# Multilinguality for free, or why you should care about linking vector representations to (BabelNet) synsets

Roberto Navigli

DIPARTIMENTO
DI INFORMATICA
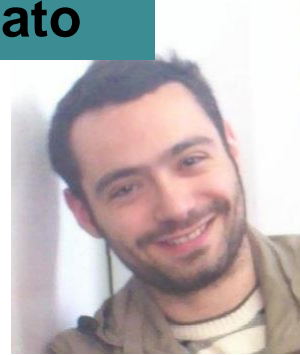
SAPIENZA
UNIVERSITÀ DI ROMA

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

19th October 2017 – Hannover, Germany
4th Alexandria workshop

Alessandro Raganato

Claudio Delli Bovi

Taher Pilehvar

Ignacio Iacobacci

José Camacho Collados

Massimiliano Mancini

**Multilinguality for free, or why you should care about linking vector representations to (BabelNet) synsets**
Roberto Navigli

20/10/2017
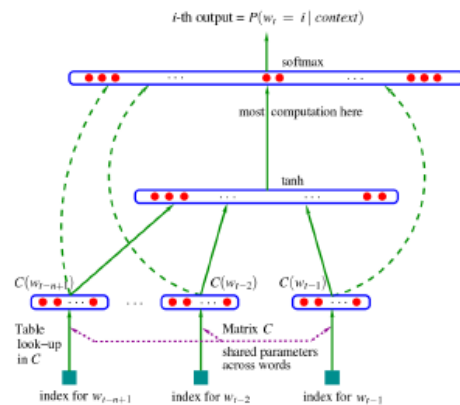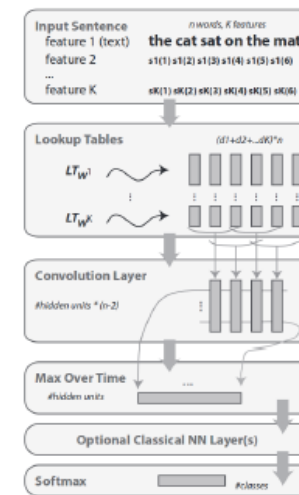
2

# How to represent words and word senses?

- Vectors provide a representation which is easy to use, visualize and combine
  - Excellent survey (Turney and Pantel, 2010)

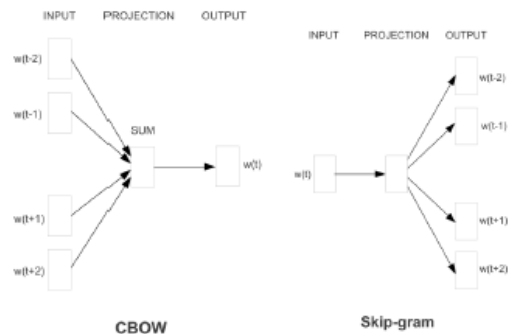**Multilinguality for free, or why you should care about linking vector representations to (BabelNet) synsets**
Roberto Navigli

20/10/2017

3

# Much work on vector representations of meaning



Bengio et al. (2003)



Collobert & Weston (2008)



CBOW          Skip-gram

Mikolov et al. (2013)

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

Pennington et al. (2014)

**Multilinguality for free, or why you should care about**   20/10/2017                                                  4
**linking vector representations to (BabelNet) synsets**
Roberto Navigli

# Problem: word representations cannot capture polysemy



**Multilinguality for free, or why you should care about** 20/10/2017
**linking vector representations to (BabelNet) synsets**
Roberto Navigli

5

# Problem: word representations cannot capture polysemy

**Multilinguality for free, or why you should care about**   20/10/2017
**linking vector representations to (BabelNet) synsets**
Roberto Navigli

6

# Problem: word representations cannot capture polysemy

**Multilinguality for free, or why you should care about linking vector representations to (BabelNet) synsets**
Roberto Navigli

20/10/2017

7

# Why should we care?

- With **word embeddings** we would have:

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

pollen ⟷ refinery

plant



(Neelakantan et al. 2014)

**Multilinguality for free, or why you should care about** 20/10/2017
**linking vector representations to (BabelNet) synsets**
Roberto Navigli

8

# Why should we care?

- With **sense embeddings**, instead:

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

pollen $\longleftrightarrow$ refinery

plant[1] $\longleftrightarrow$ plant[2]





(Neelakantan et al. 2014)

**Multilinguality for free, or why you should care about linking vector representations to (BabelNet) synsets**
Roberto Navigli

20/10/2017

9

# Solution: distinct representation for each word's meaning



**Word** vector space model

**Sense** vector space model

**Multilinguality for free, or why you should care about linking vector representations to (BabelNet) synsets**
Roberto Navigli

20/10/2017

12

# Solution: distinct representation for each word's meaning

# Where are we?

- Motivation for our work: word vs. sense representations
- Approach 1: SensEmbed (latent)
  - monolingual, but replicable in any language
- Approach 2: NASARI (explicit and latent versions)
  - monolingual and multilingual
- Approach 3: SW2V - Modeling words and senses jointly (latent)
  - in between
- Industrial applications @ Babelscape
- Conclusions

**Multilinguality for free, or why you should care about linking vector representations to (BabelNet) synsets**
Roberto Navigli

20/10/2017

22

# Latent representation of word senses: SensEmbed

Iacobacci, Pilehvar and Navigli (ACL 2015)

**Multilinguality for free, or why you should care about linking vector representations to (BabelNet) synsets**
Roberto Navigli

20/10/2017

23

# Starting point: the CBOW architecture [Mikolov et al., 2013]

**Multilinguality for free, or why you should care about** 20/10/2017
**linking vector representations to (BabelNet) synsets**
Roberto Navigli

24

# Step 1: select a large corpus



...survey on the relationship between the banks and our industry , in preparation for a forthcoming forum.
...and it stands on the right bank  of the Drava River , bounded by the river to the north...
... If you have dividend or receive bank  or building society interest on which tax has been paid ,
...workplaces and unions. Corporations, banks and trusts controlled a great deal and , although machines...
...The critical decision for the banks will come if their own adviser sticks to his view of the costs.
countryside of high hedges and tall earth banks with trees on top. The heavily wooded area was criss-crossed...

