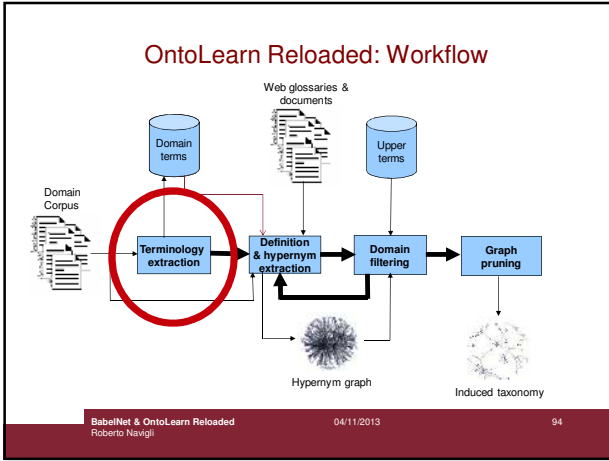
OntoLearn Reloaded: Workflow

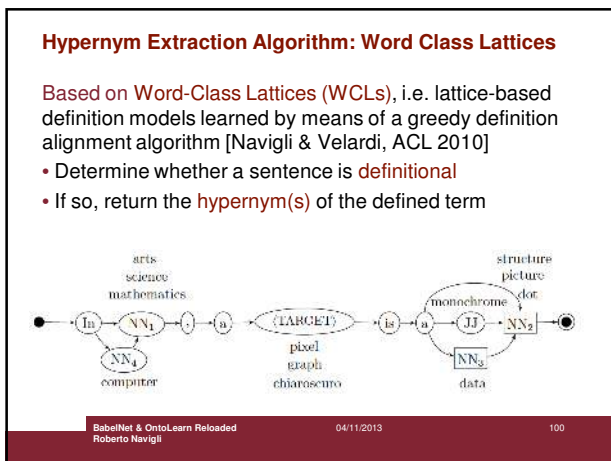
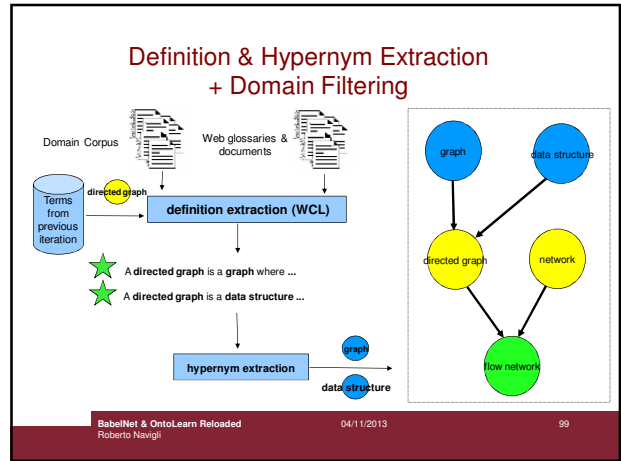
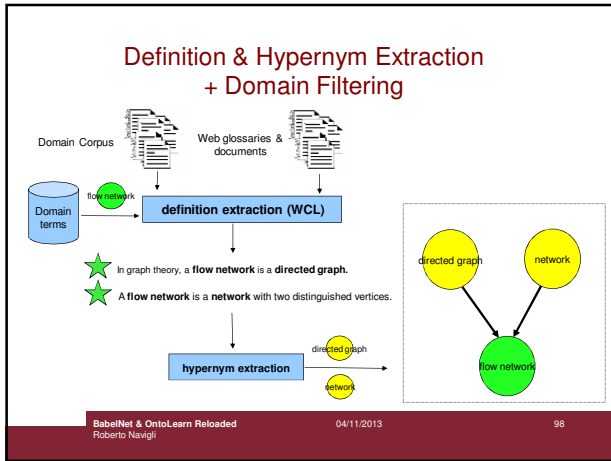


The workflow consists of the following steps:

- Domain Corpus** (Input)
- Terminology extraction** (Produces **Domain terms**)
- Definition & hyponym extraction** (Produces **Web glossaries & documents** and **Upper terms**)
- Domain filtering** (Produces **Hyponym graph**)
- Graph pruning** (Produces **Induced taxonomy**)

