

# Semantic Representations of Concepts and Entities and their Applications

Jose Camacho-Collados



SAPIENZA  
UNIVERSITÀ DI ROMA

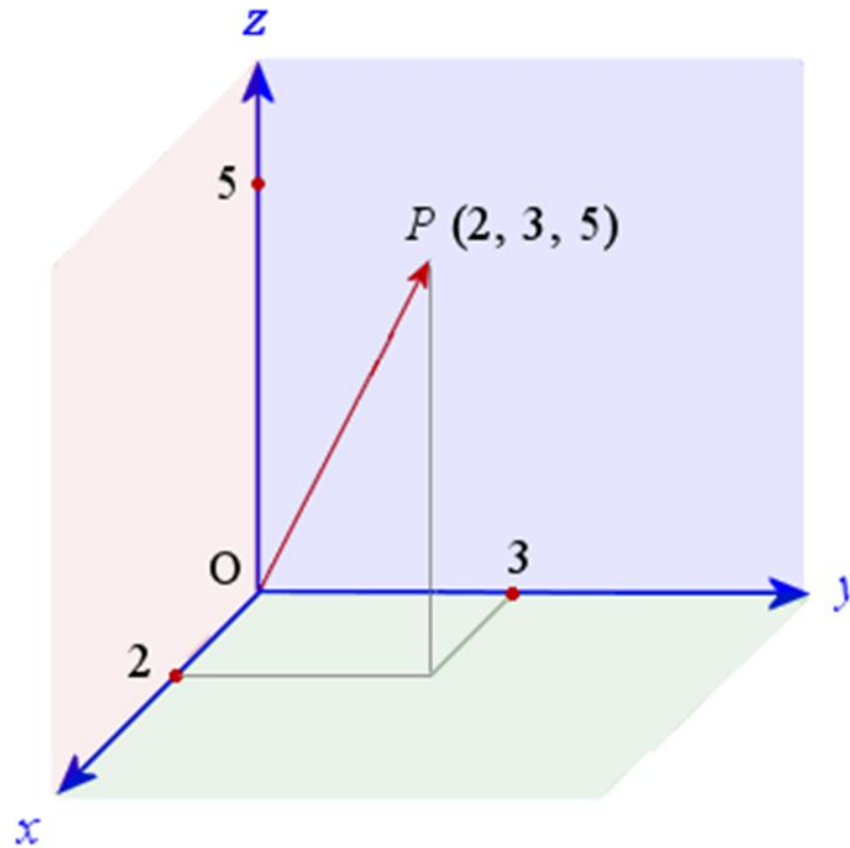
University of Cambridge, 20 April 2017

# Outline

- Background: Vector Space Models
- Semantic representations for Senses, Concepts and Entities -> NASARI
- Applications
- Conclusions

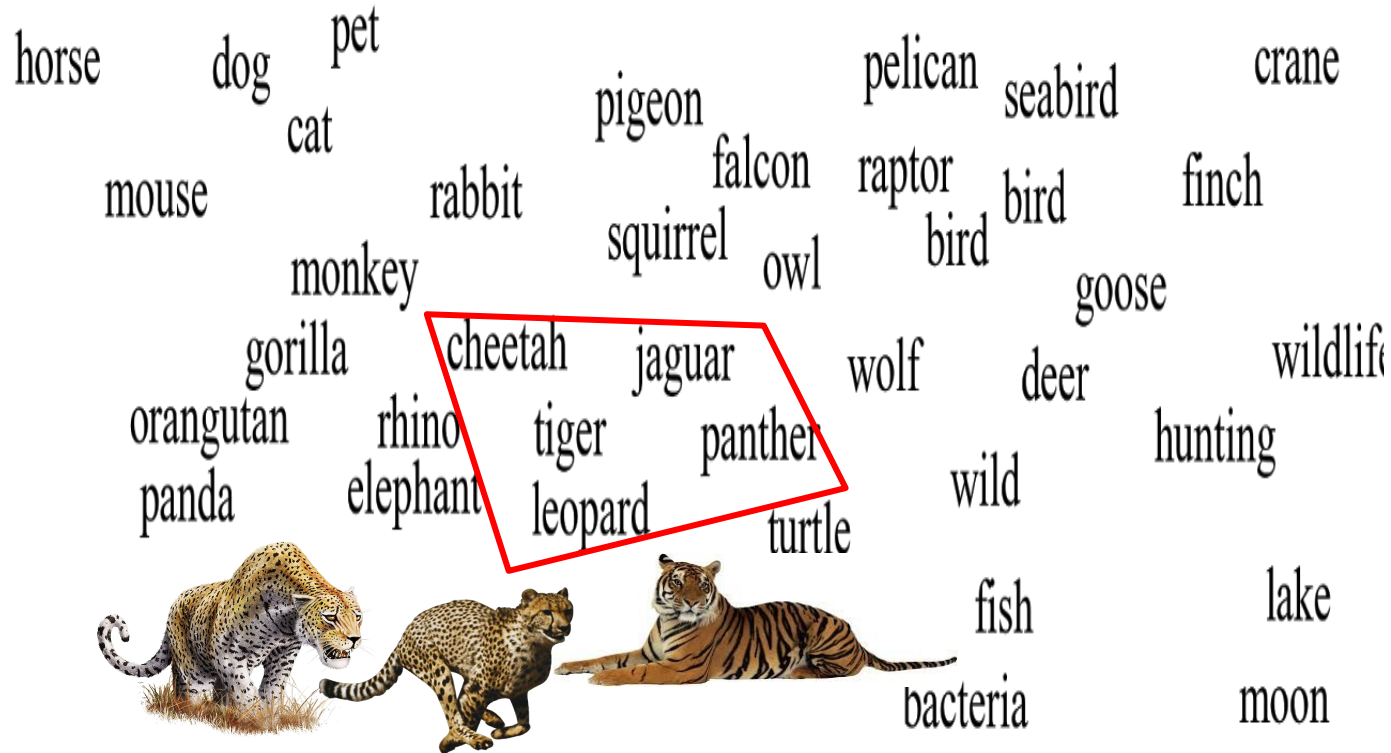
# Vector Space Model

Turney and Pantel (2010): Survey on Vector Space Model of semantics

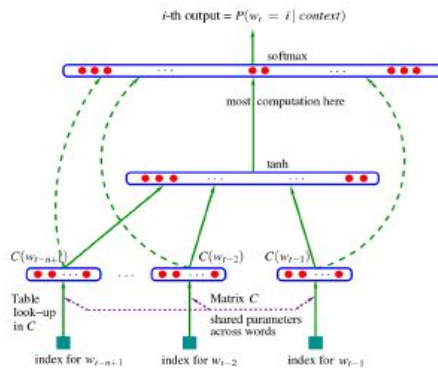


# Word vector space models

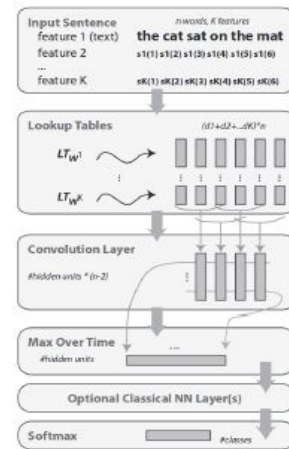
Words are represented as vectors: semantically similar words are close in the vector space



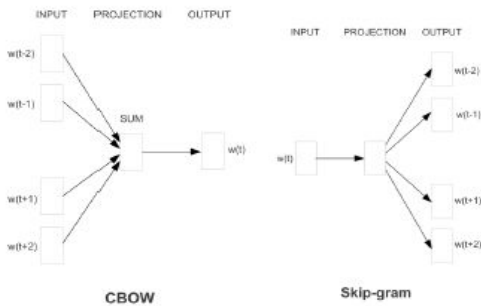
# Neural networks for learning word vector representations from text corpora -> word embeddings



Bengio et al. (2003)



Collobert & Weston (2008)



Mikolov et al. (2013)

Probability and Ratio	$k = \text{solid}$	$k = \text{gas}$	$k = \text{water}$
$P(k \text{ice})$	$1.9 \times 10^{-4}$	$6.6 \times 10^{-5}$	$3.0 \times 10^{-3}$
$P(k \text{steam})$	$2.2 \times 10^{-5}$	$7.8 \times 10^{-4}$	$2.2 \times 10^{-3}$
$P(k \text{ice})/P(k \text{steam})$	8.9	$8.5 \times 10^{-2}$	1.36

Pennington et al. (2014)

# Why word embeddings?

Embedded vector representations:

- are compact and fast to compute
- preserve important relational information between words (actually, meanings):

$$\textit{king} - \textit{man} + \textit{woman} \approx \textit{queen}$$

- are geared towards general use

# Applications for word representations

- Syntactic parsing (Weiss et al. 2015)
- Named Entity Recognition (Guo et al. 2014)
- Question Answering (Bordes et al. 2014)
- Machine Translation (Zou et al. 2013)
- Sentiment Analysis (Socher et al. 2013)

... and many more!

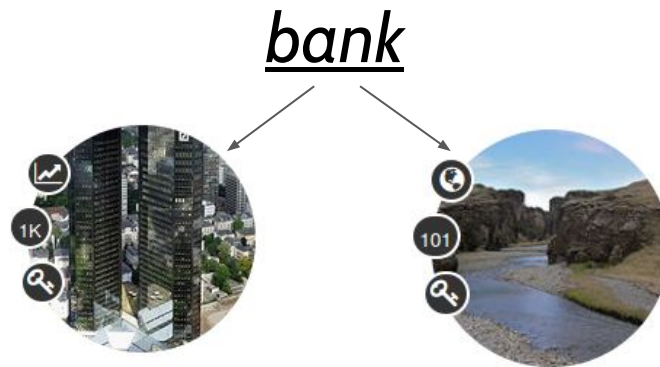
# AI goal: language understanding





# Limitations of word embeddings

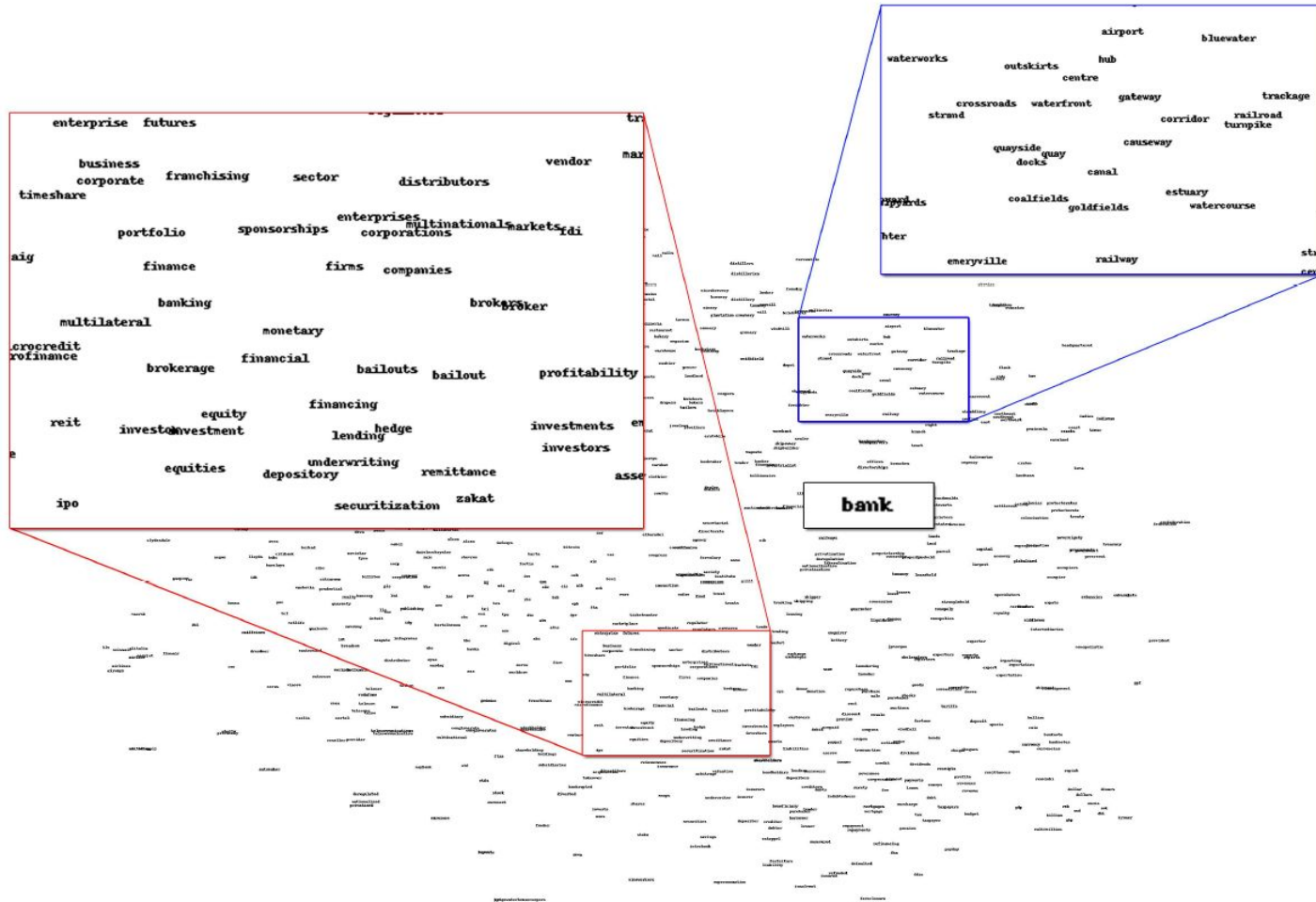
- Word representations cannot capture ambiguity. For instance,





# Problem 1:

word representations cannot capture ambiguity

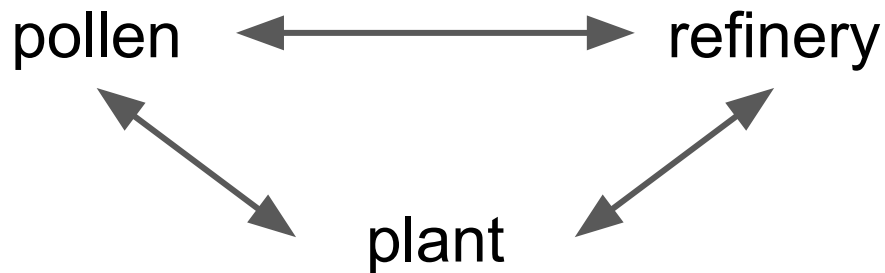




# Word representations and the triangular inequality

Example from Neelakantan et al (2014)

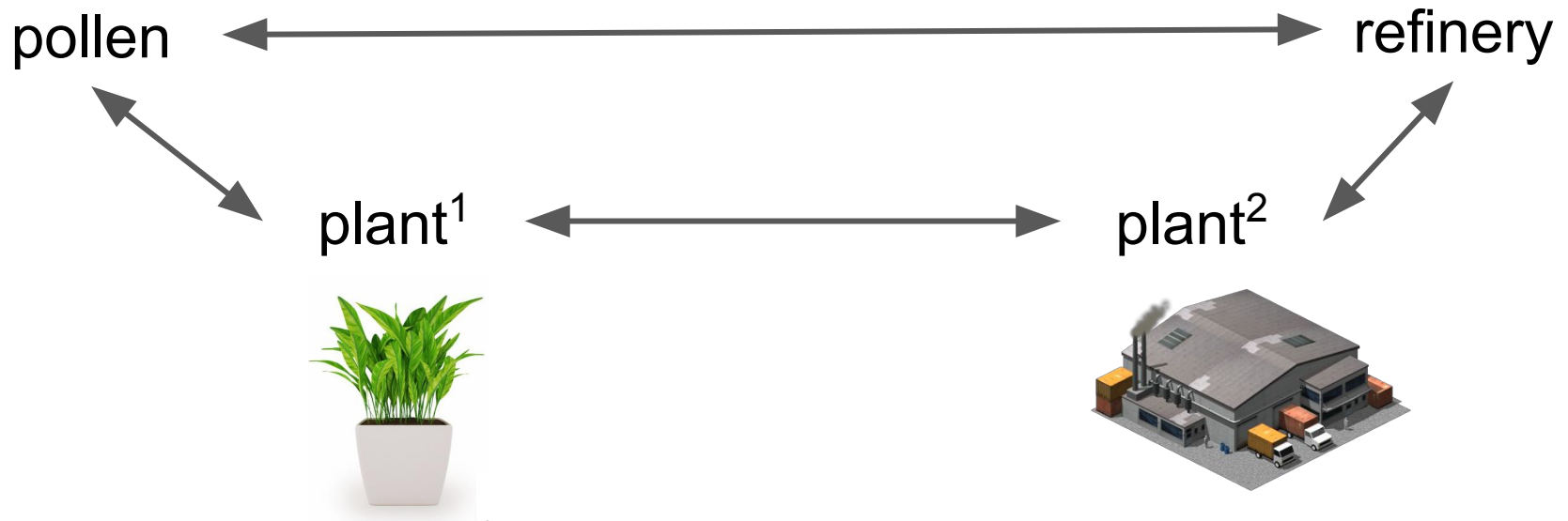
For distance  $d$ ,  $d(a, c) \leq d(a, b) + d(b, c)$ .