**Multilinguality for free, or why you should care about**   20/10/2017                                                                25
**linking vector representations to (BabelNet) synsets**
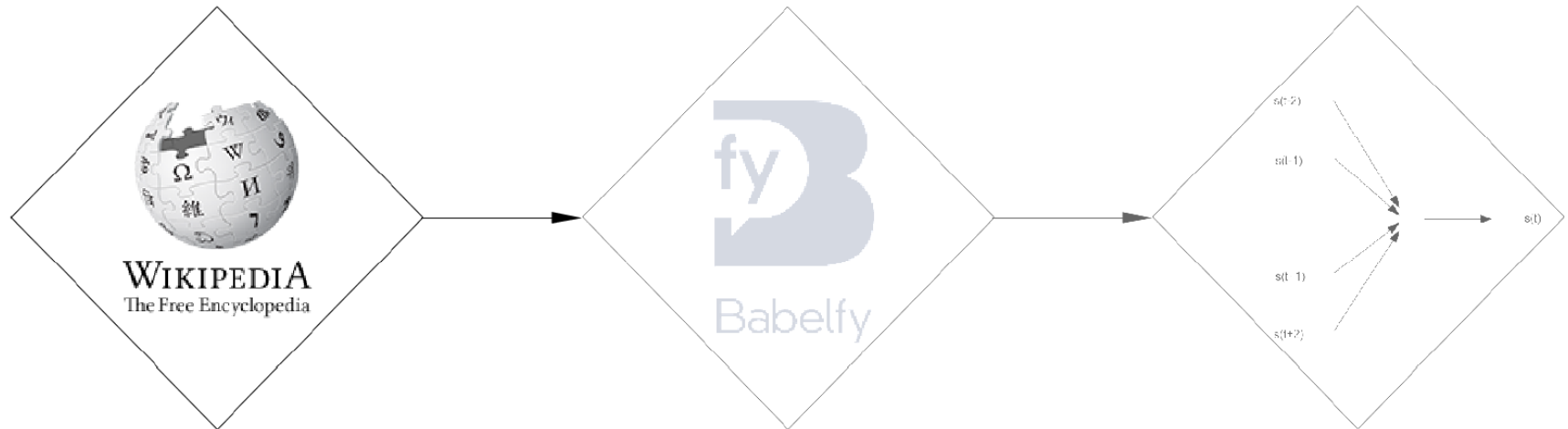Roberto Navigli
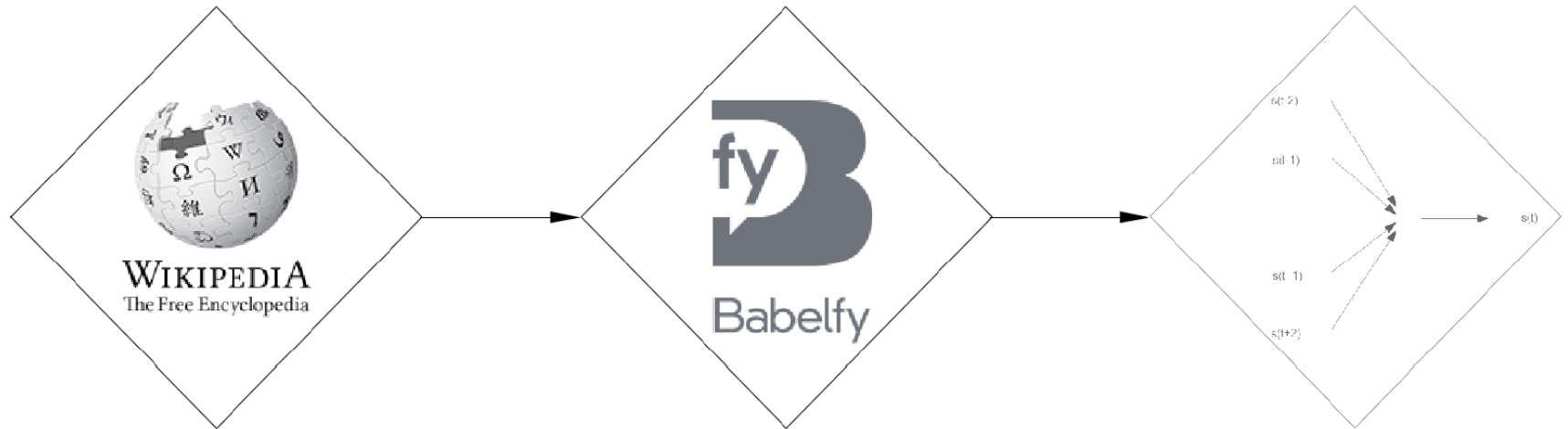
# Step 2: identify all the occurrences of a target word



...survey on the relationship between the **banks** and our industry , in preparation for a forthcoming forum.
...and it stands on the right **bank** of the Drava River , bounded by the river to the north...
... If you have dividend or receive **bank** or building society interest on which tax has been paid ,
...workplaces and unions. Corporations, **banks** and trusts controlled a great deal and , although machines...
...The critical decision for the **banks** will come if their own adviser sticks to his view of the costs.
countryside of high hedges and tall earth **banks** with trees on top. The heavily wooded area was criss-crossed...

**Multilinguality for free, or why you should care about** 20/10/2017
**linking vector representations to (BabelNet) synsets**
Roberto Navigli

26

# Step 3: disambiguate each target word occurrence



...survey on the relationship between the **banks** and our industry , in preparation for a forthcoming forum.
...and it stands on the right **bank** of the Drava River , bounded by the river to the north...
... If you have dividend or receive **bank** or building society interest on which tax has been paid ,
...workplaces and unions. Corporations, **banks** and trusts controlled a great deal and , although machines...
...The critical decision for the **banks** will come if their own adviser sticks to his view of the costs.
countryside of high hedges and tall earth **banks** with trees on top. The heavily wooded area was criss-crossed...

**Multilinguality for free, or why you should care about** 20/10/2017
**linking vector representations to (BabelNet) synsets**
Roberto Navigli

27

# Step 4: train CBOW with senses as targets



...survey on the relationship between the **banks** and our industry , in preparation for a forthcoming forum.
...and it stands on the right **bank** of the Drava River , bounded by the river to the north...
... If you have dividend or receive **bank** or building society interest on which tax has been paid ,
...workplaces and unions. Corporations, **banks** and trusts controlled a great deal and , although machines...
...The critical decision for the **banks** will come if their own adviser sticks to his view of the costs.
countryside of high hedges and tall earth **banks** with trees on top. The heavily wooded area was criss-crossed...

-2.19067  1.16642  -1.91385  -0.269672  0.712771  -0.623024  -3.20115  0.560895  0.891554  0.145258  1.26956  -0.221078
-0.0733777 2.08072  -3.30558  -0.727272  -0.902202  -1.84578  -1.38985  -0.0791954 0.989769  -1.34631 1.10242  -1.59836
-1.37341  -1.42038  0.238941  -2.98729  -0.730938  0.267584  0.0560677  -0.722721 2.23752  -2.99094  -1.45598  -0.645446
0.278277  2.28877  -0.926191  2.89934  -1.17254  1.38449  2.38617  -0.0838845  -1.80698  0.622097  0.223875  0.870654
-0.33808 -0.41957

1.16672 0.811884 -0.115492 -2.59049 -1.50286 1.2536 1.44281 0.0136615 0.131499 2.04445 -0.425782 1.29676 0.0996086
1.52687 -0.0951281 -0.715488 -0.71172 0.453871 1.08481 1.55074 0.385158 -0.116754 -0.582987 -1.56923 -0.488404
-1.07999 0.0447149 -0.733387 0.765212 2.67995 2.51105 0.192151 1.49743 2.91849 1.86901 0.23101 0.381663 1.20355
0.126758 1.57204 -0.372069 -2.45076 0.514557 -1.4028 -1.20396 0.726036 2.41265 -0.104843 2.26862 1.21729

**Multilinguality for free, or why you should care about** 20/10/2017
**linking vector representations to (BabelNet) synsets**
Roberto Navigli

28

# Setup – Sense inventory:
# BabelNet (Navigli and Ponzetto, AI Journal 2012)

- We used BabelNet, a merger of WordNet, Wikipedia, Wiktionary, OmegaWiki and other knowledge resources
- Why?
  - An extension of the lexical-semantic knowledge model of WordNet
  - Wide coverage: 271 languages (multilingual synsets),14M synsets

## Setup – Sense inventory:
## BabelNet (Navigli and Ponzetto, AI Journal 2012)

- We used BabelNet, a merger of WordNet, Wikipedia, Wiktionary, OmegaWiki and other knowledge resources

- Why?

  – An extension of the lexical-semantic knowledge model of WordNet

  – Wide coverage: 271 languages (multilingual synsets),14M synsets

  – It integrates concepts (6M) and named entities (7.7M) seamlessly

# BabelNet is now live!

- 284 languages

- 15 million concepts and named entities

- 1.8 billion semantic relations

## BabelNet goes **live.**

**BabelNet live** (beta) is the next evolutionary stage of BabelNet, today's most far-reaching **multilingual resource** that covers **hundreds of languages** and, according to need, can be used as either an **encyclopedic dictionary**, or a **semantic network**, or a huge **knowledge base**. BabelNet live (beta) is growing continuously, thanks to being fed with **daily updates** from all the sources that go to make it up, including Wikipedia, Wiktionary, users' input, etc.