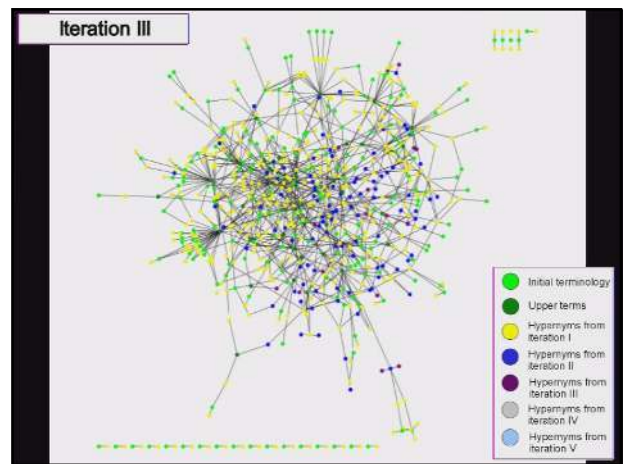
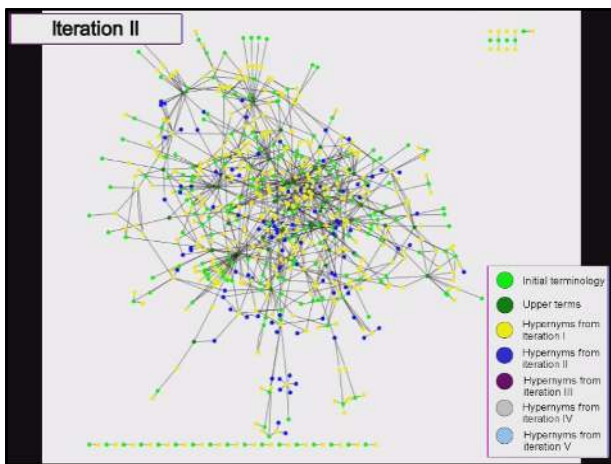
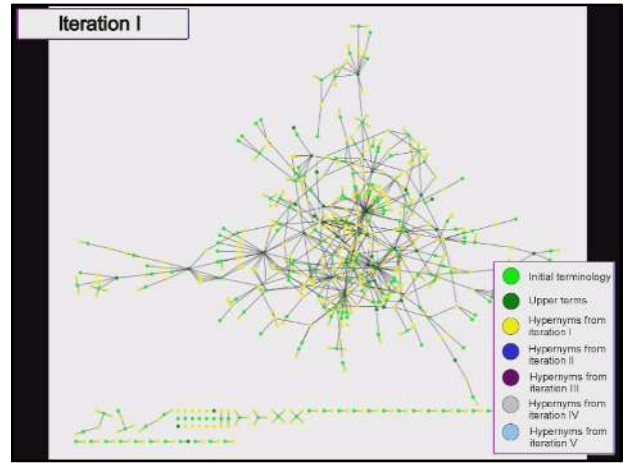
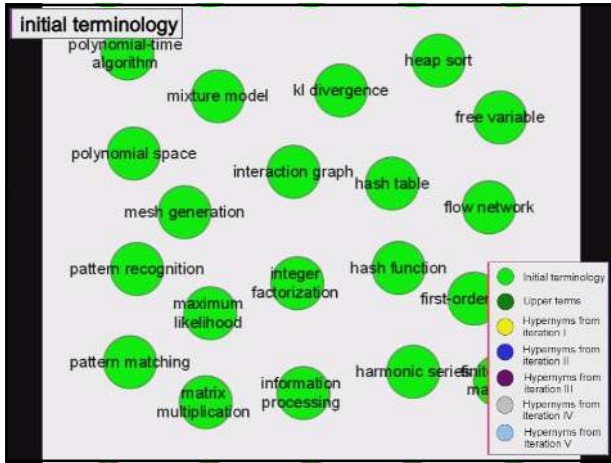
BabelNet & OntoLearn Reloaded 04/11/2013 93