# Word representations and the triangular inequality

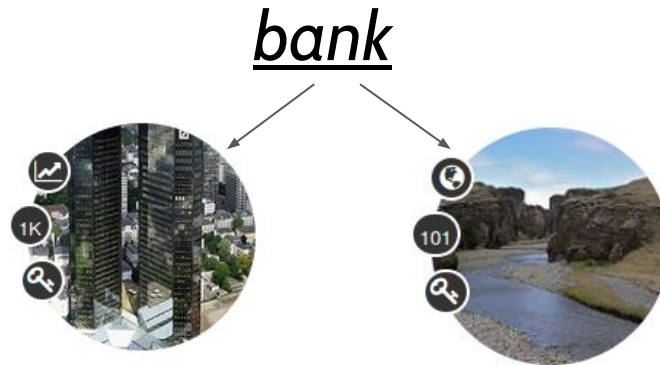
Example from Neelakantan et al (2014)

For distance  $d$ ,  $d(a, c) \leq d(a, b) + d(b, c)$ .



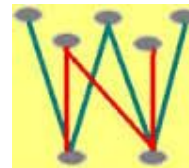
# Limitations of word representations

- They cannot capture ambiguity. For instance,



-> They neglect rare senses and infrequent words

- Word representations do not exploit knowledge from existing lexical resources.





# a Novel Approach to a Semantically-Aware Representations of Items

<http://lcl.uniroma1.it/nasari/>



# NASARI semantic representations

- NASARI 1.0 (April 2015): *Lexical and unified vector representations for WordNet synsets and Wikipedia pages for English.*

José Camacho Collados, Mohammad Taher Pilehvar and Roberto Navigli. *NASARI: a Novel Approach to a Semantically-Aware Representation of Items*. **NAACL 2015**, Denver, USA, pp. 567-577.

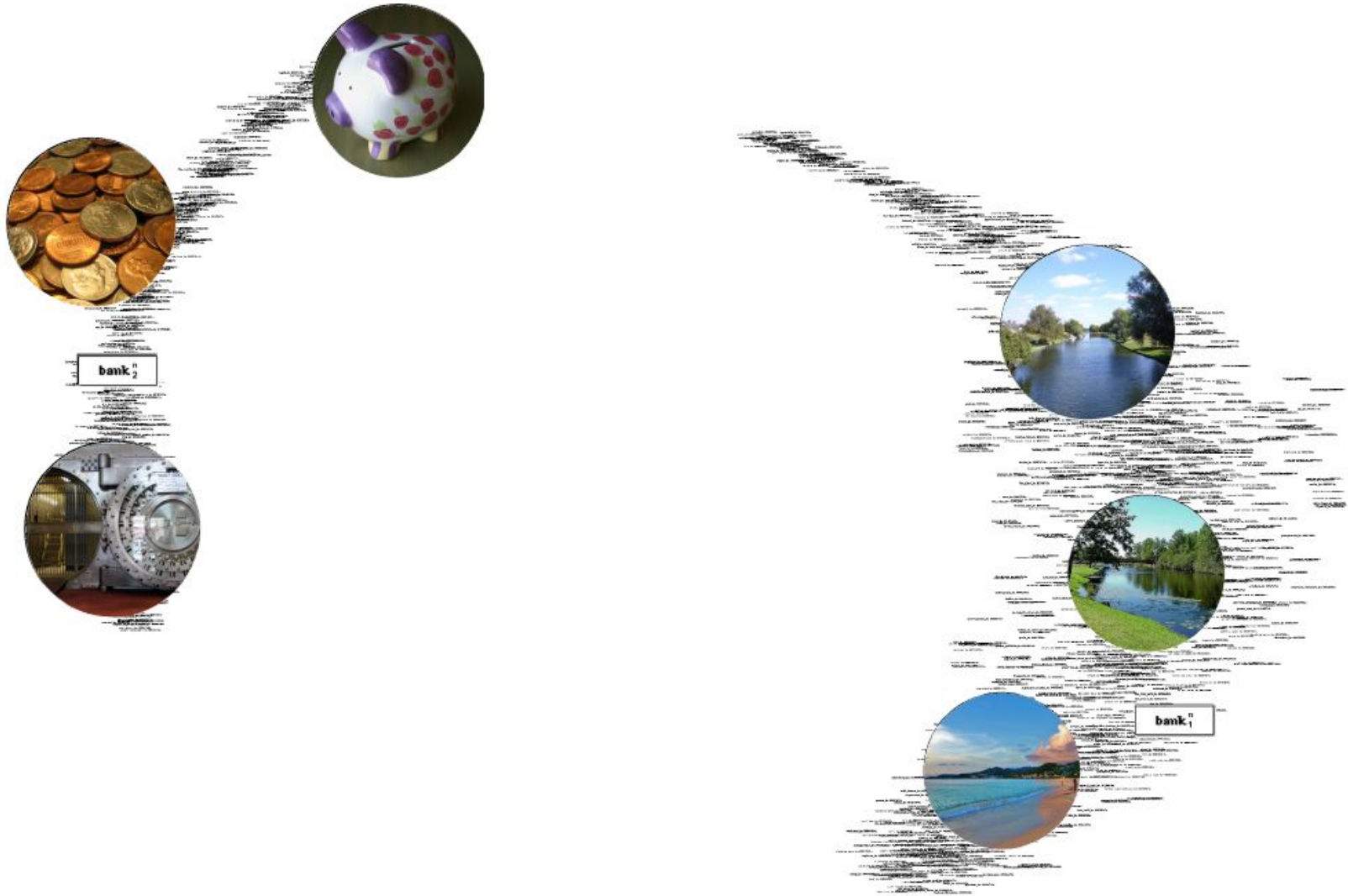
- NASARI 2.0 (August 2015): + Multilingual extension.

José Camacho Collados, Mohammad Taher Pilehvar and Roberto Navigli. *A Unified Multilingual Semantic Representation of Concepts*. **ACL 2015**, Beijing, China, pp. 741-751.

- NASARI 3.0 (March 2016): + Embedded representations, new applications.

José Camacho Collados, Mohammad Taher Pilehvar and Roberto Navigli. *Nasari: Integrating explicit knowledge and corpus statistics for a multilingual representation of concepts and entities*. **Artificial Intelligence Journal**, 2016, 240, 36-64.

# Key goal: obtain sense representations



# Key goal: obtain sense representations

- Nome
- Verbo

**Nome**

	<b>bank, streambank</b> Sloping land (especially the slope beside a body of water) ID: 00008363n   Concetto	<b>AR</b> ضفة, حافة <b>ZH</b> 岸, 河边 <b>FR</b> berge, rive <b>IT</b> riva, argine, sponda
	<b>bank, depository financial institution, banking company</b> A financial institution that accepts deposits and channels the money into lending activities ID: 00008364n   Concetto	<b>AR</b> مصرف (أموال), بنك, البنك <b>ZH</b> 銀行, 银行, 存放款金融机构 <b>FR</b> banque, institution financière de dépôt, établissement bancaire <b>IT</b> banca, banco, cassa
	<b>bank</b> A long ridge or pile ID: 00008365n   Concetto	<b>FR</b> banc <b>IT</b> banco
	<b>bank</b> An arrangement of similar objects in a row ID: 00008366n   Concetto	
	<b>bank</b> A supply or stock held in reserve for future use (especially in emergencies) ID: 00008367n   Concetto	<b>ZH</b> 储备金 <b>FR</b> banque <b>IT</b> banca

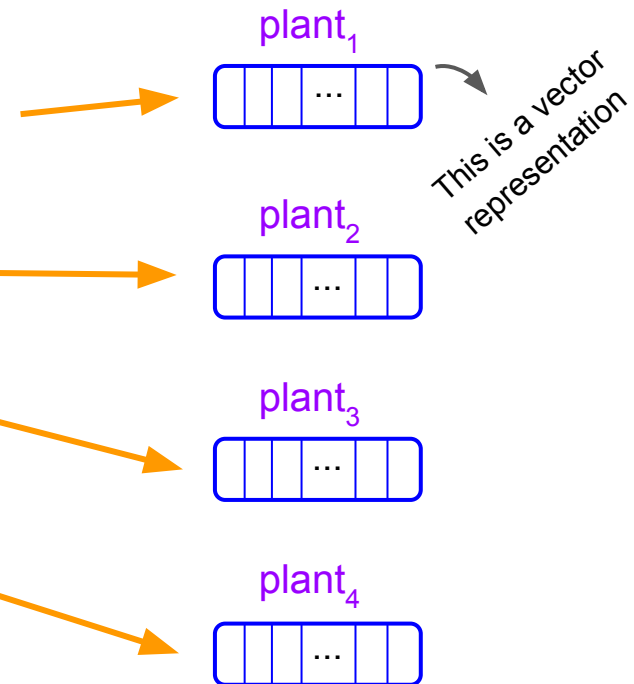
We want to create a separate representation for each entry of a given word

# Knowledge-based Sense Representations

Represent word senses as defined by sense inventories

## plant

- **plant, works, industrial plant** (buildings for carrying on industrial labor)
- **plant, flora, plant life** ((botany) a living organism lacking the power of locomotion)
- **plant** (an actor situated in the audience whose acting is rehearsed but seems spontaneous to the audience)
- **plant** (something planted secretly for discovery by another)