☐ Don't show me again.

| **CURRENT** VERSION (3.7) | TEST THE **LIVE** VERSION (BETA) |

brought to you by
Babelscape

# Setup – Disambiguation: Babelfy [Moro et al., TACL 2014]

- We used Babelfy for disambiguating the Wikipedia corpus
- Why?
  - The first (and only) system that performs Word Sense Disambiguation (common nouns, verbs, adjectives, adverbs) and Entity Linking (names) **jointly**



I was so  lucky  I could  drive  a  Ferrari Testarossa  !

lucky
Occurring by chance

drive
Operate or control a vehicle

Ferrari Testarossa
The Ferrari Testarossa is a 12-cylinder mid-engine sports car

We  wrote  PageRank  in  Java  .

wrote
Create code, write a computer program

PageRank
PageRank is an algorithm used by Google Search to rank websites in their

Java
A platform-independent object-oriented programming language

Legend:  Named Entities  ·  Concepts

**Multilinguality for free, or why you should care about linking vector representations to (BabelNet) synsets**
Roberto Navigli

20/10/2017

34

# Setup – Disambiguation: Babelfy [Moro et al., TACL 2014]

- We used Babelfy for disambiguating the Wikipedia corpus
- Why?
  - The first (and only) system that performs Word Sense Disambiguation (common nouns, verbs, adjectives, adverbs) and Entity Linking (names) **jointly**
  - Knowledge-based: does not need millions of sentences annotated in each language (Pilehvar and Navigli, 2015)
  - Works in arbitrary languages (271 languages)
  - Can disambiguate texts written in mixed languages (language-agnostic setting)

  - [Demo on recent news]

**Multilinguality for free, or why you should care about**   20/10/2017                                   35
**linking vector representations to (BabelNet) synsets**
Roberto Navigli

# Qualitative Evaluation

- Closest senses to different senses of ambiguous words:

| $bank_1^n$ (geographical) | $bank_2^n$ (financial) | $number_4^n$ (phone) | $number_3^n$ (acting) | $hood_1^n$ (gang) | $hood_{12}^n$ (convertible car) |
|---|---|---|---|---|---|
| $upstream_1^r$ | $commercial\_bank_1^n$ | $calls_1^n$ | $appearing_6^v$ | $tortures_5^n$ | $taillights_1^n$ |
| $downstream_1^r$ | $financial\_institution_1^n$ | $dialled_1^v$ | $minor\_roles_1^n$ | $vengeance_1^n$ | $grille_2^n$ |
| $runs_6^v$ | $national\_bank_1^n$ | $operator_{20}^n$ | $stage\_production_1^n$ | $badguy_1^n$ | $bumper_2^n$ |
| $confluence_1^n$ | $trust\_company_1^n$ | $telephone\_network_1^n$ | $supporting\_roles_1^n$ | $brutal_1^a$ | $fascia_2^n$ |
| $river_1^n$ | $savings\_bank_1^n$ | $telephony_1^n$ | $leading\_roles_1^n$ | $execution_1^n$ | $rear\_window_1^n$ |
| $stream_1^n$ | $banking_1^n$ | $subscriber_2^n$ | $stage\_shows_1^n$ | $murders_1^n$ | $headlights_1^n$ |

# Quantitative Evaluation: word similarity – results

| Measure | Dataset | | | | |
|---|---|---|---|---|---|
| | RG-65 | WS-Sim | WS-Rel | YP-130 | MEN |
| Pilehvar et al. (2013) | 0.868 | 0.677 | 0.457 | 0.710 | 0.690 |
| Zesch et al. (2008) | 0.820 | — | — | 0.710 | — |
| Collobert and Weston (2008) | 0.480 | 0.610 | 0.380 | — | 0.570 |
| Word2vec (Baroni et al., 2014) | 0.840 | 0.800 | 0.700 | — | 0.800 |
| GloVe | 0.769 | 0.666 | 0.559 | 0.577 | 0.763 |
| ESA | 0.749 | — | — | — | — |
| PMI-SVD | 0.738 | 0.659 | 0.523 | 0.337 | 0.726 |
| Word2vec | 0.732 | 0.707 | 0.476 | 0.343 | 0.665 |
| $\textsc{SensEmbed}_{closest}$ | **0.894** | 0.756 | 0.645 | **0.734** | 0.779 |
| $\textsc{SensEmbed}_{weighted}$ | 0.871 | **0.812** | **0.703** | 0.639 | **0.805** |

- State-of-the-art performance + sense-level vectors in the same space as word vectors

# Explicit representation of concepts: NASARI

Camacho-Collados, Pilehvar and Navigli

(NAACL 2015; ACL 2015;

Artificial Intelligence Journal 2016)

**Multilinguality for free, or why you should care about**   20/10/2017
**linking vector representations to (BabelNet) synsets**
Roberto Navigli

48

# Motivation

# Idea 1: collect documents about a concept/entity

- For a **given concept/entity**, the initial idea is to collect a corpus of documents (Wikipedia pages) about it

**Multilinguality for free, or why you should care about linking vector representations to (BabelNet) synsets**
Roberto Navigli

20/10/2017

50

# Idea 2: we can create 3 different vector representations

- The collected corpus will be a **subcorpus** of a given **reference corpus** (the whole Wikipedia)
- The goal is to create a vector that represents the semantics of the **concept of interest**
- Three variants:
  - Lexical vectors (having words as components)
  - Unified vectors (**language-independent**, having BabelNet synsets as components)
  - Embedded vectors (having latent dimensions)

**Multilinguality for free, or why you should care about linking vector representations to (BabelNet) synsets**
Roberto Navigli

20/10/2017

51

# Calculating lexical specificity

- Given:
  - a reference corpus of T words (Wikipedia)
  - a subcorpus of t words (our set of Wikipedia pages)
- **Goal**: find a set of terms that are peculiar to the subcorpus, but not to the whole reference corpus.
- Given a word w that occurs F and f times in the corpus and subcorpus, respectively, compute the relevance of w to the subcorpus as a function of P (X ≥ f), X being a random variable following a hypergeometric distribution with parameters F, t and T.

$$spec(T, t, F, f) = -\log_{10} P(X \geq f)$$

**Multilinguality for free, or why you should care about** 20/10/2017
**linking vector representations to (BabelNet) synsets**
Roberto Navigli
53

# NASARI: the lexical vector

- Conventional vector with words as dimensions
- Individual weights calculated using lexical specificity, by contrasting the frequencies in the subcorpus and the overall corpus (whole Wikipedia)
- **Pruning:** we keep only components with $P(X \geq f) \leq 0.01$
- **Example:** top-ranking components of 2 meanings of bank:

| Bank (financial institution) | | | Bank (geography) | | |
|---|---|---|---|---|---|
| English | French | Spanish | English | French | Spanish |
| bank | banque | banco | river | eau | banco |
| banking | bancaire | bancario | stream | castor | limnología |
| deposit | crédit | banca | bank | berge | ecología |
| credit | financier | financiero | riparian | canal | barrera |
| money | postal | préstamo | creek | barrage | estuarios |
| loan | client | entidad | flow | zone | isla |
| commercial_bank | dépôt | déposito | water | perchlorate | interés |
| central_bank | billet | crédito | watershed | humide | laguna |

**Multilinguality for free, or why you should care about linking vector representations to (BabelNet) synsets**
Roberto Navigli

20/10/2017

55

# NASARI: the unified vector

- Cluster together similar dimensions in the lexical vector
- Then re-compute the weights for the new dimensions