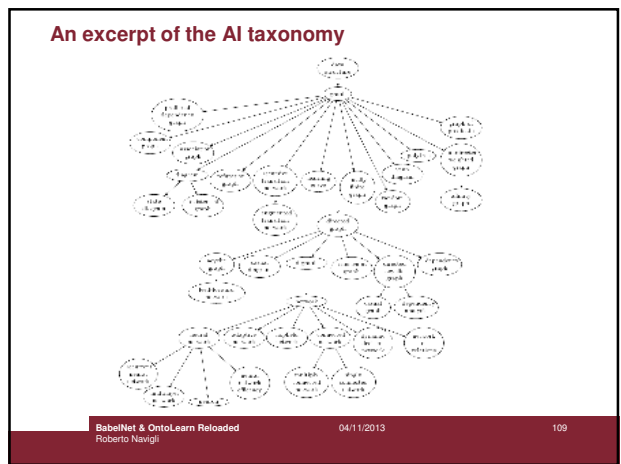
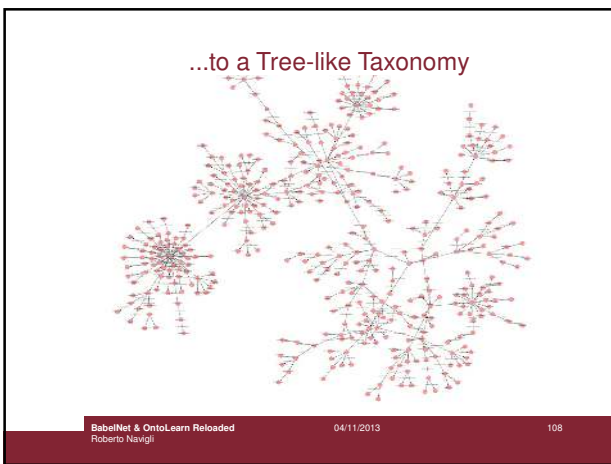
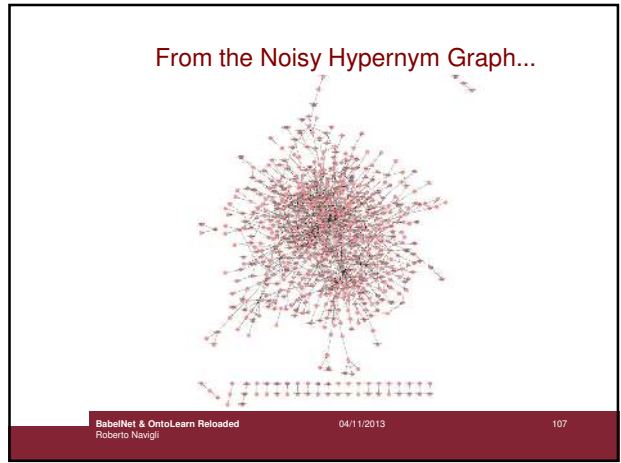
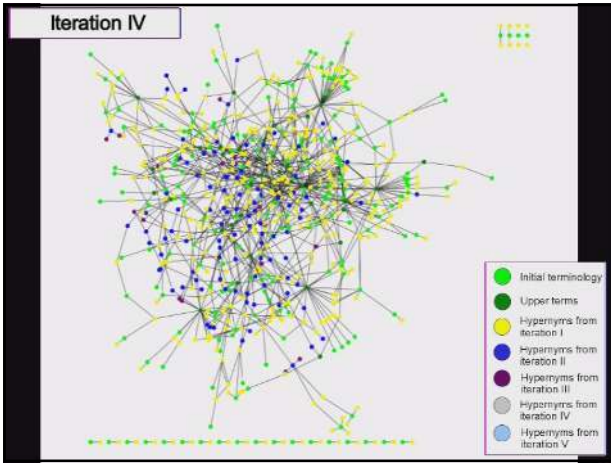




The iterative growth of the hypernym graph

BabelNet & OntoLearn Reloaded 04/11/2013 101
Roberto Navigli





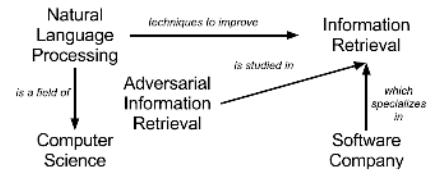
Semantically-Enhanced Open Information Extraction [Moro and Navigli, IJCAI 2013]