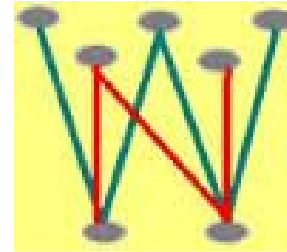
# Idea

Encyclopedic knowledge



**WIKIPEDIA**  
The Free Encyclopedia

Lexicographic knowledge



**WordNet**



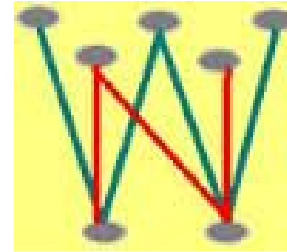
# Idea

Encyclopedic knowledge



**WIKIPEDIA**  
The Free Encyclopedia

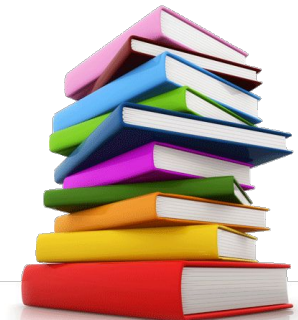
Lexicographic knowledge



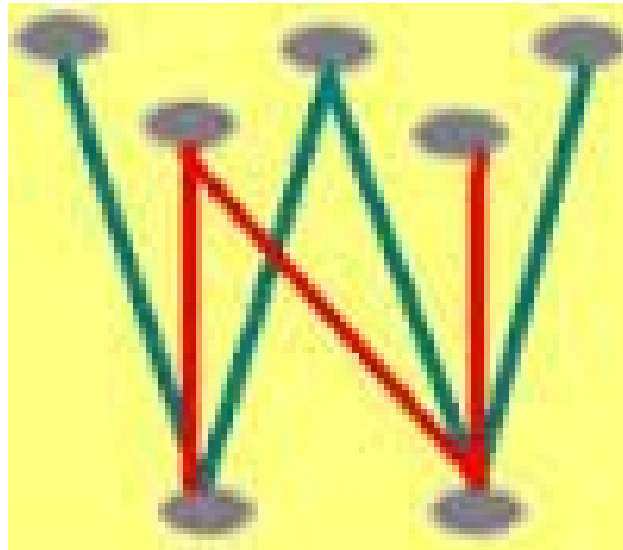
**WordNet**



**Information from text corpora**

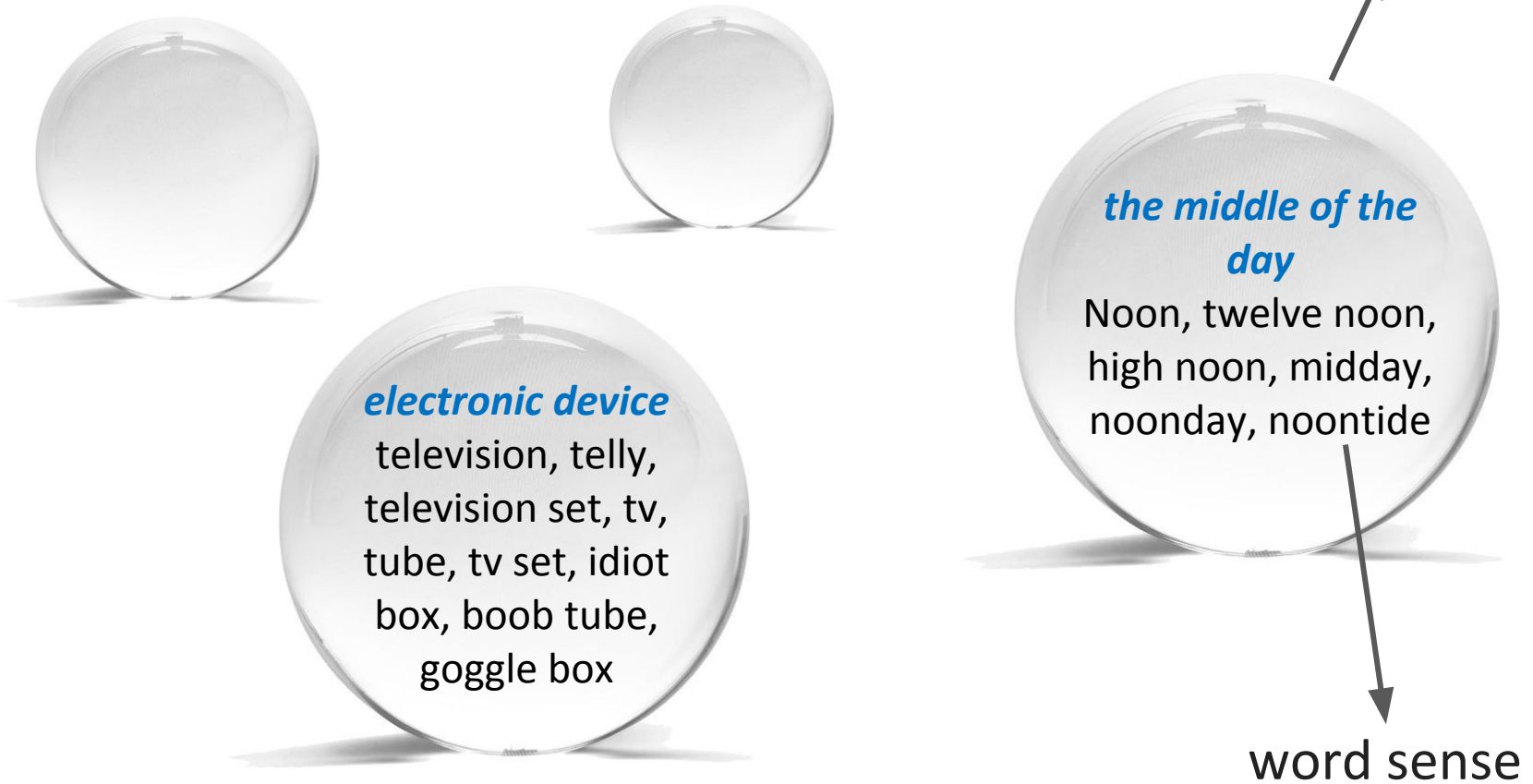


# WordNet



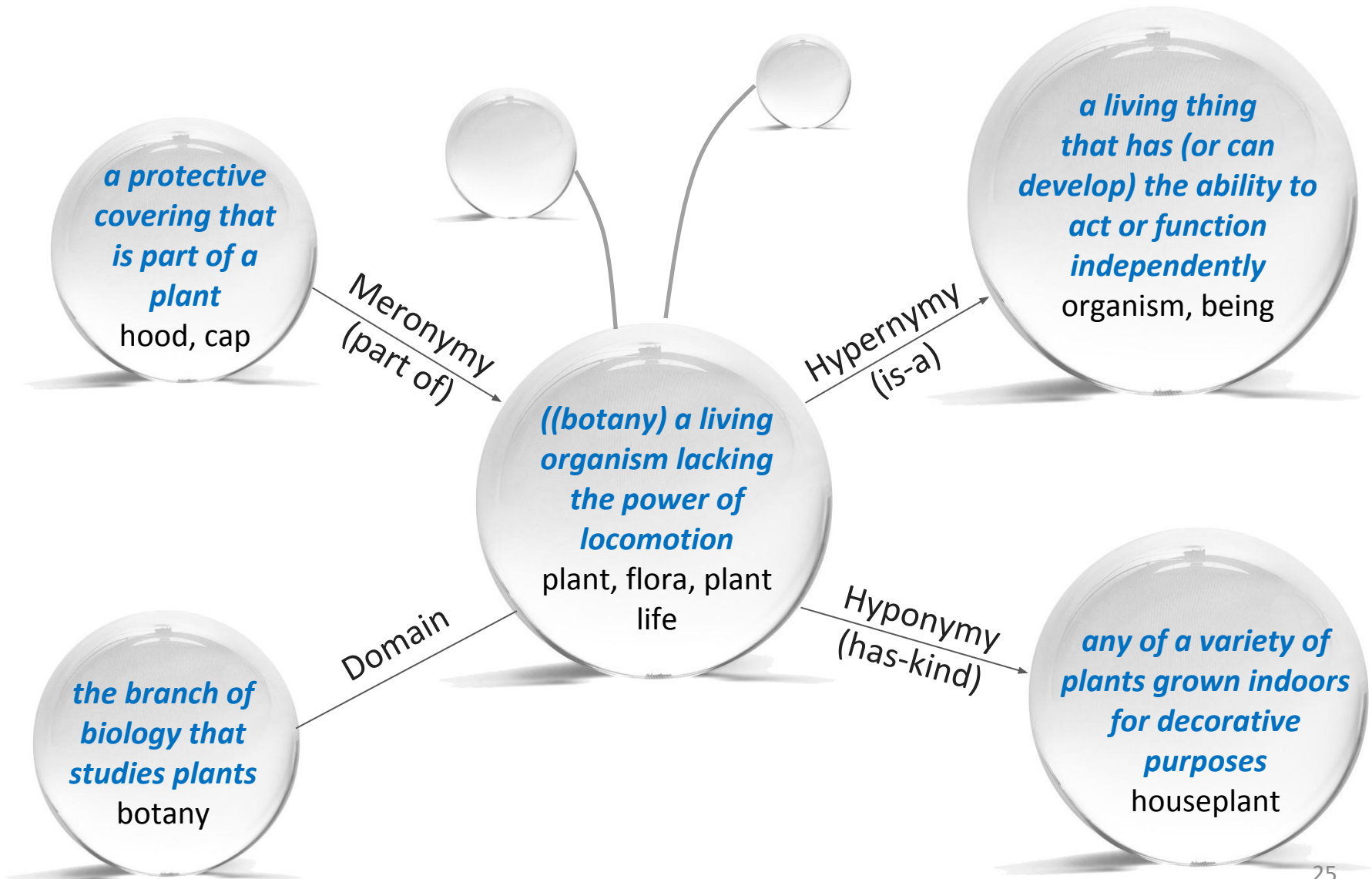
# WordNet

Main unit: synset (concept)





# WordNet semantic relations



# WordNet

## WordNet Search - 3.1

- [WordNet home page](#) - [Glossary](#) - [Help](#)

Word to search for:

Display Options:

Key: "S:" = Show Synset (semantic) relations, "W:" = Show Word (lexical) relations

Display options for sense: (gloss) "an example sentence"

### Noun

- [S:](#) (n) **plant**, [works](#), [industrial plant](#) (buildings for carrying on industrial labor) *"they built a large plant to manufacture automobiles"*
- [S:](#) (n) **plant**, [flora](#), [plant life](#) ((botany) a living organism lacking the power of locomotion)
- [S:](#) (n) **plant** (an actor situated in the audience whose acting is rehearsed but seems spontaneous to the audience)
- [S:](#) (n) **plant** (something planted secretly for discovery by another) *"the police used a plant to trick the thieves"; "he claimed that the evidence against him was a plant"*

### Verb

- [S:](#) (v) **plant**, [set](#) (put or set (seeds, seedlings, or plants) into the ground) *"Let's plant flowers in the garden"*
- [S:](#) (v) [implant](#), [engraft](#), [embed](#), [imbed](#), **plant** (fix or set securely or deeply) *"He planted a knee in the back of his opponent"; "The dentist implanted a tooth in the gum"*
- [S:](#) (v) [establish](#), [found](#), **plant**, [constitute](#), [institute](#) (set up or lay the groundwork for) *"establish a new department"*
- [S:](#) (v) **plant** (place into a river) *"plant fish"*
- [S:](#) (v) **plant** (place something or someone in a certain position in order to secretly observe or deceive) *"Plant a spy in Moscow"; "plant bugs in the dissident's apartment"*
- [S:](#) (v) **plant**, [implant](#) (put firmly in the mind) *"Plant a thought in the students' minds"*

[Link to online browser](#)

# Knowledge-based Sense Representations using WordNet

M. T. Pilehvar, D. Jurgens and R. Navigli: **Align, Disambiguate and Walk: A Unified Approach for Measuring Semantic Similarity** (ACL 2013)

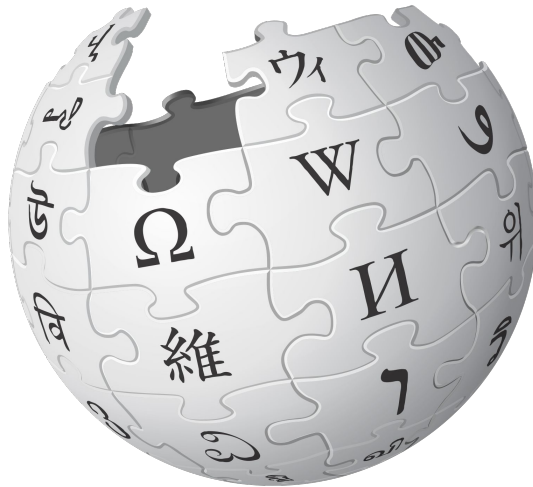
X. Chen, Z. Liu, M. Sun: **A Unified Model for Word Sense Representation and Disambiguation** (EMNLP 2014)

S. Rothe and H. Schutze: **AutoExtend: Extending Word Embeddings to Embeddings for Synsets and Lexemes** (ACL 2015)

S. K. Jauhar, C. Dyer, E. Hovy: **Ontologically Grounded Multi-sense Representation Learning for Semantic Vector Space Models** (NAACL 2015)

M. T. Pilehvar and N. Collier: **De-Conflated Semantic Representations** (EMNLP 2016)

# Wikipedia



**WIKIPEDIA**  
The Free Encyclopedia

# Wikipedia

High coverage of **named entities** and **specialized concepts** from different domains



Article Talk

Read Edit View history Search

## University of California, Los Angeles

From Wikipedia, the free encyclopedia

Coordinates: 34°04′20.00″N 118°26′38.75″W﻿ / ﻿﻿ / ﻿

*"UCLA", "Ucla", and "U.C.L.A." redirect here. For other uses, see UCLA (disambiguation).*

The **University of California, Los Angeles (UCLA)** is a public research university located in the Westwood neighborhood of Los Angeles, California, United States. It became the University of California Southern Branch in 1919, making it the second-oldest undergraduate campus of the ten-campus system after the original University of California campus in Berkeley (1873).<sup>[11]</sup> It offers 337 undergraduate and graduate degree programs in a wide range of disciplines.<sup>[12]</sup> With an approximate enrollment of 30,000 undergraduate and 12,000 graduate students, UCLA has the highest enrollment of any university in California<sup>[6]</sup> and is the most applied to university in the United States with over 112,000 applications for fall 2015.<sup>[13]</sup>

The university is organized into five undergraduate colleges, seven professional schools, and four professional health science schools. The undergraduate colleges are the College of Letters and Science; Henry Samueli School of Engineering and Applied Science (HSSEAS); School of the Arts and Architecture; School of Theater, Film, and Television; and School of Nursing. Fifteen<sup>[14]</sup><sup>[15]</sup> Nobel laureates, one Fields Medalist,<sup>[16]</sup> and three Turing Award winners<sup>[17]</sup> have been faculty, researchers, or alumni. Among the current faculty members, 55 have been elected to the National Academy of Sciences, 28 to the National Academy of Engineering, 39 to the Institute of Medicine, and 124 to the American Academy of Arts and Sciences.<sup>[18]</sup> The university was elected to the Association of American Universities in 1974.<sup>[19]</sup>

UCLA student-athletes compete as the Bruins in the Pacific-12 Conference. The Bruins have won 125 national championships, including 112 NCAA team championships.<sup>[20]</sup><sup>[21]</sup> UCLA student-athletes have won 250 Olympic medals: 125 gold, 65 silver and 60 bronze.<sup>[22]</sup> The Bruins have competed in every Olympics since 1920 with one exception (1924), and have won a gold medal in every Olympics that the United States has participated in since 1932.<sup>[23]</sup>

**Contents** [hide]

1 History

### University of California, Los Angeles

**UCLA official seal**

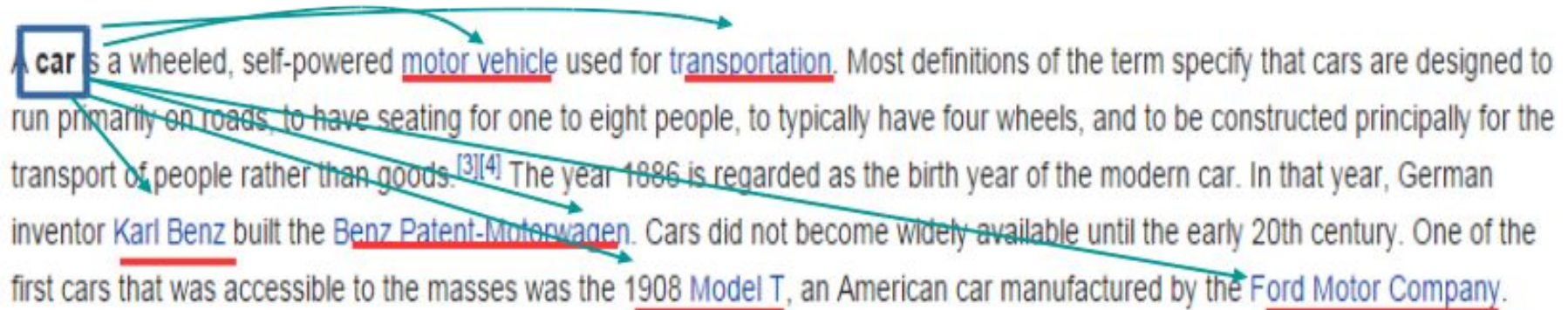
<b>Former names</b>	State Normal School at Los Angeles (1882-1919) University of California Southern Branch (1919–1927) University of California at Los Angeles (1927–1958)
<b>Motto</b>	<i>Fiat lux</i> (Latin)
<b>Motto in English</b>	Let there be light

# Wikipedia hyperlinks

A **car** is a wheeled, self-powered [motor vehicle](#) used for [transportation](#). Most definitions of the term specify that cars are designed to run primarily on roads, to have seating for one to eight people, to typically have four wheels, and to be constructed principally for the transport of people rather than goods.<sup>[3][4]</sup> The year 1886 is regarded as the birth year of the modern car. In that year, German inventor [Karl Benz](#) built the [Benz Patent-Motorwagen](#). Cars did not become widely available until the early 20th century. One of the first cars that was accessible to the masses was the 1908 [Model T](#), an American car manufactured by the [Ford Motor Company](#).

# Wikipedia hyperlinks

A **car** is a wheeled, self-powered motor vehicle used for transportation. Most definitions of the term specify that cars are designed to run primarily on roads, to have seating for one to eight people, to typically have four wheels, and to be constructed principally for the transport of people rather than goods.<sup>[3][4]</sup> The year 1886 is regarded as the birth year of the modern car. In that year, German inventor Karl Benz built the Benz Patent-Motorwagen. Cars did not become widely available until the early 20th century. One of the first cars that was accessible to the masses was the 1908 Model T, an American car manufactured by the Ford Motor Company.

A diagram illustrating hyperlinks from the word "car" in the text. The word "car" is enclosed in a blue box. Five teal arrows originate from the right side of this box and point to the following underlined terms in the text: "motor vehicle", "transportation", "Karl Benz", "Benz Patent-Motorwagen", and "Ford Motor Company".



BabelNet

Thanks to an automatic mapping algorithm, **BabelNet integrates Wikipedia and WordNet**, among other resources (Wiktionary, OmegaWiki, WikiData...).

Key feature: **Multilinguality** (271 languages)



# BabelNet



BabelNet

ENTRA REGISTRATI

jaguar|

ENGLISH

4 SELEZIONATE

TRADUCI

PREFERENZE

Tutti

Concetti

Entità nominate



21 risultati

Nome

Nome

Concept

Entity



jaguar, panther, Felis onca

A large spotted feline of tropical America similar to the leopard; in some classifications considered a member of the genus Felis

ID: 00033987n | Concetto

ZH 美洲豹

FR jaguar, panthère

IT giaguaro, Panthera onca, pantera

ES jaguar, panthera onca, pantera



Jaguar Cars, Jaguar

Jaguar Cars is a brand of Jaguar Land Rover, a British multinational car manufacturer headquartered in Whitley, Coventry, England, owned by Tata Motors since 2008.

ID: 00688731n | Entità

ZH 捷豹

FR Jaguar (automobile)

IT Jaguar

ES Jaguar Cars, Jaguar



Atari Jaguar, Jaguar (video game console)

The Atari Jaguar is a home video game console that was released by Atari Corporation in 1993.

ID: 02142312n | Entità

ZH Atari Jaguar, 雅达利Jaguar

FR Jaguar (console)

IT Atari Jaguar

ES Atari Jaguar



Mac OS X v10.2, Jaguar (macos)

Mac OS X version 10.2 Jaguar is the third major release of Mac OS X, Apple's desktop and server operating system.