$v_{lex}($ bank $) = ( 0, 0, ..., \mathbf{0.9}, 0.1, 0.3, \mathbf{0.6}, ..., 0.1, \mathbf{0.8}, ..., 0 )$ <span style="color:red">lexical</span>

ocean    sea    lake

flow    body of water    watercraft

$v_u($ bank $) = (0.1, 0.3, ..., 0, 0, 0.5, 0.95, ..., 0.5, ..., 0.15 )$ <span style="color:green">semantic</span>

**Multilinguality for free, or why you should care about** 20/10/2017
**linking vector representations to (BabelNet) synsets**
Roberto Navigli

56

# We use the BabelNet taxonomy

- BabelNet provides a full-fledged taxonomy: is-a relations are available for millions of concepts and named entities (**Wikipedia Bitaxonomy**, Flati et al. ACL 2014; AIJ 2016)
  - Ferrari Testarossa *is-a* sports car
  - BabelNet *is-a* semantic network & encyclopedic dictionary

**Multilinguality for free, or why you should care about** 20/10/2017
**linking vector representations to (BabelNet) synsets**
Roberto Navigli

57

# NASARI: the unified vector

- Unified vectors have BabelNet synsets as dimensions
- Two key benefits:
  - Disambiguated dimensions
  - Smoothing
- Enables
  - Transfer of semantic knowledge across languages
  - Cross-lingual semantic comparison

| Bank (financial institution) | | | Bank (geography) | | |
|---|---|---|---|---|---|
| English | French | Spanish | English | French | Spanish |
| ‡$\text{bank}_n^2$ | ‡$\text{banque}_n^1$ | ‡$\text{banco}_n^1$ | ★$\text{stream}_n^1$ | $\text{eau}_n^1$ | $\text{inclinación}_n^9$ |
| $\text{reserve}_n^2$ | ●$\text{fonds}_n^2$ | ★$\text{Institución\_financiera}_n^1$ | $\text{river}_n^1$ | $\text{eau}^{15}$ | $\text{lago}_n^1$ |
| ★$\text{financial\_institution}_n^1$ | ◇$\text{dépôt}_n^9$ | ◇$\text{depósito}_n^{15}$ | ‡$\text{body\_of\_water}_n^1$ | $\text{excrément}_n^1$ | ‡$\text{cuerpo\_de\_agua}_n^1$ |
| ◇$\text{deposit}_n^8$ | ○$\text{emprunt}_n^2$ | †$\text{Finanzas}_n^1$ | $\text{flow}_n^1$ | $\text{castor}_n^1$ | ★$\text{arroyo}_n^1$ |
| $\text{banking}_n^2$ | $\text{paiement}_n^1$ | ●$\text{dinero}_n^2$ | $\text{course}_n^2$ | ‡$\text{étendue\_d'eau}_n^1$ | $\text{tierra}_n^{11}$ |
| †$\text{finance}_n^1$ | $\text{argent}_n^2$ | ○$\text{préstamo}_n^2$ | $\text{bank}_n^1$ | $\text{fourrure}_n^1$ | $\text{costa}_n^1$ |

# NASARI: embedded representation

- We calculate a weighted average of the word embeddings of the lexical components of the vector for a given subcorpus T (corresponding to a concept of interest):

$$E(\mathcal{T}) = \frac{\sum_{w \in \vec{v}_{lex}(\mathcal{T})} \left( \frac{1}{rank(w, \vec{v}_{lex}(\mathcal{T}))} E(w) \right)}{\sum_{w \in \vec{v}_{lex}(\mathcal{T})} \frac{1}{rank(w, \vec{v}_{lex}(\mathcal{T}))}}$$

- **Key feature:** words and senses in the same space!
- Example of closest embedded vectors:

| Bank (financial institution) | | Bank (geography) | | bank | |
|---|---|---|---|---|---|
| Closest senses | Cosine | Closest senses | Cosine | Closest senses | Cosine |
| Deposit account | 0.99 | Stream bed | 0.98 | Bank (financial institution) | 0.86 |
| Universal bank | 0.99 | Current (stream) | 0.97 | Universal bank | 0.86 |
| British banking | 0.98 | River engineering | 0.97 | British banking | 0.86 |
| German banking | 0.98 | Braided river | 0.97 | German banking | 0.85 |
| Commercial bank | 0.98 | Fluvial terrace | 0.97 | Branch (banking) | 0.85 |
| Banking in Israel | 0.98 | Bar (river morphology) | 0.97 | McFadden Act | 0.85 |
| Financial institution | 0.98 | River | 0.97 | Four Northern Banks | 0.84 |
| Community bank | 0.97 | Perennial stream | 0.96 | State bank | 0.84 |

# Experiments (Camacho-Collados et al., AI Journal 2016)

- Word similarity
- Cross-lingual similarity
  - RG-65 in English, German and French
- Word Sense Disambiguation (WSD)
  - Multilingual WSD
- Domain labeling
  - "BabelDomains: Large-Scale Domain Labeling of Lexical Resources" (Camacho-Collados and Navigli, EACL 2017)
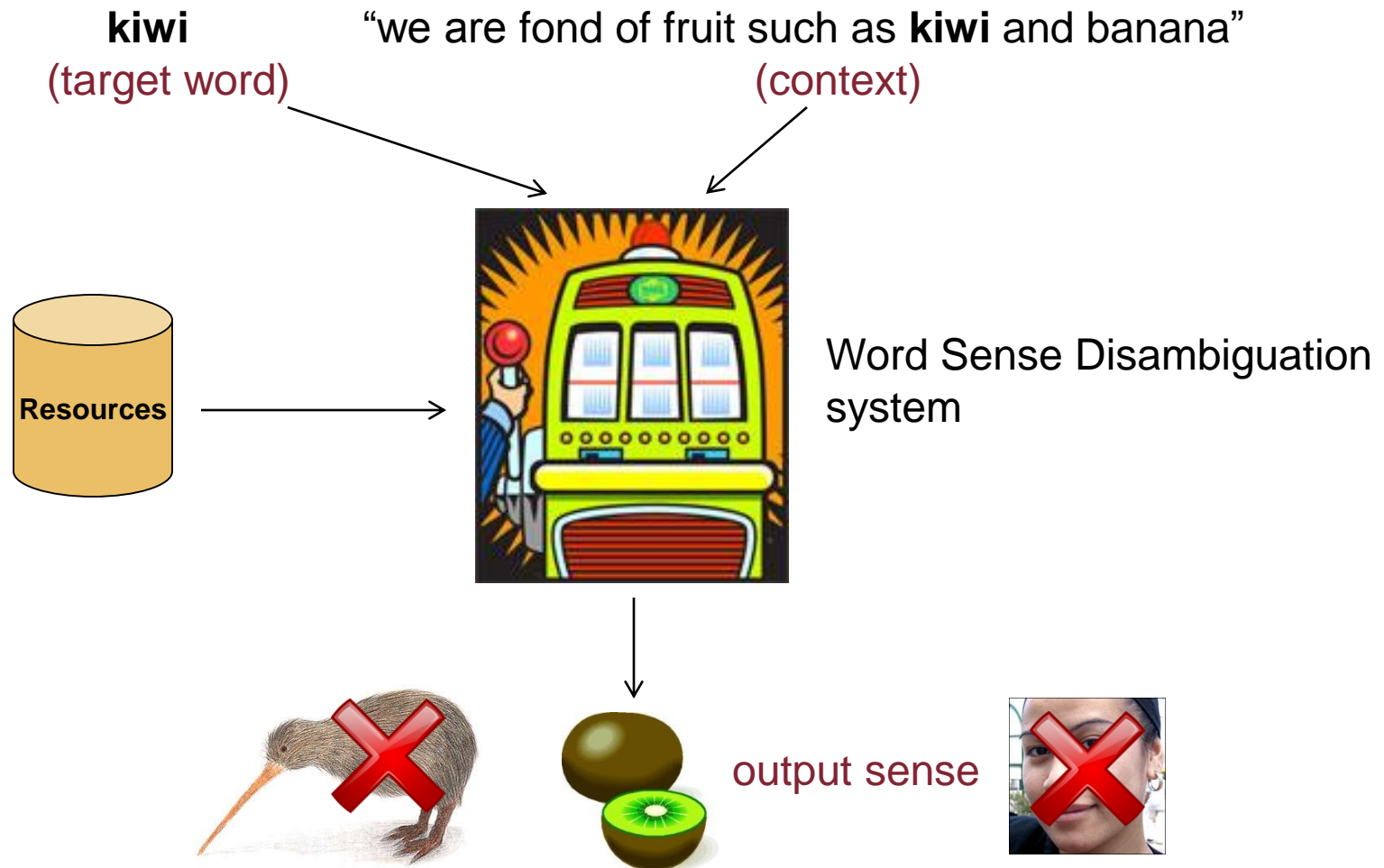- Sense clustering