Semantic Network	Semantic Relations	Open set of relations
ReVerb	✗	✓
YAGO2	✓	✗
WikiNet	✓	✗
WiSeNet	✓	✓

A. Fader, S. Soderland and O. Etzioni. *Identifying Relations for Open Information Extraction*. In Proc. of EMNLP, 2011.
 J. Hoffart et al. *YAGO2: A spatially and temporally enhanced knowledge base from Wikipedia*. Journal of Artificial Intelligence, 2012.
 V. Nastase and M. Strube. *Transforming Wikipedia into a large scale multilingual concept network*. Journal of Artificial Intelligence, 2012.
 A. Moro and R. Navigli. *WiSeNet: building a wikipedia-based semantic network with ontologized relations*. In Proc. of CIKM, 2012.

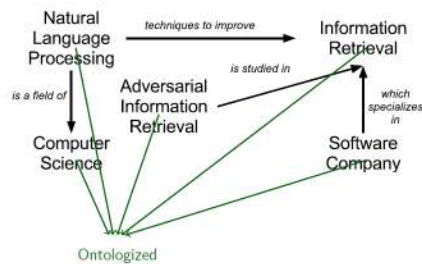
Semantically-Enhanced Open Information Extraction [Moro and Navigli, IJCAI 2013]

- Moving to semantically enhanced OIE:



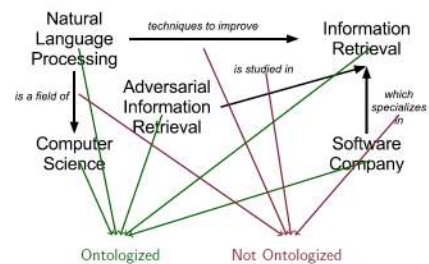
Semantically-Enhanced Open Information Extraction [Moro and Navigli, IJCAI 2013]

- Moving to semantically enhanced OIE:



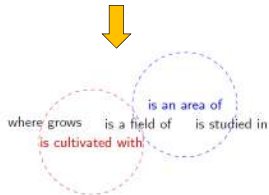
Semantically-Enhanced Open Information Extraction [Moro and Navigli, IJCAI 2013]

- Moving to semantically enhanced OIE:



Ontologizing the relation strings

where grows is cultivated with is an area of is studied in is a field of

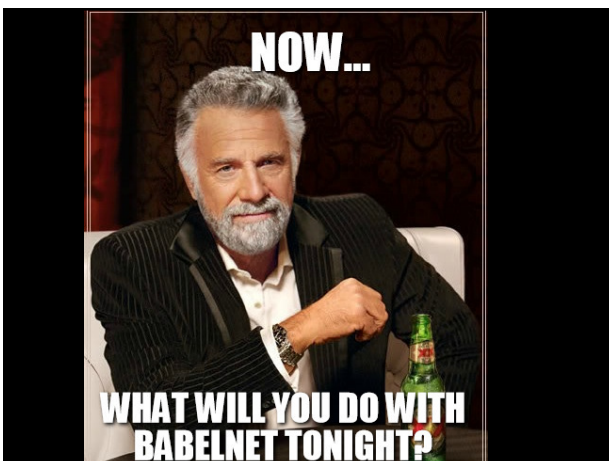
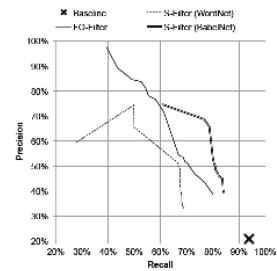


Relation Synset: {where grows, is cultivated with, is a field of}
 Relation Synset: {is studied in, is an area of, is a field of}

Improving traditional IE with BabelNet [Moro et al., ISWC 2013]



- Using rich semantics increases precision while keeping recall high



Thanks or...





SAPIENZA
UNIVERSITÀ DI ROMA

Roberto Navigli

Linguistic Computing Laboratory
<http://lcl.uniroma1.it>