ZH Mac OS X Jaguar, Mac OS X v10.2

FR Mac OS X v10.2

# BabelNet

It follows the same structure of WordNet:  
**synsets** are the main units

Nome



jaguar, panther, Felis onca

A large spotted feline of tropical America similar to the leopard; in some classifications considered a member of the genus Felis

ID: 00033987n | Concetto

ZH 美洲豹

FR jaguar, panthère

IT giaguaro, Panthera onca, pantera

ES jaguar, panthera onca, pantera

# BabelNet

In this case, **synsets are multilingual**

Nome



jaguar, panther, Felis onca

A large spotted feline of tropical America similar to the leopard; in some classifications considered a member of the genus Felis

ID: [00033987n](#) | Concetto

- ZH 美洲豹
- FR jaguar, panthère
- IT giaguaro, Panthera onca, pantera
- ES jaguar, panthera onca, pantera

# NASARI: Integrating Explicit Knowledge and Corpus Statistics for a Multilingual Representation of Concepts and Entities

(Camacho-Collados et al., AIJ 2016)

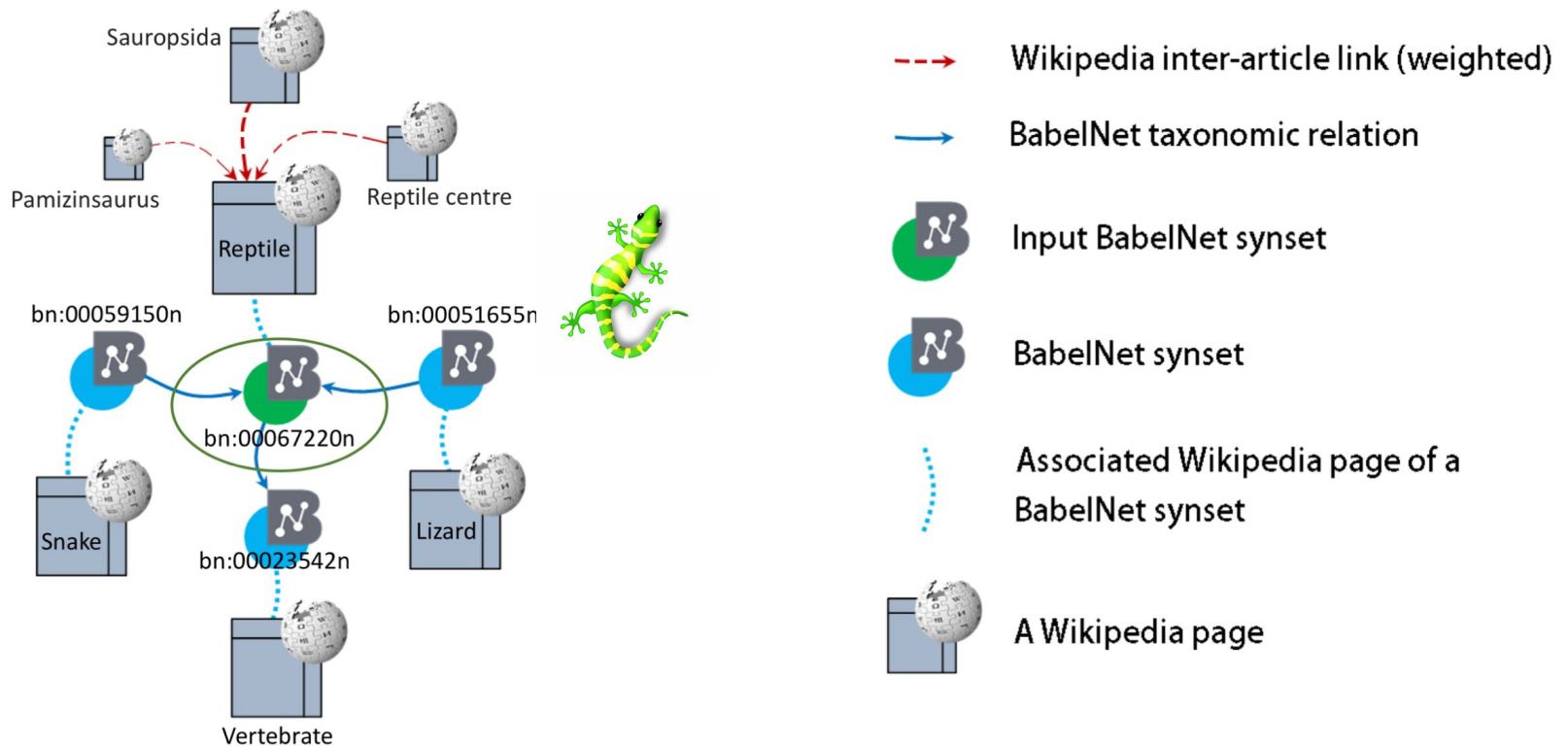
## Goal

Build vector representations for multilingual BabelNet synsets.

## How?

We exploit **Wikipedia semantic network** and **WordNet taxonomy** to construct a subcorpus (contextual information) for any given BabelNet synset.

# Pipeline



Process of obtaining contextual information for a BabelNet synset exploiting BabelNet taxonomy and Wikipedia as a semantic network

# Three types of vector representations

Three types of vector representations:

- **Lexical** (dimensions are words)
  
  
  
  
  
  
  
  
  
  
- **Unified** (dimensions are multilingual BabelNet synsets)
  
  
  
  
  
  
  
  
  
  
- **Embedded** (latent dimensions)

# Three types of vector representations

Three types of vector representations:

- **Lexical** (dimensions are words): Dimensions are weighted via **lexical specificity**, a statistical measure based on the hypergeometric distribution.
- **Unified** (dimensions are multilingual BabelNet synsets)
- **Embedded** (latent dimensions)

# Lexical specificity

It is a statistical measure based on the **hypergeometric distribution**, particularly suitable for term extraction tasks.

Thanks to its statistical nature, it is **less sensitive to corpus sizes** than the conventional *tf-idf* (in our setting, it **consistently outperforms *tf-idf* weighting**).



# Three types of vector representations

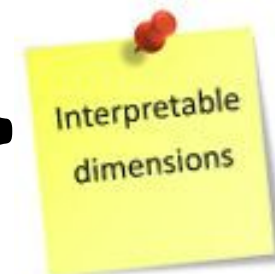
Three types of vector representations:

- **Lexical** (dimensions are words):
  
- **Unified** (dimensions are multilingual BabelNet synsets): This representation uses a **hypernym-based clustering technique** and can be used in **cross-lingual applications**
  
- **Embedded** (latent dimensions)

# Three types of vector representations

Three types of vector representations:

- **Lexical** (dimensions are words):
- **Unified** (dimensions are multilingual BabelNet synsets): This representation uses a **hypernym-based clustering technique** and can be used in **cross-lingual applications**
- **Embedded** (latent dimensions)

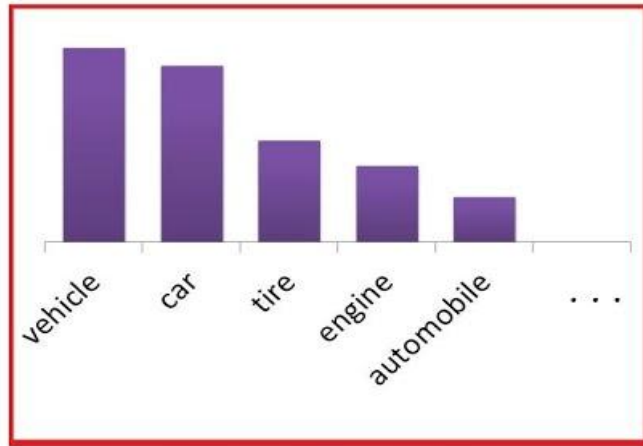


# Lexical and unified vector representations

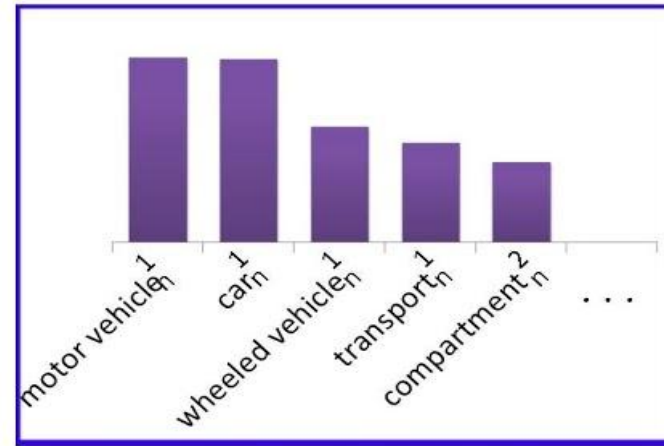
Interpretable dimensions



## EXAMPLE



Word-based representation



Synset-based representation

# From a lexical vector to a unified vector

**Lexical vector**= (automobile, car, engine, vehicle, motorcycle, ...)

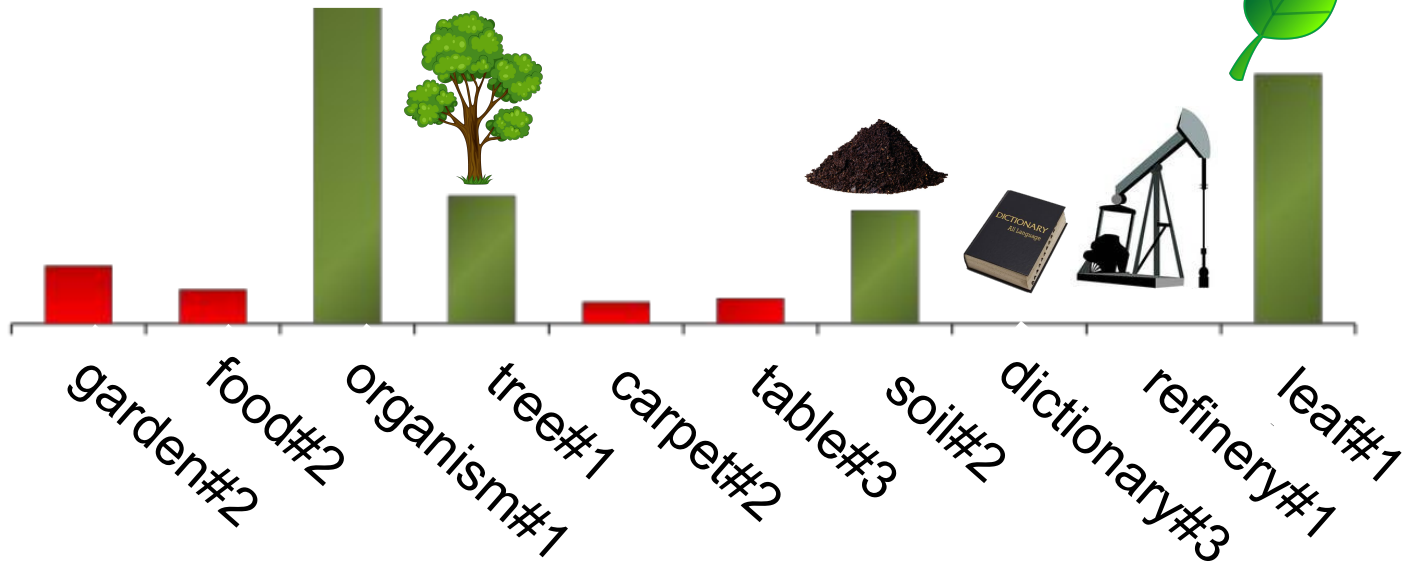


**motor\_vehicle<sub>n</sub><sup>1</sup>**

**Unified vector**= (motor\_vehicle<sub>n</sub><sup>1</sup>, ... )

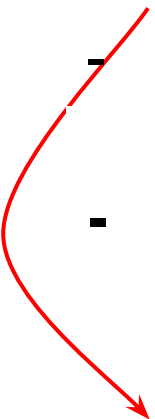
# Human-interpretable dimensions

plant (living organism)



# Three types of vector representations

Three types of vector representations:

- **Lexical** (dimensions are words)
  - **Unified** (dimensions are multilingual BabelNet synsets)
  - **Embedded**: Low-dimensional vectors (latent) exploiting **word embeddings** obtained from **text corpora**. This representation is obtained by plugging word embeddings on the lexical vector representations.
- 

# Three types of vector representations

Three types of vector representations:

- **Lexical** (dimensions are words)
- **Unified** (dimensions are multilingual BabelNet synsets)
- **Embedded**: Low-dimensional vectors (latent) exploiting **word embeddings** obtained from **text corpora**. This representation is obtained by plugging word embeddings on the lexical vector representations.

Word and synset embeddings share the same vector space!

# Sense-based Semantic Similarity

Based on the semantic similarity between senses.

Two main measures:

- **Cosine similarity** for low-dimensional vectors
- **Weighted Overlap** for sparse high-dimensional vectors (interpretable)

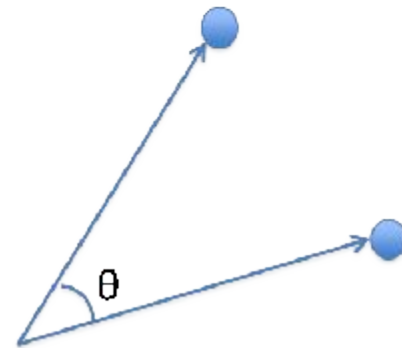


# Vector Comparison

## Cosine Similarity

The most commonly used measure for the similarity of vector space model (sense) representations

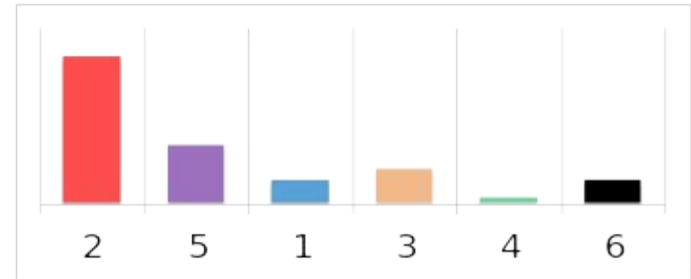
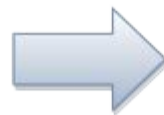
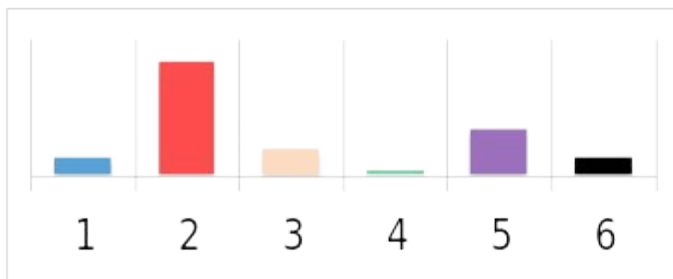
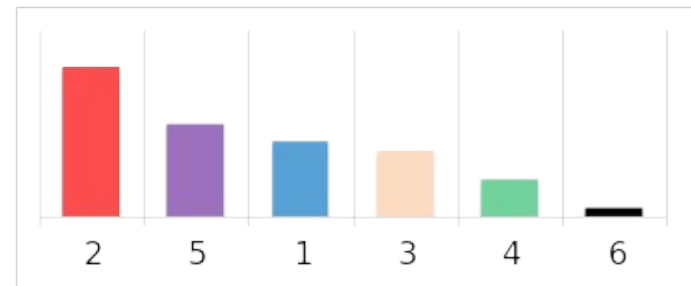
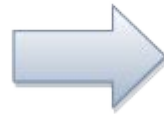
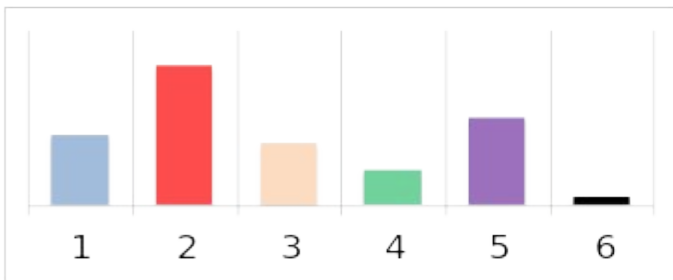
$$\text{sim}(A, B) = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|}$$



# Vector Comparison

## Weighted Overlap

$$WO(v_1, v_2) = \frac{\sum_{q \in O} \left( \text{rank}(q, v_1) + \text{rank}(q, v_2) \right)^{-1}}{\sum_{i=1}^{|O|} (2i)^{-1}}$$



# Embedded vector representation

## *Closest senses*



### Bank (financial institution)

Closest senses	Cosine
Deposit account	0.99
Universal bank	0.99
British banking	0.98
German banking	0.98
Commercial bank	0.98
Banking in Israel	0.98
Financial institution	0.98
Community bank	0.97

### Bank (geography)

Closest senses	Cosine
Stream bed	0.98
Current (stream)	0.97
River engineering	0.97
Braided river	0.97
Fluvial terrace	0.97
Bar (river morphology)	0.97
River	0.97
Perennial stream	0.96

### *bank*

Closest senses	Cosine
Bank (financial institution)	0.86
Universal bank	0.86
British banking	0.86
German banking	0.85
Branch (banking)	0.85
McFadden Act	0.85
Four Northern Banks	0.84
State bank	0.84

# NASARI semantic representations

## *Summary*

- Three types of semantic representation: **lexical, unified and embedded.**
- **High coverage of concepts and named entities** in multiple languages (all Wikipedia pages covered).