**Multilinguality for free, or why you should care about linking vector representations to (BabelNet) synsets** 20/10/2017
Roberto Navigli

62

# Cross-lingual Word similarity

| English | $r$ | $\rho$ | French | $r$ | $\rho$ | German | $r$ | $\rho$ | Spanish | $r$ | $\rho$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Nasari | 0.81 | 0.78 | Nasari | **0.82** | 0.73 | Nasari | 0.69 | 0.65 | Nasari | **0.85** | 0.79 |
| Nasari$_{lexical}$ | 0.80 | 0.78 | Nasari$_{lexical}$ | 0.80 | 0.70 | Nasari$_{lexical}$ | 0.69 | 0.67 | Nasari$_{lexical}$ | **0.85** | 0.79 |
| Nasari$_{unified}$ | 0.80 | 0.76 | Nasari$_{unified}$ | **0.82** | **0.76** | Nasari$_{unified}$ | **0.71** | **0.68** | Nasari$_{unified}$ | 0.82 | 0.77 |
| Nasari$_{embed}$ | **0.82** | **0.80** | – | – | – | – | – | – | Nasari$_{embed}$ | 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 | **0.80** |
| Retrofitting | – | 0.77 | Retrofitting | – | 0.61 | Retrofitting | – | 0.60 | – | – | – |
| Nasari$_{poly-embed}$ | 0.74 | 0.77 | Nasari$_{poly-embed}$ | 0.60 | 0.69 | Nasari$_{poly-embed}$ | 0.46 | 0.52 | Nasari$_{poly-embed}$ | 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 | - |

Spearman (ρ) and Pearson (r) correlation performance of different systems on multilingual editions of the RG-65 datasets.

Comparison systems:

- SOC-PMI and PMI (Joubarne and Inkpen, 2011) – 1st and 2nd order co-occ.
- Retrofitting (Faruqui et al., 2015)
- Wiki-wup (Ponzetto and Strube, 2015)
- LSA-Wiki (Granada et al., 2014)
- Polyglot-embed (Al-Rfou et al., 2013) – emb. on wikipedias in many languages

# Understanding text: Word Sense Disambiguation

**kiwi**
(target word)

"we are fond of fruit such as **kiwi** and banana"
(context)

**Resources**

Word Sense Disambiguation system

output sense

**Multilinguality for free, or why you should care about linking vector representations to (BabelNet) synsets**
Roberto Navigli

20/10/2017

64

# Multilingual Word Sense Disambiguation

- **Sense choice:** the best sense is given by the NASARI vector closest to the text vector:

$$\hat{s} = \underset{s \in \mathcal{L}_w}{\operatorname{argmax}} WO(\vec{v}_{lex}(\mathcal{T}), \overrightarrow{\text{NASARI}}_{lex}(s))$$

- **Dataset:** the Wikipedia sense inventory for the SemEval-2013 all-words multilingual WSD task (Navigli et al. 2013) – from 1242 to 1039 annotated instances

- **Evaluation measure:** F1-measure

- **Results:**

| System | English | French | Italian | German | Spanish | Average |
|---|---|---|---|---|---|---|
| NASARI | 86.3 | **76.2** | 83.7 | **83.2** | 82.9 | **82.5** |
| MUFFIN | 84.5 | 71.4 | 81.9 | 83.1 | **85.1** | 81.2 |
| Babelfy | **87.4** | 71.6 | **84.3** | 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 |

# Latent representation of words AND senses together: SW2V

Mancini, Camacho, Iacobacci and Navigli
(CoNLL 2017)

**Multilinguality for free, or why you should care about**    20/10/2017        68
**linking vector representations to (BabelNet) synsets**
Roberto Navigli

# Objective

- Other approaches model either senses (SensEmbed) or obtain embeddings as a result of postprocessing word embeddings

- **Goal:** modeling words and senses in the same vector space

- **How:** exploiting the explicit relationships between words and senses available in BabelNet *for the words in context*

# Example



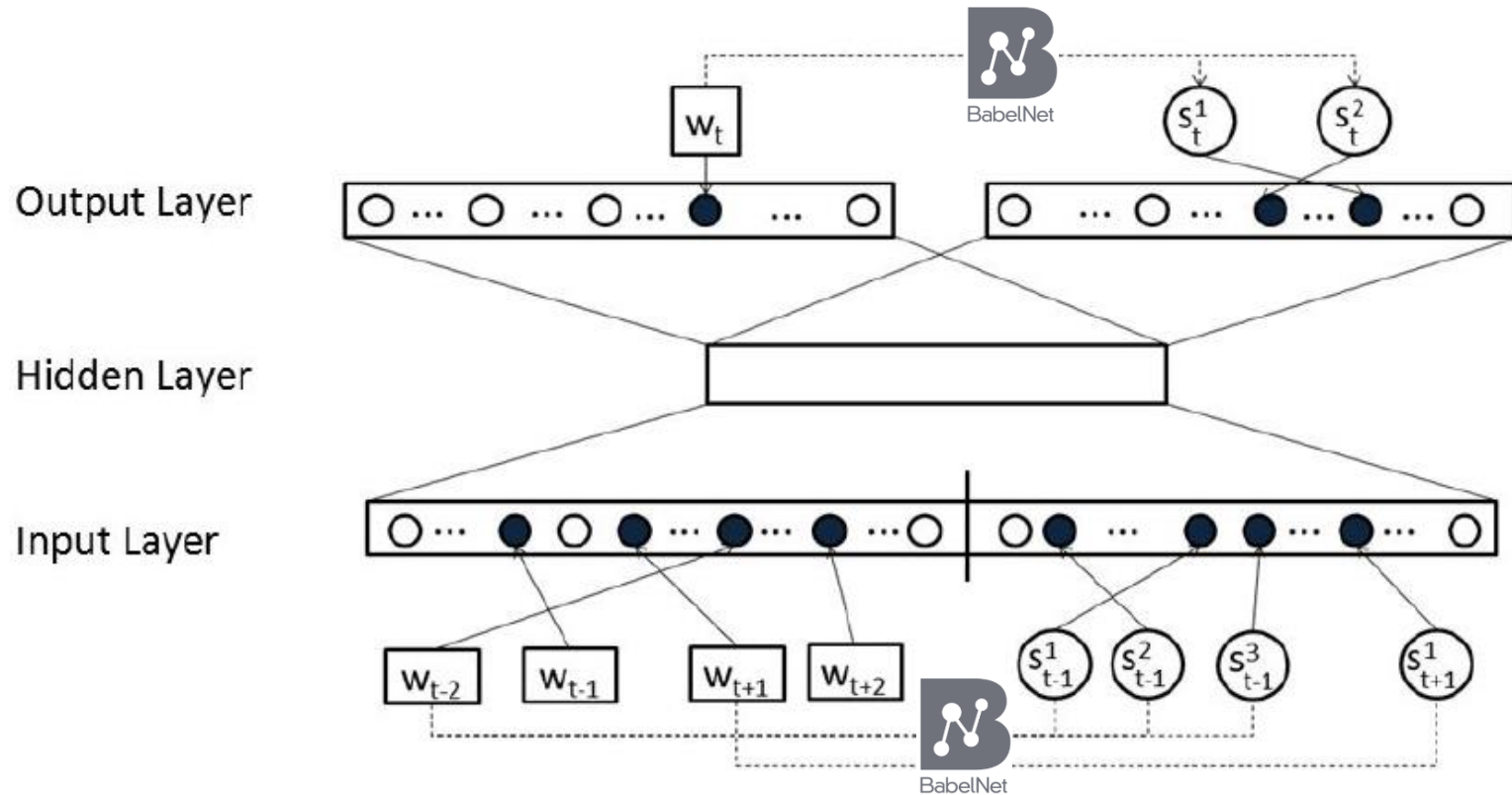He **withdrew** **money** from the **bank**