# NASARI semantic representations

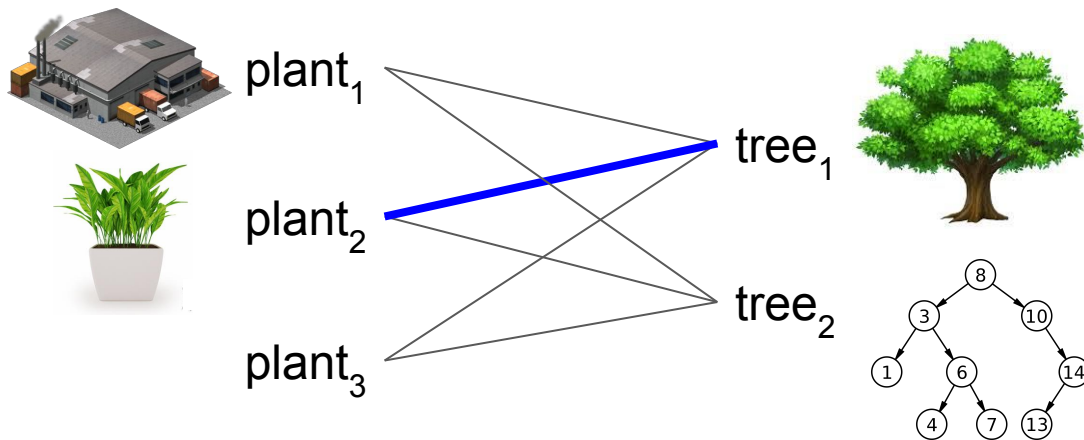
## *Summary*

- Three types of semantic representation: **lexical, unified and embedded.**
- **High coverage of concepts and named entities** in multiple languages (all Wikipedia pages covered).

**What's next?** **Evaluation** and **use** of these semantic representations in **NLP applications.**

# How are sense representations used for word similarity?

1- **MaxSim**: similarity between the most similar senses across two words



$$\text{MaxSim}(w, w') \stackrel{\text{def}}{=} \max_{1 \leq j \leq K, 1 \leq k \leq K'} d(\pi_k(w), \pi_j(w'))$$

# Intrinsic evaluation

## *Monolingual semantic similarity (English)*

	MC-30		WS-Sim		SimLex-999 (nouns)		Average	
	$r$	$\rho$	$r$	$\rho$	$r$	$\rho$	$r$	$\rho$
NASARI	0.89	0.78	0.74	0.72	0.50	0.49	<b>0.71</b>	0.67
NASARI <sub>lexical</sub>	0.88	0.81	0.74	0.73	0.51	0.49	<b>0.71</b>	<b>0.68</b>
NASARI <sub>unified</sub>	0.88	0.78	0.72	0.70	0.49	0.48	0.70	0.65
NASARI <sub>embed</sub>	<b>0.91</b>	0.83	0.68	0.68	0.48	0.46	0.69	0.66
ESA	0.59	0.65	0.45	0.53	0.16	0.23	0.40	0.47
Lin	0.76	0.72	0.66	0.62	<b>0.58</b>	<b>0.58</b>	0.67	0.64
ADW	0.79	0.83	0.63	0.67	0.44	0.45	0.62	0.65
Chen	0.82	0.82	0.63	0.64	0.48	0.44	0.64	0.63
Word2Vec	0.80	0.80	<b>0.76</b>	<b>0.77</b>	0.46	0.45	0.67	0.67
Best-Word2Vec	0.83 <sup>‡</sup>	0.83 <sup>‡</sup>	0.76 <sup>‡</sup>	0.78 <sup>‡</sup>	0.48	0.49	0.69	0.70
Best-PMI-SVD	0.76 <sup>‡</sup>	0.71 <sup>‡</sup>	0.68 <sup>‡</sup>	0.66 <sup>‡</sup>	0.40	0.40	0.61	0.59
SensEmbed	0.89	<b>0.88</b>	0.65	0.75	0.46 <sup>†</sup>	0.47 <sup>†</sup>	0.67	0.70

# Intrinsic evaluation

(Camacho-Collados et al., ACL 2015)

Most current approaches are developed for English only and there are no many datasets to evaluate multilinguality.

To this end, we developed a semi-automatic framework to extend English datasets to other languages (and across languages):

*Data available at*

<http://lcl.uniroma1.it/similarity-datasets/>



# Intrinsic evaluation

## *Multilingual semantic similarity*

<b>English</b>	<i>r</i>	$\rho$	<b>French</b>	<i>r</i>	$\rho$	<b>German</b>	<i>r</i>	$\rho$	<b>Spanish</b>	<i>r</i>	$\rho$
NASARI	0.81	0.78	NASARI	<b>0.82</b>	0.73	NASARI	0.69	0.65	NASARI	<b>0.85</b>	0.79
NASARI <sub>lexical</sub>	0.80	0.78	NASARI <sub>lexical</sub>	0.80	0.70	NASARI <sub>lexical</sub>	0.69	0.67	NASARI <sub>lexical</sub>	<b>0.85</b>	0.79
NASARI <sub>unified</sub>	0.80	0.76	NASARI <sub>unified</sub>	<b>0.82</b>	<b>0.76</b>	NASARI <sub>unified</sub>	<b>0.71</b>	<b>0.68</b>	NASARI <sub>unified</sub>	0.82	0.77
NASARI <sub>embed</sub>	<b>0.82</b>	<b>0.80</b>	–	–	–	–	–	–	NASARI <sub>embed</sub>	0.79	0.77
SOC-PMI	0.61	–	SOC-PMI	0.19	–	SOC-PMI	0.27	–	–	–	–
PMI	0.41	–	PMI	0.34	–	PMI	0.40	–	–	–	–
LSA-Wiki	0.65	0.69	LSA-Wiki	0.57	0.52	–	–	–	–	–	–
Wiki-wup	0.59	–	–	–	–	Wiki-wup	0.65	–	–	–	–
Word2Vec	–	0.73	Word2Vec	–	0.47	Word2Vec	–	0.53	Best-Word2Vec	0.80	<b>0.80</b>
Retrofitting	–	0.77	Retrofitting	–	0.61	Retrofitting	–	0.60	–	–	–
NASARI <sub>poly-embed</sub>	0.74	0.77	NASARI <sub>poly-embed</sub>	0.60	0.69	NASARI <sub>poly-embed</sub>	0.46	0.52	NASARI <sub>poly-embed</sub>	0.68	0.74
Polyglot-embed	0.51	0.55	Polyglot-embed	0.38	0.35	Polyglot-embed	0.18	0.15	Polyglot-embed	0.51	0.56
IAA	0.85°	-	IAA	-	-	IAA	0.81	-	IAA	0.83	-

# Intrinsic evaluation

## *Cross-lingual semantic similarity*

Measure	EN-FR		EN-DE		EN-ES		FR-DE		FR-ES		DE-ES		Average	
	$r$	$\rho$	$r$	$\rho$	$r$	$\rho$	$r$	$\rho$	$r$	$\rho$	$r$	$\rho$	$r$	$\rho$
NASARI <sub>unified</sub>	<b>0.84</b>	0.79	<b>0.79</b>	0.79	<b>0.84</b>	0.82	0.75	0.70	<b>0.86</b>	0.78	<b>0.81</b>	<b>0.80</b>	<b>0.82</b>	0.78
CL-MSR-2.0	0.30	–	–	–	–	–	–	–	–	–	–	–	–	–
NASARI <sub>pivot</sub>	0.79	0.69	0.78	0.76	0.80	0.74	<b>0.79</b>	0.70	0.80	0.67	0.72	0.68	0.78	0.71
ADW <sub>pivot</sub>	0.80	0.82	0.73	<b>0.82</b>	0.78	<b>0.84</b>	0.72	<b>0.77</b>	0.81	<b>0.81</b>	0.68	0.72	0.75	<b>0.80</b>
Word2Vec <sub>pivot</sub>	0.77	0.82	0.70	0.73	0.76	0.80	0.65	0.70	0.75	0.76	0.64	0.63	0.71	0.74
Best-Word2Vec <sub>pivot</sub>	0.75	<b>0.84</b>	0.69	0.76	0.75	0.82	0.77	0.73	0.74	0.79	0.64	0.64	0.72	0.76
Best-PMI-SVD <sub>pivot</sub>	0.76	0.76	0.72	0.74	0.77	0.77	0.65	0.69	0.76	0.74	0.62	0.61	0.71	0.72

# **NEW: SemEval 2017** task on multilingual and cross-lingual semantic word similarity

Large datasets to evaluate semantic similarity in **five languages** (within and across languages): English, Farsi, German, Italian and Spanish.

Additional challenges:

- **Multiwords:** *black hole*
- **Entities:** *Microsoft*
- **Domain-specific terms:** *chemotherapy*

*Data available at*

<http://alt.qcri.org/semeval2017/task2/>

# Applications

- Domain labeling/adaptation
- Word Sense Disambiguation
- Sense Clustering
- Topic categorization and sentiment analysis

# Domain labeling

(Camacho-Collados et al., AIJ 2016)

Annotate each **concept/entity** with its corresponding **domain of knowledge**.

To this end, we use the [Wikipedia featured articles page](#), which includes 34 domains and a number of Wikipedia pages associated with each domain (*Biology, Geography, Mathematics, Music, etc.* ).

# Domain labeling

## Wikipedia featured articles

### Chemistry and mineralogy

Acetic acid • Antioxidant • Astatine • Caesium • Californium • Cyclol • Diamond • DNA nanotechnology • Enzyme • Enzyme inhibitor • Enzyme kinetics • Fluorine • Francium • Germanium • Helium • Hydrochloric acid • Hydrogen • Iridium • Lead(II) nitrate • Metalloid • Nicotinamide adenine dinucleotide • Niobium • Noble gas • Oxidative phosphorylation • Oxygen • Periodic table • Plutonium • Psilocybin • Rhodocene • Synthetic diamond • Technetium • Titanium • Ununocium • Ununseptium • Uranium • Xenon • Yogo sapphire • Yttrium • Zinc

### Chemistry and mineralogy biographies

James Bryant Conant • Joseph Priestley

### Computing

4chan • Acid2 • Delrina • Folding@home • Macintosh Classic • Manchester Mark 1 • Manchester Small-Scale Experimental Machine • Microsoft Security Essentials • The Million Dollar Homepage • NeXT • Parallel computing • PowerBook 100 • Rosetta@home • ROT13 • Scene7

### Culture and society

Aggie Bonfire • Hadji Ali • The Livestock Conservancy • Anna Anderson • Marshall Applewhite • Baden-Powell House • Isabella Beeton • Biddenden Maids • William D. Boyce • Guy Bradley • Burke and Hare murders • William Henry Burry • "The Bus Uncle" • Josephine Butler • The Chaser APEC pranks • Cleveland Street scandal • Cock Lane ghost • D. B. Cooper • Daylight saving time • Disco Demolition Night • Charles Domery • Dorset Ooser • Marjory Stoneman Douglas • Montague Druitt • W. E. B. Du Bois • Monroe Edwards • Female genital mutilation • Terry Fox • Ursula Franklin • Free Association of German Trade Unions • Margaret Fuller • E. Urner Goodman • Debora Green • Stanley Green • Green children of Woolpit • Gropecunt Lane • Guy Fawkes Night • Hanged, drawn and quartered • William Hillcourt • Fanny Imlay • Indigenous people of the Everglades region • *An Introduction to Animals and Political Theory* • Jack the Ripper • *Jack the Ripper: The Final Solution* • *Ketuanan Melayu* • Akhtar Hameed Khan • Kylfings • Daniel Lambert • Liberty Bell • Lynching of Jesse Washington • Macedonia (terminology) • Mantra-Rock Dance • Bob Marshall (wilderness activist) • Murder of Dwayne Jones • Florence Nagle • The Negro Motorist Green Book • Emmeline Pankhurst • Pig-faced women • Polish culture during World War II • Postage stamps of Ireland • Ramblin' Wreck • Rosewood massacre • Royal baccarat scandal • Same-sex marriage in Spain • Mark Satin • Scouting • John Martin Scripps • Sexuality after spinal cord injury • Grace Sherwood • Society of the Song dynasty • Stonewall riots • Taiwanese aborigines • Mary Toft • Toraja • Truthiness • Voluntary Human Extinction Movement • Whitechapel murders • Wife selling (English custom) • Wonderbra • Wood Badge • Robert Sterling Yard • Zong massacre

### Education

Alpha Kappa Alpha • Amador Valley High School • ANAK Society • Avery Coonley School • Baltimore City College • Boden Professor of Sanskrit election, 1860 • James E. Boyd (scientist) • C. R. M. F. Cruttwell • Dartmouth College • Duke University • Florida Atlantic University • Georgetown University • The Green (Dartmouth College) • *The Guardian of Education* • History of Baltimore City College • History of Texas A&M University • The Judd School • Kappa Kappa Psi • *Lessons for Children* • Michigan State University • Ohio Wesleyan University • Oriel College, Oxford • Romney Literary Society • Royal National College for the Blind • School for Creative and Performing Arts • Shimer College • *Some Thoughts Concerning Education* • Stuyvesant High School • Texas A&M University • Texas Tech University • *Thoughts on the Education of Daughters* • Tuck School of Business • United States Academic Decathlon • United States Military Academy • University of California, Riverside • University of Michigan • Vkhutemas

### Engineering and technology

2013 Rosario gas explosion • Apollo 8 • Atomic line filter • Caesar cipher • Calutron • CFM International CFM56 • Construction of the World Trade Center • Distributed element filter • Draining and development of the Everglades • Gas metal arc welding • Gas tungsten arc welding • Grand Coulee Dam • Halkett boat • Hanford Site • History of timekeeping devices • Hoover Dam • Mechanical filter • Oil shale • Panavision • Pigeon photography • Rampart Dam • Renewable energy in Scotland • Restoration of the Everglades • Rolls-Royce Merlin • Rolls-Royce R • Scout Moor Wind Farm • Shale oil extraction • Shielded metal arc welding • Shoe polish • Sholes and Glidden typewriter • Shuttle-*Mir* Program • Science and technology of the Song dynasty • Waveguide filter • Webley Revolver • Welding • World Science Festival, 2008

# Domain labeling

## How to associate a synset with a domain?

- We first construct a **NASARI lexical vector** for the **concatenation of all Wikipedia pages associated with a given domain** in the featured article page.
- Then, we calculate the **semantic similarity** between the corresponding NASARI vectors of the synset and all domains:

$$\hat{d}(s) = \arg \max_{d \in D} \text{WO}(\vec{\text{NASARI}}_{lex}(s), \vec{v}_{lex}(d))$$

# Domain labeling

This results in **over 1.5M synsets** associated with a domain of knowledge.

This domain information has already been integrated in the last version of BabelNet.



# Domain labeling



LOG IN REGISTER

eclipse ENGLISH TRANSLATE INTO... SEARCH

PREFERENCES

All Concepts Named Entities 48 results

🎵 🖱️ 🌐 ⭐ 🎮 ⚽ 📊 📄 🖼️ 🌍 🎬 📈

- Noun
- Verb

## Noun

Physics and astronomy

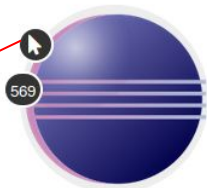


### eclipse, occultation

One celestial body obscures another

ID: 00029648n | Concept

Computing



### Eclipse (software)

In computer programming, Eclipse is an integrated development environment.

ID: 01457115n | Named Entity

Media



### The Twilight Saga: Eclipse, Eclipse (2010 film)

The Twilight Saga: Eclipse, commonly referred to as Eclipse, is a 2010 American romantic fantasy film based on Stephenie Meyer's 2007 novel Eclipse.

ID: 01455414n | Named Entity

# Domain labeling

	WordNet dataset			BabelNet dataset		
	Precision	Recall	F-Measure	Precision	Recall	F-Measure
NASAR <sub>lexical</sub>	<b>77.9</b>	70.1	<b>73.8</b>	<b>62.3</b>	<b>40.5</b>	<b>49.1</b>
Wikipedia-TF	25.4	16.4	19.9	3.4	2.5	2.9
Wikipedia-TFidf	45.9	29.7	36.1	8.8	6.5	7.5
Taxo-Prop (WN)	71.3	70.7	71.0	–	–	–
Taxo-Prop (BN)	73.5	<b>73.5</b>	73.5	48.3	37.2	42.0
WN-Domains-3.2	93.6	64.4	76.3	–	–	–

## Domain labeling results on WordNet and BabelNet

# BabelDomains

(Camacho-Collados and Navigli, EACL 2017)

**As a result:**

Unified resource with information about domains of knowledge

*BabelDomains available for **BabelNet**, **Wikipedia** and **WordNet** available at*

*<http://icl.uniroma1.it/babeldomains>*

*Already integrated into BabelNet (online interface and API)*

# Domain filtering for supervised distributional hypernym discovery

(Espinosa-Anke et al., EMNLP 2016;  
Camacho-Collados and Navigli, EACL 2017)



*Apple*

is | a

*Fruit*

**Task:** Given a term, predict its hypernym(s)

**Model:** Distributional supervised system based on the transformation matrix of Mikolov et al. (2013).

**Idea:** Training data filtered by domain of knowledge

# Domain filtering for supervised distributional hypernym discovery

**Domain-filtered training data**

	art			biology			education			geography			health		
Train	MRR	MAP	R-P	MRR	MAP	R-P	MRR	MAP	R-P	MRR	MAP	R-P	MRR	MAP	R-P
5k	0.12	0.12	0.12	0.63	0.63	0.59	0.00	0.00	0.00	0.08	0.07	0.07	0.08	0.08	0.07
15k	0.21	0.20	0.18	<b>0.84</b>	0.72	0.79	0.22	0.22	0.21	0.15	0.14	0.14	0.08	0.07	0.07
25k	<b>0.29</b>	0.27	<b>0.26</b>	<b>0.84</b>	<b>0.83</b>	<b>0.81</b>	0.33	0.32	0.30	<b>0.23</b>	<b>0.22</b>	<b>0.21</b>	0.09	0.09	0.08
25k+ $K_{1k}^d$	<b>0.29</b>	<b>0.28</b>	<b>0.26</b>	<b>0.84</b>	0.80	0.79	0.32	0.29	0.27	0.22	<b>0.22</b>	<b>0.21</b>	0.09	0.09	0.08
25k+ $K_{25k}^d$	0.26	0.24	0.22	0.70	0.63	0.56	<b>0.38</b>	<b>0.36</b>	<b>0.33</b>	0.15	0.13	0.12	0.11	0.11	0.10
25k+ $K_{50k}^r$	0.28	0.26	0.24	0.82	0.77	0.72	0.36	0.33	0.30	0.17	0.16	0.16	<b>0.12</b>	0.11	0.10
100k $_{wd}^r$	0.00	0.00	0.00	<b>0.84</b>	0.81	0.77	0.00	0.00	0.00	0.01	0.01	0.01	0.07	0.06	0.06
100k $_{kbu}^r$	0.01	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.01	0.01	0.01	<b>0.12</b>	0.12	0.11
Baseline	0.13	0.12	0.10	0.58	0.57	<b>0.57</b>	0.10	0.10	0.09	0.12	0.09	0.05	0.07	<b>0.13</b>	<b>0.14</b>

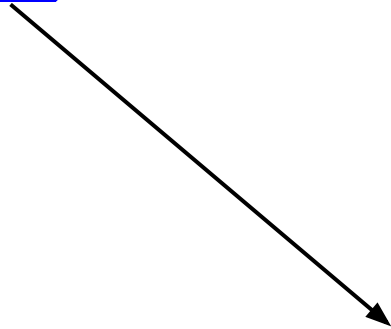
**Non-filtered training data**

Results on the hypernym discovery task for five domains

**Conclusion:** Filtering training data by domains prove to be clearly beneficial

# Word Sense Disambiguation

Kobe, which is one of Japan's largest cities, [...]



# Word Sense Disambiguation

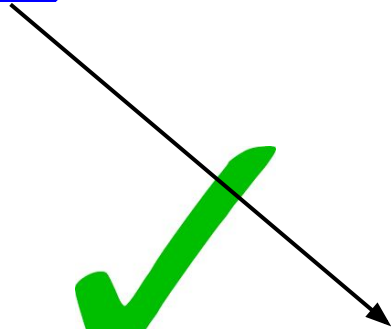
Kobe, which is one of Japan's largest cities, [...]

**X**



# Word Sense Disambiguation

Kobe, which is one of Japan's largest cities, [...]





# Word Sense Disambiguation

(Camacho-Collados et al., AIJ 2016)

## Basic idea

Select the sense which is semantically closer to the semantic representation of the whole document  
**(global context).**

$$\hat{d}(s) = \operatorname{argmax}_{d \in D} WO(\vec{N}_{ASARI_{lex}}(s), \vec{v}_{lex}(d))$$

# Word Sense Disambiguation

System	English	French	Italian	German	Spanish	Average
NASARI	86.3	<b>76.2</b>	83.7	<b>83.2</b>	82.9	<b>82.5</b>
MUFFIN	84.5	71.4	81.9	83.1	<b>85.1</b>	81.2
Babelfy	<b>87.4</b>	71.6	<b>84.3</b>	81.6	83.8	81.7
UMCC-DLSI	54.8	60.5	58.3	61.0	58.1	58.5
MFS	80.2	74.9	82.2	83.0	82.1	79.3

**Multilingual Word Sense Disambiguation using Wikipedia  
as sense inventory (F-Measure)**

# Word Sense Disambiguation

System	SemEval-2013	SemEval-2007
NASARI	66.7	66.7
NASARI+IMS	67.0	<b>68.5</b>
MUFFIN	66.0	66.0
Babelfy	65.9	62.7
UKB	61.3	56.0
UMCC-DLSI	64.7	–
Multi-Objective	<b>72.8</b>	66.0
IMS	65.3	67.3
MFS	63.2	65.8

**All-words Word Sense Disambiguation using WordNet  
as sense inventory (F-Measure)**

# Word Sense Disambiguation

System	SemEval-2013	SemEval-2007
NASARI	66.7	66.7
NASARI+IMS	67.0	<b>68.5</b>
MUFFIN	66.0	66.0
Babelfy	65.9	62.7
UKB	61.3	56.0
UMCC-DLSI	64.7	–
Multi-Objective	<b>72.8</b>	66.0
IMS	65.3	67.3
MFS	63.2	65.8

**All-words Word Sense Disambiguation using WordNet  
as sense inventory (F-Measure)**

# Word Sense Disambiguation: Empirical Comparison

(Raganato et al., EACL 2017)

- Supervised systems clearly outperform knowledge-based systems, but they **only exploit local context** (future direction -> integration of both)
- Supervised systems perform well when trained on **large amounts of sense-annotated data** (even if not manually annotated).

*Data and results available at*

<http://lcl.uniroma1.it/wsdeval/>

# Word Sense Disambiguation on textual definitions

(Camacho-Collados et al., LREC 2016)

Combination of a graph-based disambiguation system (Babelify) with NASARI to **disambiguate** the concepts and named entities of **over 35M definitions** in **256 languages**.

*Sense-annotated corpus freely available at*

<http://lcl.uniroma1.it/disambiguated-glosses/>

# Context-rich WSD



**castling** (*chess*)



*Interchanging the positions of the **king** and a **rook**.*

# Context-rich WSD



## castling (*chess*)



*Interchanging the positions of the **king** and a **rook**.*



**Castling** is a move in the game of **chess** involving a player's **king** and either of the player's original **rooks**.



A move in which the **king** moves two **squares** towards a **rook**, and the **rook** moves to the other side of the **king**.



# Context-rich WSD



## castling (*chess*)



*Interchanging the positions of the **king** and a **rook**.*



***Castling** is a move in the game of **chess** involving a player's **king** and either of the player's original **rooks**.*



*A move in which the **king** moves two **squares** towards a **rook**, and the **rook** moves to the other side of the **king**.*



*Manœuvre du jeu  
d'**échecs***



***Rošáda** je zvláštní tah v  
šachu, při kterém táhne  
zároveň **král** a **věž**.*



*Spielzug im **Schach**, bei  
dem **König** und **Turm**  
einer Farbe bewegt  
werden*



*El **enroque** es un movimiento especial  
en el juego de **ajedrez** que involucra al  
**rey** y a una de las **torres** del jugador.*



***Rokade** er et  
spesialtrekk i  
**sjakk**.*



***Rok** İngilizce'de kaleye **rook**  
denmektedir.*



*Το ροκέ είναι μια ειδική **κίνηση** στο  
**σκάκι** που συμμετέχουν ο βασιλιάς  
και ένας από τους δυο **πύργους**.*

# Context-rich WSD



## castling (chess)



Interchanging the positions of the **king** and a **rook**.



**Castling** is a move in the game of **chess** involving a player's **king** and either of the player's **rooks**.



which the **king** moves two squares towards a **rook**, and the **rook** moves to the other side of the **king**.



Mancu  
d'échec



**Rošáda** je zvláštní tah v šachu, při kterém táhne zároveň **král** a **věž**.



Spielzug im **Schach**, bei dem **König** und **Turm** einer Farbe bewegt werden



El **enroque** es un movimiento especial en el juego de **ajedrez** que involucra al **rey** y a una de las **torres** del jugador.



**Rokade** er et spesialtrekk i **sjakk**.



**Rok** İngilizce'de kaleye **rook** denmektedir.



Το ροκέ είναι μια ειδική **κίνηση** στο **σκάκι** που συμμετέχουν ο βασιλιάς και ένας από τους δυο **πύργους**.

# Context-rich WSD exploiting parallel corpora

(Delli Bovi et al., ACL 2017)

Applying the same method to provide **high-quality sense annotation from parallel corpora** (Europarl): 120M+ sense annotations for 21 languages.

**Extrinsic evaluation:** Improved performance of a standard supervised WSD system using this automatically sense-annotated corpora.

# Sense Clustering

- Current sense inventories suffer from the **high granularity** of their sense inventories.
- A meaningful clustering of senses would help **boost the performance on downstream applications** (Hovy et al., 2013)

Example:

- Parameter (computer programming) - Parameter



# Sense Clustering

## Idea

Using a clustering algorithm based on the **semantic similarity between sense vectors**

# Sense Clustering

(Camacho-Collados et al., AIJ 2016)

Measure	System type	500-pair		SemEval	
		Acc.	F1	Acc.	F1
NASARI	unsupervised	83.8	<b>70.5</b>	<b>87.4</b>	<b>63.1</b>
NASARI <sub>lexical</sub>	unsupervised	81.6	65.4	85.7	57.4
NASARI <sub>unified</sub>	unsupervised	82.6	69.5	87.2	<b>63.1</b>
NASARI <sub>embed</sub>	unsupervised	81.2	65.9	86.3	45.5
SVM-monolingual	supervised	77.4	-	83.5	-
SVM-multilingual	supervised	<b>84.4</b>	-	85.5	-
Baseline <sub>no-cluster</sub>	-	71.4	0.0	82.5	0.0
Baseline <sub>cluster</sub>	-	28.6	44.5	17.5	29.8

**Clustering of Wikipedia pages**

# Towards a seamless integration of senses in downstream NLP applications

(Pilehvar et al., ACL 2017)

**Question:** What if we apply WSD and inject sense embeddings to a standard neural classifier?

**Problems:**

# Towards a seamless integration of senses in downstream NLP applications

(Pilehvar et al., ACL 2017)

**Question:** What if we apply WSD and inject sense embeddings to a standard neural classifier?

## **Problems:**

- WSD is not perfect



# Towards a seamless integration of senses in downstream NLP applications

(Pilehvar et al., ACL 2017)

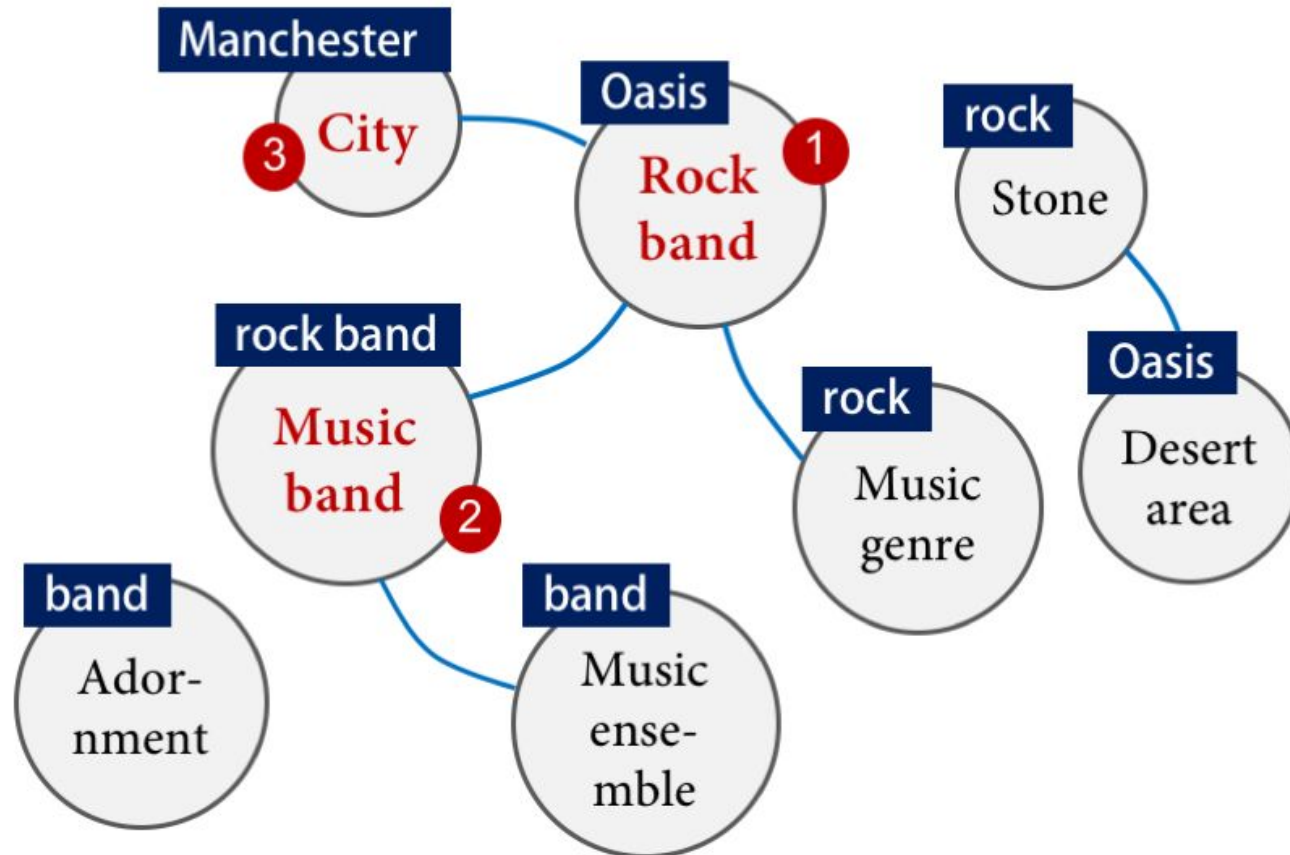
**Question:** What if we apply WSD and inject sense embeddings to a standard neural classifier?