retire ✗

take out

cash

geography ✗

financial institution

building

**Multilinguality for free, or why you should care about** 20/10/2017
**linking vector representations to (BabelNet) synsets**
Roberto Navigli

70

# Extending Word2Vec with senses

$$E = -\log(p(w_t|W^t, S^t)) - \sum_{s \in St} \log(p(s|W^t, S^t))$$



Words and their associated senses used in the input and output layers

**Multilinguality for free, or why you should care about** 20/10/2017
**linking vector representations to (BabelNet) synsets**
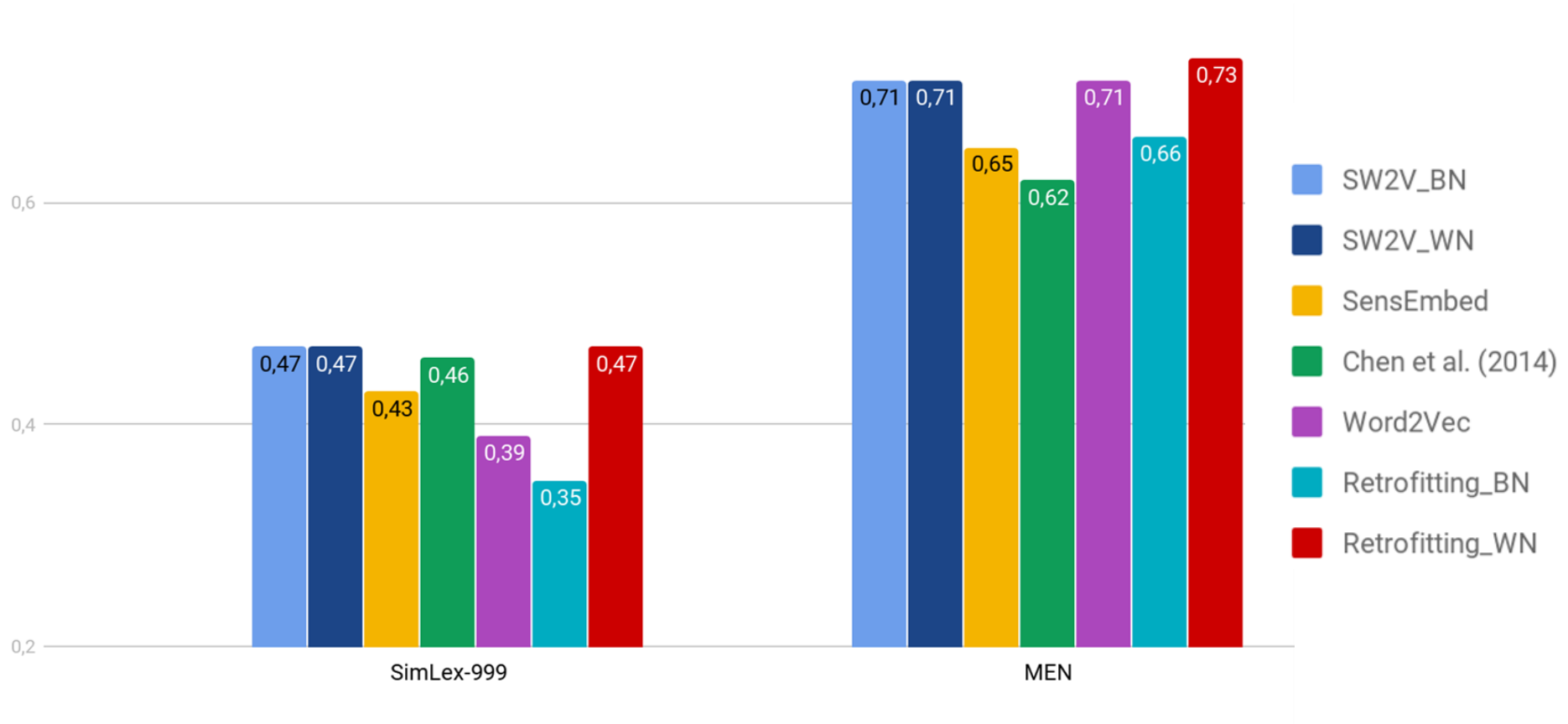Roberto Navigli

71

# Words+Senses as input and output in SW2V

- The **best configuration** is with senses only as input and words+senses as output
  - On: WS-Sim and RG-65

| | | Output | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Words | | | | Senses | | | | Both | | | | |
| | | WS-Sim | | RG-65 | | WS-Sim | | RG-65 | | WS-Sim | | RG-65 | | |
| | | $r$ | $\rho$ | $r$ | $\rho$ | $r$ | $\rho$ | $r$ | $\rho$ | $r$ | $\rho$ | $r$ | $\rho$ |
| Input | Words | 0.49 | 0.48 | 0.65 | 0.66 | 0.56 | 0.56 | 0.67 | 0.67 | 0.54 | 0.53 | 0.66 | 0.65 |
| | Senses | 0.69 | 0.69 | 0.70 | 0.71 | 0.69 | 0.70 | 0.70 | **0.74** | **0.72** | **0.71** | **0.71** | **0.74** |
| | Both | 0.60 | 0.65 | 0.67 | 0.70 | 0.62 | 0.65 | 0.66 | 0.67 | 0.65 | **0.71** | 0.68 | 0.70 |

**Multilinguality for free, or why you should care about linking vector representations to (BabelNet) synsets**
Roberto Navigli

20/10/2017

72

# Evaluation: word similarity

- All models using Wikipedia corpus (Pearson correlation)

**Multilinguality for free, or why you should care about linking vector representations to (BabelNet) synsets**
Roberto Navigli

20/10/2017

73

# Evaluation: word similarity

- All models using UMBC corpus (Pearson correlation)