## Problems:

- WSD is not perfect -> **Solution:** High-confidence disambiguation

# High confidence graph-based disambiguation

Oasis was a rock band formed in Manchester.



# Towards a seamless integration of senses in downstream NLP applications

(Pilehvar et al., ACL 2017)

**Question:** What if we apply WSD and inject sense embeddings to a standard neural classifier?

## Problems:

- WSD is not perfect -> **Solution:** High-confidence disambiguation
- Senses in WordNet are too fine-grained

# Towards a seamless integration of senses in downstream NLP applications

(Pilehvar et al., ACL 2017)

**Question:** What if we apply WSD and inject sense embeddings to a standard neural classifier?

## Problems:

- WSD is not perfect -> **Solution:** High-confidence disambiguation
- Senses in WordNet are too fine-grained -> **Solution:** Supersenses

# Towards a seamless integration of senses in downstream NLP applications

(Pilehvar et al., ACL 2017)

**Question:** What if we apply WSD and inject sense embeddings to a standard neural classifier?

## Problems:

- WSD is not perfect -> **Solution:** High-confidence disambiguation
- Senses in WordNet are too fine-grained -> **Solution:** Supersenses
- WordNet lacks coverage

# Towards a seamless integration of senses in downstream NLP applications

(Pilehvar et al., ACL 2017)

**Question:** What if we apply WSD and inject sense embeddings to a standard neural classifier?

## Problems:

- WSD is not perfect -> **Solution:** High-confidence disambiguation
- Senses in WordNet are too fine-grained -> **Solution:** Supersenses
- WordNet lacks coverage -> **Solution:** Use of Wikipedia

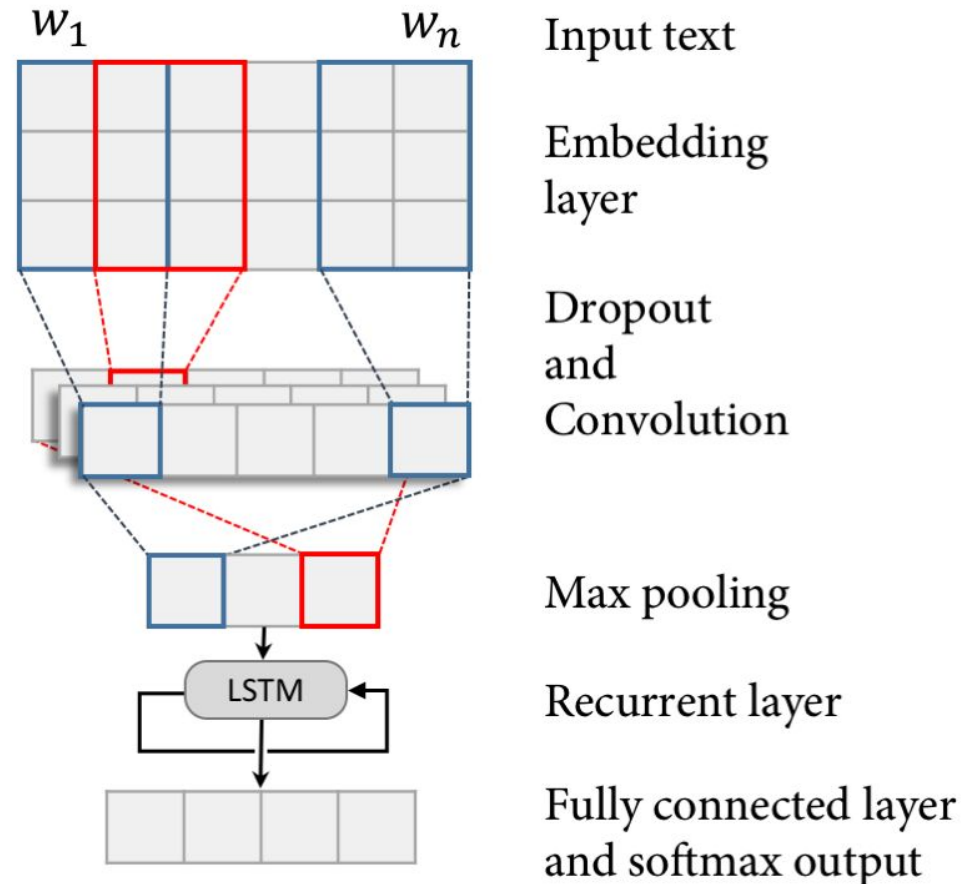
# Tasks: Topic categorization and sentiment analysis (polarity detection)

**Topic categorization:** Given a text, assign it a label (i.e. topic).

**Polarity detection:** Predict the sentiment of the sentence/review as either positive or negative.

# Classification model

Standard CNN classifier  
inspired by Kim (2014)





# Sense-based vs. word-based: Conclusions

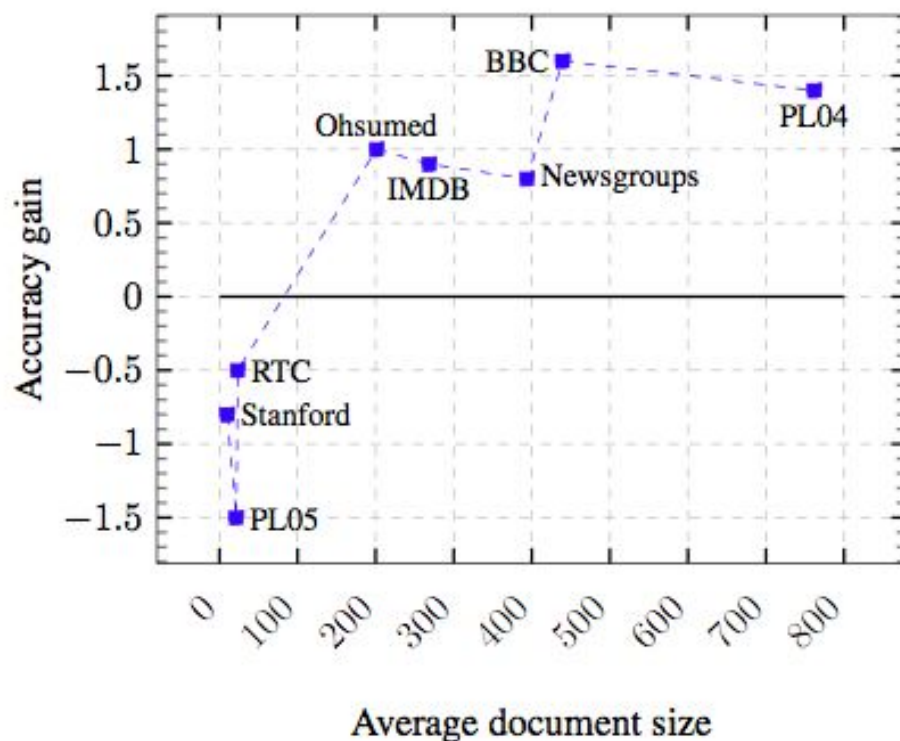
- Coarse-grained senses (*supersenses*) better than fine-grained senses.

# Sense-based vs. word-based: Conclusions

- Coarse-grained senses (*supersenses*) better than fine-grained senses.
- Sense-based **better** than word-based... when the **input text is large enough**

# Sense-based vs. word-based:

Sense-based **better** than word-based... when the **input text is large enough**:



# Why does the input text size matter?

- Graph-based WSD works better in larger texts (Moro et al. 2014; Raganato et al. 2017)
- Disambiguation increases sparsity

# Conclusions of the talk

- Novel approach to **represent concepts and entities in a multilingual vector space (NASARI)**.
- These knowledge-based sense representations can be easily integrated in several applications, acting as a glue for **combining corpus-based information and knowledge from lexical resources**, while enabling:
  - **Multilinguality**
  - **Work at the deeper sense level**

For more information on other sense-based representations and their applications:

- ACL 2016 Tutorial on “**Semantic representations of word senses and concepts**”: [http://acl2016.org/index.php?article\\_id=58](http://acl2016.org/index.php?article_id=58)
- EACL 2017 workshop on “Sense, Concept and Entity Representations and their Applications”: <https://sites.google.com/site/senseworkshop2017/>

Thank you!

Questions please!

*collados@di.uniroma1.it*

**CLASSIFIED**

# Secret Slides

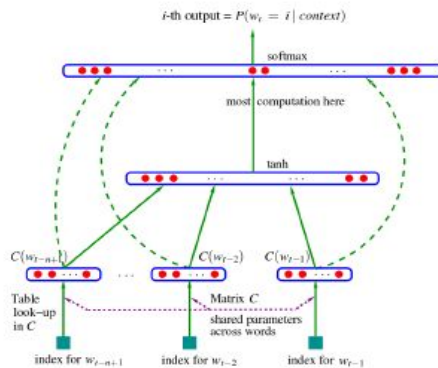


# Word vector space models

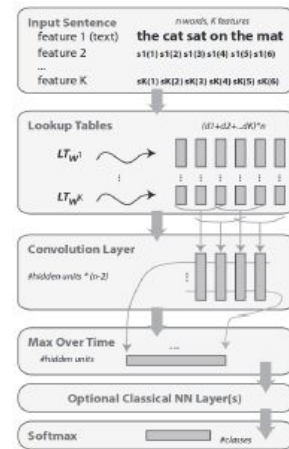
Words are represented as vectors: semantically similar words are close in the space



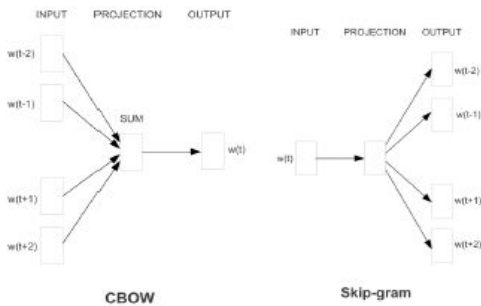
# Neural networks for learning word vector representations from text corpora -> word embeddings



Bengio et al. (2003)



Collobert & Weston (2008)

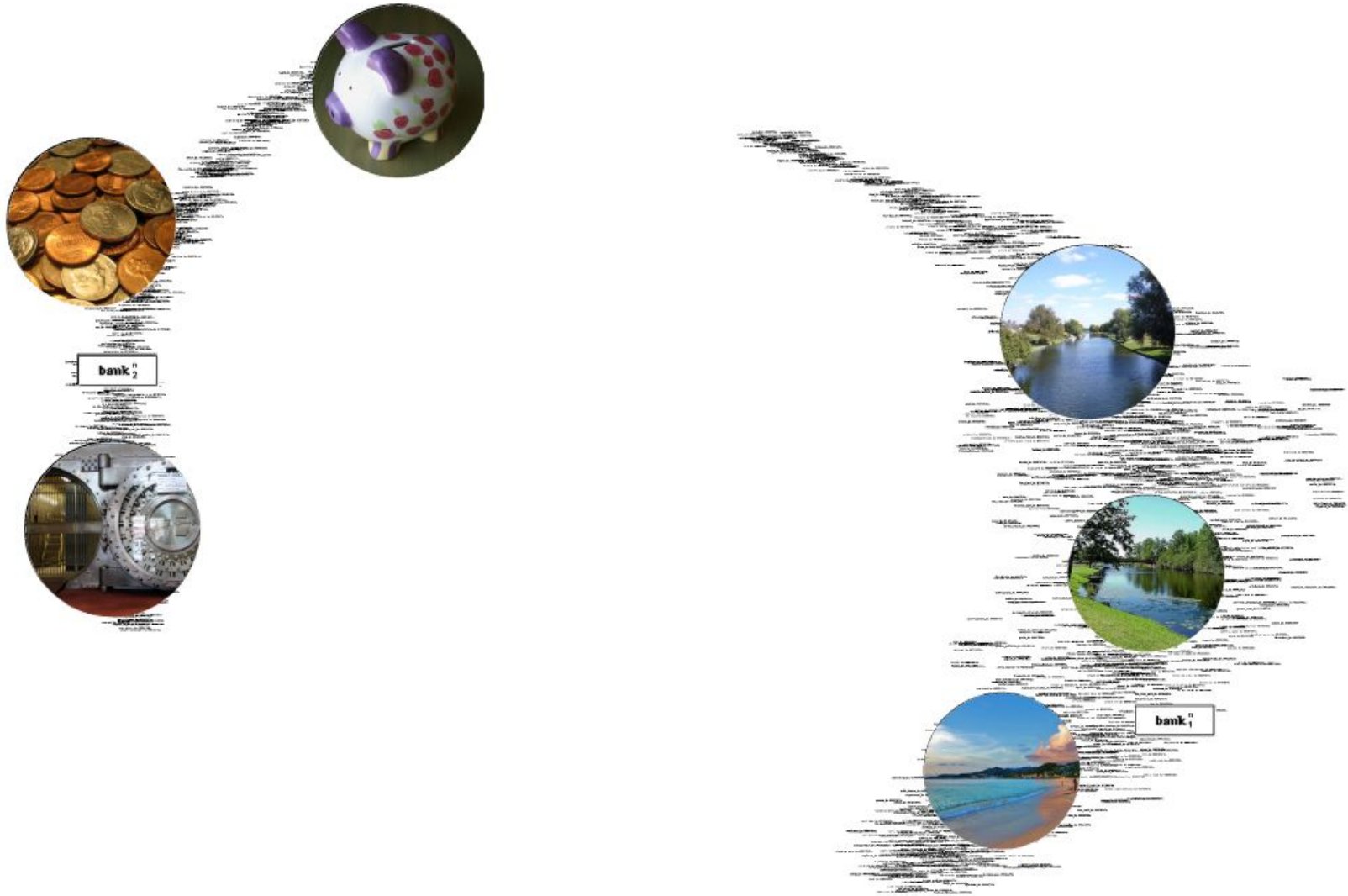


Mikolov et al. (2013)

Probability and Ratio	$k = \text{solid}$	$k = \text{gas}$	$k = \text{water}$
$P(k \text{ice})$	$1.9 \times 10^{-4}$	$6.6 \times 10^{-5}$	$3.0 \times 10^{-3}$
$P(k \text{steam})$	$2.2 \times 10^{-5}$	$7.8 \times 10^{-4}$	$2.2 \times 10^{-3}$
$P(k \text{ice})/P(k \text{steam})$	8.9	$8.5 \times 10^{-2}$	1.36

Pennington et al. (2014)

# Key goal: obtain sense representations



# NASARI semantic representations

- NASARI 1.0 (April 2015): *Lexical and unified vector representations for WordNet synsets and Wikipedia pages for English.*

**José Camacho Collados**, Mohammad Taher Pilehvar and Roberto Navigli. *NASARI: a Novel Approach to a Semantically-Aware Representation of Items*. **NAACL 2015**, Denver, USA, pp. 567-577.

# NASARI semantic representations

- NASARI 1.0 (April 2015): *Lexical and unified vector representations for WordNet synsets and Wikipedia pages for English.*

**José Camacho Collados**, Mohammad Taher Pilehvar and Roberto Navigli. *NASARI: a Novel Approach to a Semantically-Aware Representation of Items*. **NAACL 2015**, Denver, USA, pp. 567-577.