**Multilinguality for free, or why you should care about** 20/10/2017
**linking vector representations to (BabelNet) synsets**
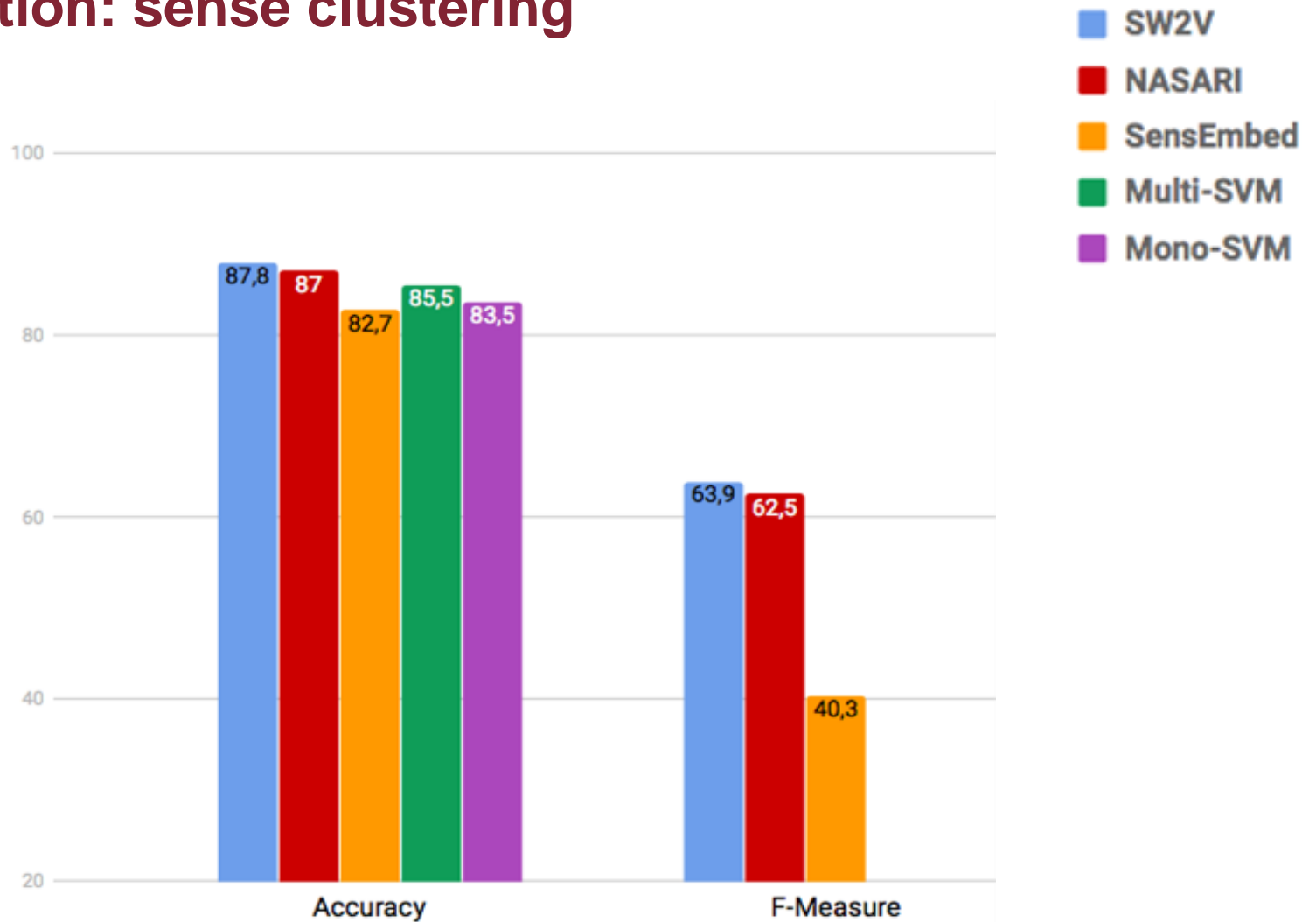Roberto Navigli

74

# Evaluation: Most Frequent Sense for WSD

- **Evaluation:** use closeness of sense vectors to word vectors to determine sense frequency
  - We can calculate the Most Frequent Sense for each word
- **Test:** Semeval-2007 and Semeval-2013 all-words WSD



F-Measure

Random Baseline: SemEval-07 24,8 / SemEval-13 34,9
AutoExtend: SemEval-07 17,6 / SemEval-13 31
SW2V: SemEval-07 39,9 / SemEval-13 54

**Multilinguality for free, or why you should care about linking vector representations to (BabelNet) synsets**
Roberto Navigli

20/10/2017
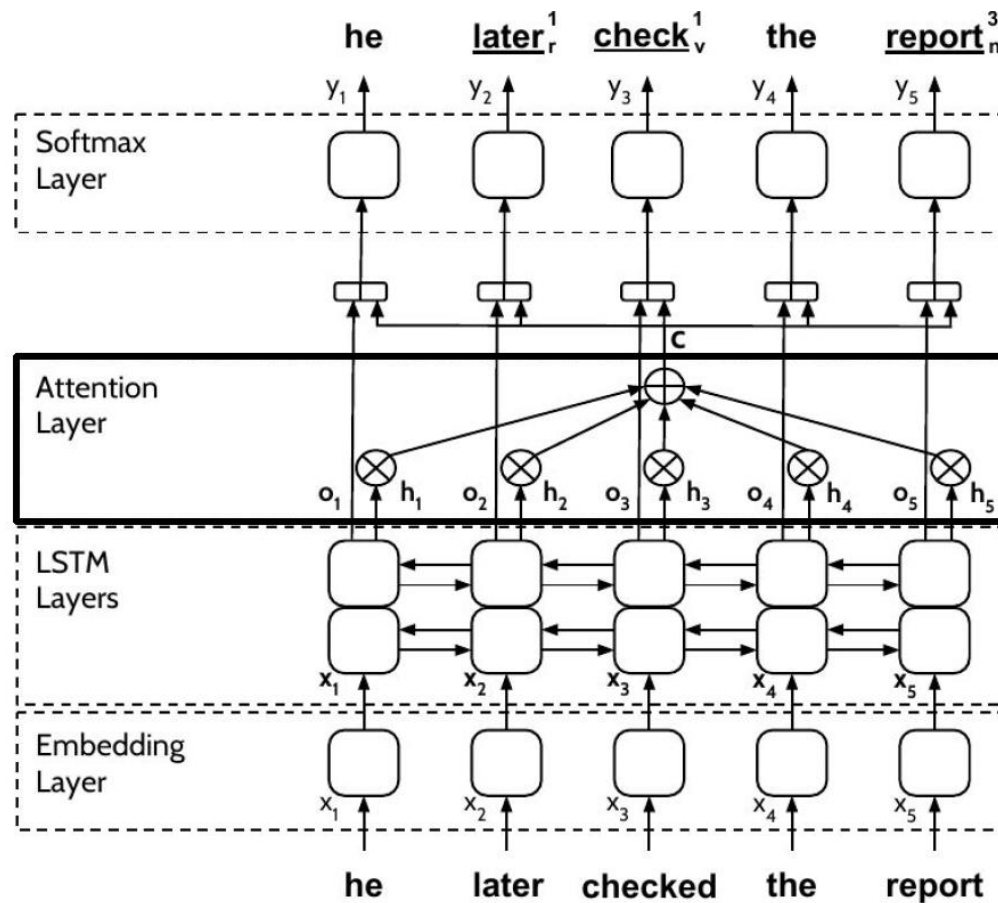
76

# Evaluation: sense clustering

- Now we can perform semantic tasks
- **Goal:** tackle the fine granularity of sense inventories
- **Evaluation datasets** from Dandala et al. (2013)
  - Highly ambiguous words from past SemEval competitions

**Multilinguality for free, or why you should care about** 20/10/2017
**linking vector representations to (BabelNet) synsets**
Roberto Navigli

77

# Evaluation: sense clustering

**Multilinguality for free, or why you should care about** 20/10/2017
**linking vector representations to (BabelNet) synsets**
Roberto Navigli

78

# Neural Models for Word Sense Disambiguation (Raganato, Delli Bovi, Navigli, EMNLP 2017)
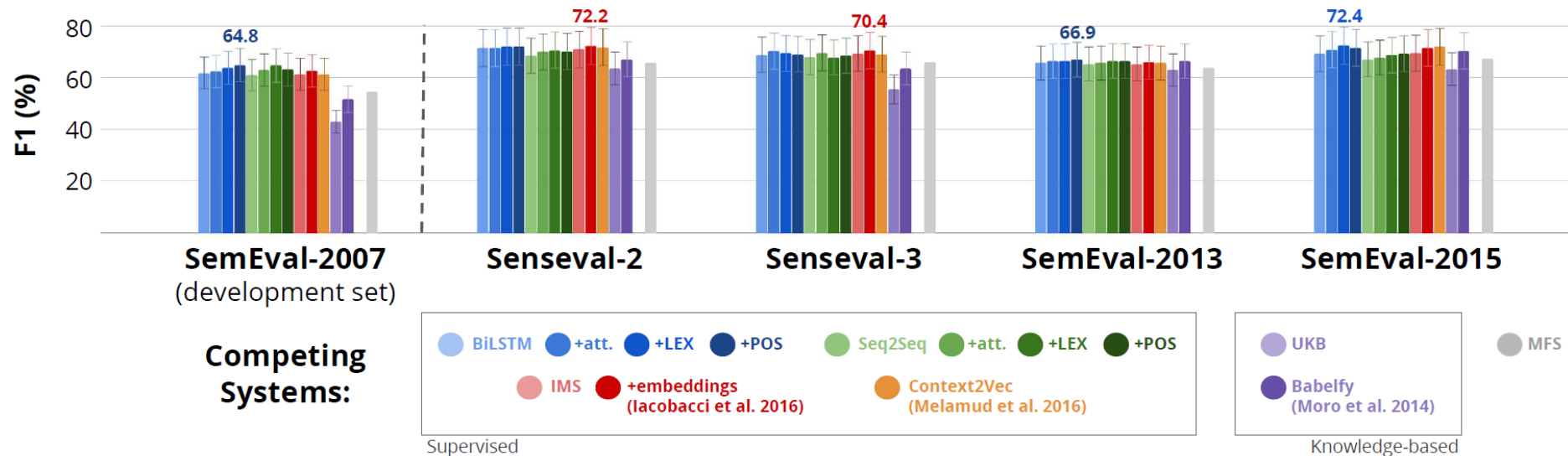
- Sequence labeling:
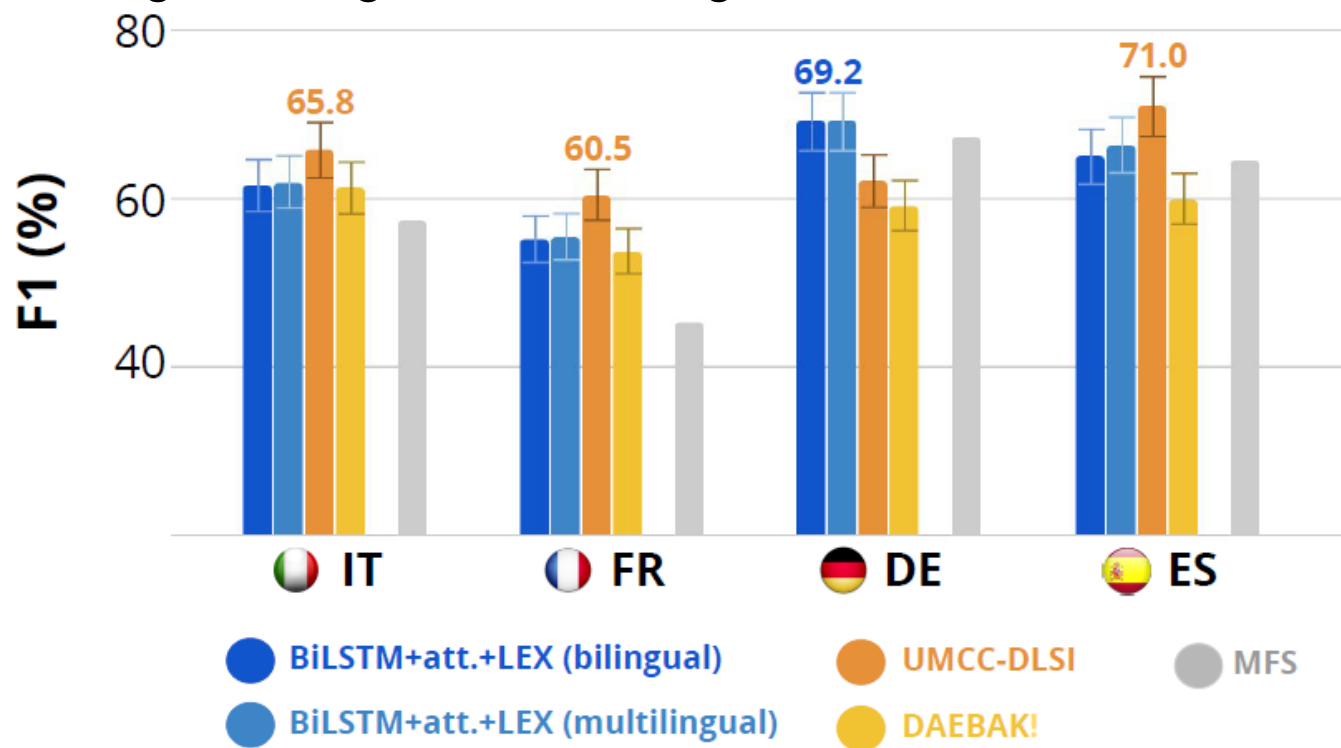


words & BabelNet concepts

# Neural Models for Word Sense Disambiguation (Raganato, Delli Bovi, Navigli, EMNLP 2017)

- Training on English (SemCor sense annotated data)
- Testing on all English Senseval & SemEval test sets

**Multilinguality for free, or why you should care about linking vector representations to (BabelNet) synsets**
Roberto Navigli

20/10/2017

80

# Neural Models for Word Sense Disambiguation (Raganato, Delli Bovi, Navigli, EMNLP 2017)

- Training on English (SemCor sense annotated data)
- **Testing on arbitrary languages** (!) – SemEval 2013
  - Using multilingual embeddings to encode words in the same space

# The future of BabelNet and related technologies

- The MultiJEDI ERC project is now over (but: the MOUSSE ERC grant just started)
  - moving to sentence representations
- However, much work still to be done in this direction
- We created a Sapienza startup, Babelscape, with the **key objective** of making BabelNet sustainable
- Income is reinvested in BabelNet and subsequent projects



Babelscape          Home   Products   About Us   Contacts

Babelscape

Multilinguality at your fingertips

# Wrapping up

- We advocated for **linking to BabelNet**
  - **SensEmbed:** lcl.uniroma1.it/sensembed
  - **NASARI:** lcl.uniroma1.it/nasari
- Monolingual vs. multilingual:
  - **Monolingual** (but no limit to which language can be used: SensEmbed, NASARI lexical/embedded)
  - **Inherently multilingual** (NASARI unified vector)
- Explicit vs. latent:
  - Explicit vectors provide **human-readable** components (NASARI lexical and unified)
  - Latent vectors are **more compact**, less sparse and **faster to process** (SensEmbed, NASARI embedded)
- Enable semantic, "translatable" output
- Move from **one language to another** seamlessly

**Multilinguality for free, or why you should care about linking vector representations to (BabelNet) synsets**
Roberto Navigli

20/10/2017

86

# Thanks or…



(grazie)

**MultiJEDI** (Starting Grant, 2011-2016) + **MOUSSE** (Consolidator Grant, 2017-2022)

**Multilinguality for free, or why you should care about** 20/10/2017
**linking vector representations to (BabelNet) synsets**
Roberto Navigli

87

**Roberto Navigli**

Linguistic Computing Laboratory
http://lcl.uniroma1.it
@RNavigli