- NASARI 2.0 (August 2015): **+ Multilingual extension.**

**José Camacho Collados**, Mohammad Taher Pilehvar and Roberto Navigli. *A Unified Multilingual Semantic Representation of Concepts*. **ACL 2015**, Beijing, China, pp. 741-751.

# NASARI semantic representations

- NASARI 1.0 (April 2015): *Lexical and unified vector representations for WordNet synsets and Wikipedia pages for English.*

**José Camacho Collados**, Mohammad Taher Pilehvar and Roberto Navigli. *NASARI: a Novel Approach to a Semantically-Aware Representation of Items*. **NAACL 2015**, Denver, USA, pp. 567-577.

- NASARI 2.0 (August 2015): + Multilingual extension.

**José Camacho Collados**, Mohammad Taher Pilehvar and Roberto Navigli. *A Unified Multilingual Semantic Representation of Concepts*. **ACL 2015**, Beijing, China, pp. 741-751.

- NASARI 3.0 (March 2016): + **Embedded representations, new applications.**

**José Camacho Collados**, Mohammad Taher Pilehvar and Roberto Navigli. *Nasari: Integrating explicit knowledge and corpus statistics for a multilingual representation of concepts and entities*. **Artificial Intelligence Journal**, 2016, 240, 36-64.



BabelNet

Nome

# BabelNet

ENTRA REGISTRATI

jaguar | ENGLISH | 4 SELEZIONATE | **TRADUCI**

PREFERENZE

Tutti | Concetti | Entità nominate | 21 risultati

## Nome



### jaguar, panther, Felis onca

A large spotted feline of tropical America similar to the leopard; in some classifications considered a member of the genus Felis

ID: 00033987n | Concetto

- ZH 美洲豹
- FR jaguar, panthère
- IT giaguaro, Panthera onca, pantera
- ES jaguar, panthera onca, pantera



### Jaguar Cars, Jaguar

Jaguar Cars is a brand of Jaguar Land Rover, a British multinational car manufacturer headquartered in Whitley, Coventry, England, owned by Tata Motors since 2008.

ID: 00688731n | Entità

- ZH 捷豹
- FR Jaguar (automobile)
- IT Jaguar
- ES Jaguar Cars, Jaguar



### Atari Jaguar, Jaguar (video game console)

The Atari Jaguar is a home video game console that was released by Atari Corporation in 1993.

ID: 02142312n | Entità

- ZH Atari Jaguar, 雅达利Jaguar
- FR Jaguar (console)
- IT Atari Jaguar
- ES Atari Jaguar



### Mac OS X v10.2, Jaguar (macos)

Mac OS X version 10.2 Jaguar is the third major release of Mac OS X, Apple's desktop and server operating system.

- ZH Mac OS X Jaguar, Mac OS X v10.2
- FR Mac OS X v10.2

# Three types of vector representations

## Three types of vector representations:

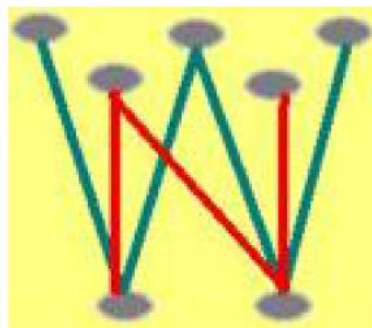
- **Lexical** (dimensions are words): Dimensions are weighted via **lexical specificity** (statistical measure based on the hypergeometric distribution)
- **Unified** (dimensions are multilingual BabelNet synsets): This representation uses a **hypernym-based clustering technique** and can be used in **cross-lingual** applications
- **Embedded** (latent dimensions)



# Key points

- **What** do we want to **represent**?
- What does "**semantic representation**" mean?
- **Why** semantic representations?
- What **problems** affect mainstream representations?
- How to **address** these problems?
- What comes **next**?

## Problem 2: word representations do not take advantage of existing semantic resources



BabelNet

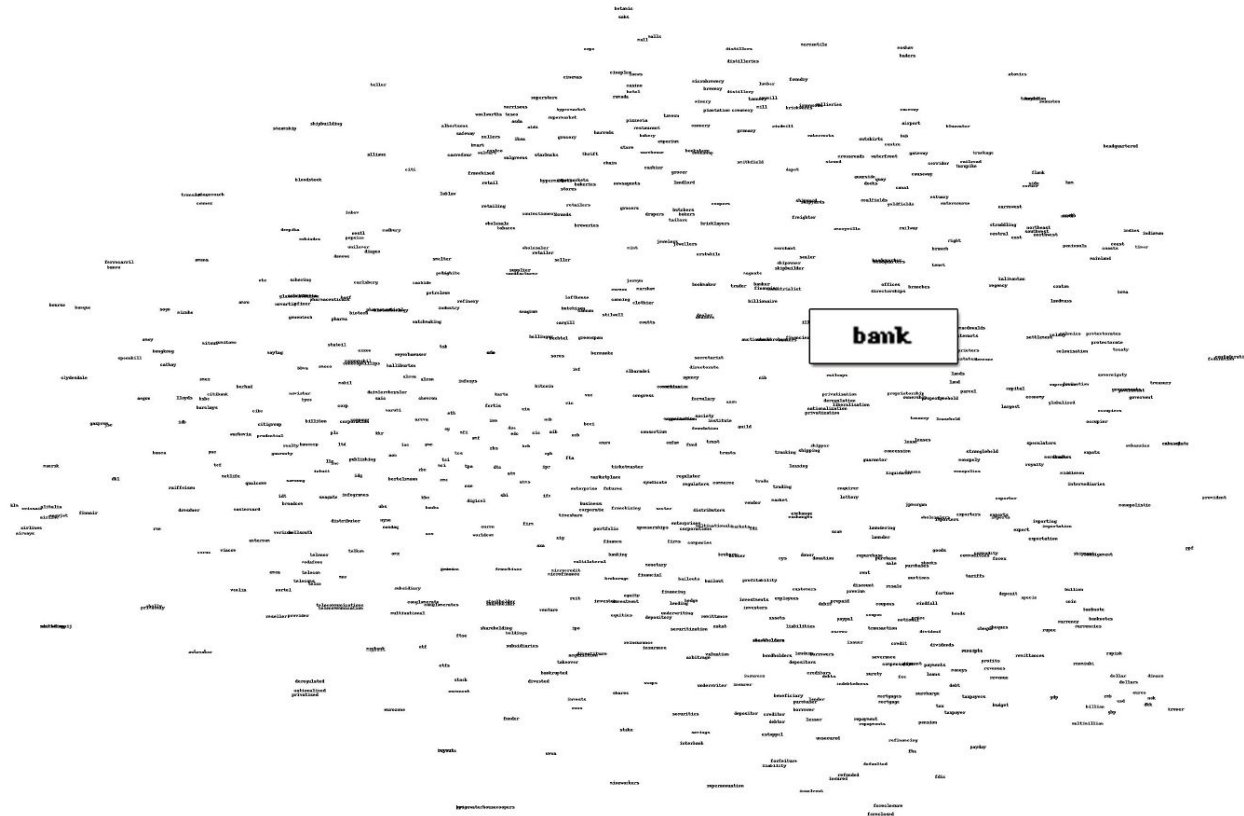


WIKIPEDIA

# Key goal: obtain sense representations

- Nome
- Verbo

Nome



on

# Named Entity Disambiguation

System	Type	F-Measure
NASARI <sub>lexical</sub>	unsupervised	87.1
DFKI	supervised	<b>88.9</b>
SUDOKU	unsupervised	87.0
e192	systems mix	86.1
MFS	–	85.7

**Named Entity Disambiguation using BabelNet as sense inventory  
on the SemEval-2015 dataset**

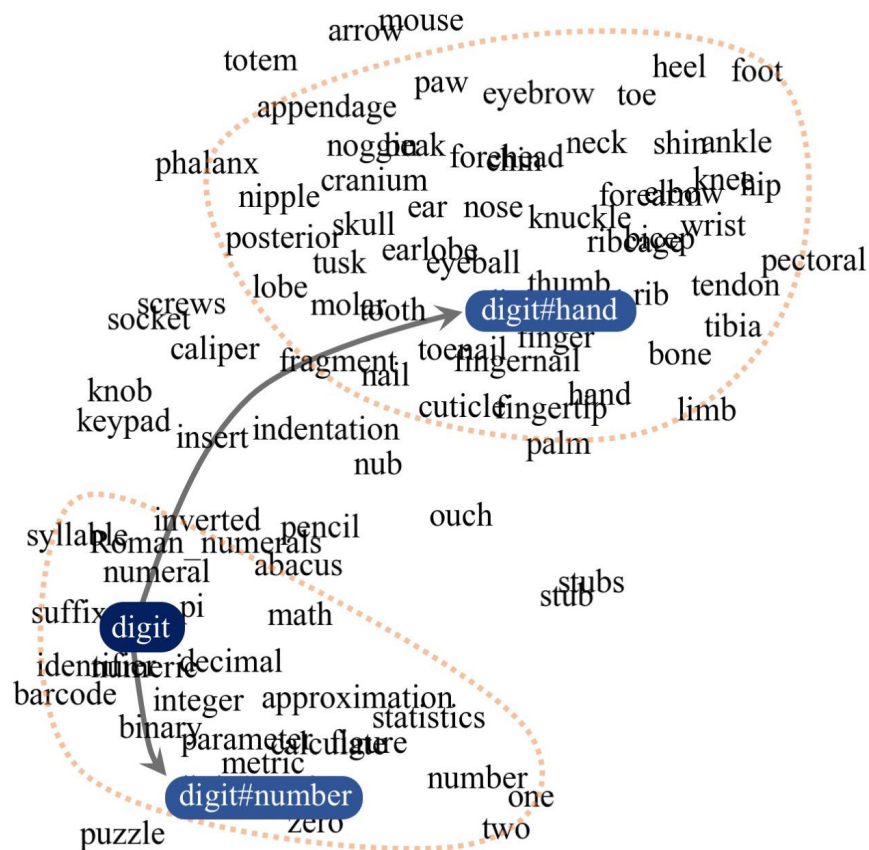
# Word Sense Disambiguation

## Open problem

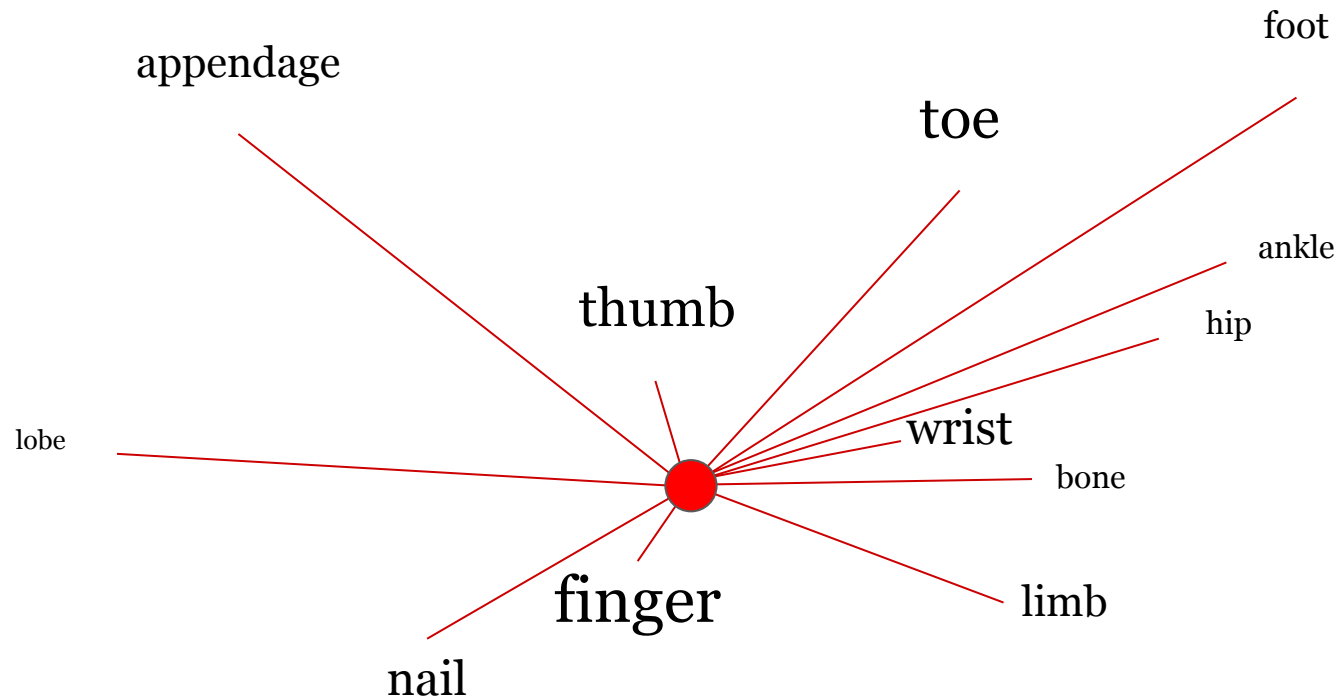
Integration of **knowledge-based** (exploiting global contexts) and **supervised** (exploiting local contexts) systems to overcome the *knowledge-acquisition bottleneck*.

# De-Conflated Semantic Representations

M. T. Pilehvar and N. Collier (EMNLP 2016)



# De-Conflated Semantic Representations



# Open Problems and Future Work

## 1. Improve evaluation

- Move from word similarity gold standards to end-to-end applications
  - Integration in Natural Language Understanding tasks (Li and Jurafsky, EMNLP 2015)
  - SemEval task? see e.g. WSD & Induction within an end user application @ SemEval 2013



# Open Problems and Future Work

2. Make semantic representations more meaningful
  - unsupervised representations are hard to inspect (clustering is hard to evaluate)
  - but also knowledge-based approaches have issues:

- e.g. top-10 closest vectors to the military sense of “company” in AutoExtend



AutoExtend
company <sub>n</sub> <sup>9</sup>
company
company <sub>n</sub> <sup>8</sup>
company <sub>n</sub> <sup>6</sup>
company <sub>n</sub> <sup>7</sup>
company <sub>v</sub> <sup>1</sup>
firm
business <sub>n</sub> <sup>1</sup>
firm <sub>n</sub> <sup>2</sup>
company <sub>n</sub> <sup>1</sup>



# Open Problems and Future Work

## 3. Interpretability

- The reason why things work or do not work is not obvious
  - E.g. avgSimC and maxSimC are based on implicit disambiguation that improves word similarity, but is not proven to disambiguate well
  - Many approaches are tuned to the task
- Embeddings are difficult to interpret and debug

# Open Problems and Future Work

4. Link the representations to rich semantic resources like WikiData and BabelNet
  - Enabling applications that can readily take advantage of huge amounts of multilinguality and information about concepts and entities
  - Improving the representation of low-frequency/isolated meanings

# Open Problems and Future Work

5. Scaling semantic representations to sentences and documents
  - Sensitivity to word order
  - Combine vectors into syntactic-semantic structures
  - Requires disambiguation, semantic parsing, etc.
  - Compositionality

# Open Problems and Future Work

6. Addressing multilinguality
  - a key trend in today's NLP research
  - We are already able to perform POS tagging and dependency parsing in dozens of languages
    - Also mixing up languages

# Open Problems and Future Work

- We can perform Word Sense Disambiguation and Entity Linking in hundreds of languages
  - Babelfy (Moro et al. 2014)
  - but with only a few sense vector representations
- Now: it is crucial that sense and concept representations are language-independent
- Enabling comparisons across languages
- Also useful in semantic parsing

# Open Problems and Future Work

- Representations are most of the time evaluated in English
  - single words only
- It is important to evaluate sense representations in other languages and across languages
  - Check out the SemEval 2017 Task 2: multilingual and cross-lingual semantic word similarity (multilwords, entities, domain-specific, slang, etc.)

# Open Problems and Future Work

7. Integrate sense representations into Neural Machine Translation
  - Previous results in the 2000s working on semantically-enhanced SMT are not very encouraging
  - However, many options have not been considered