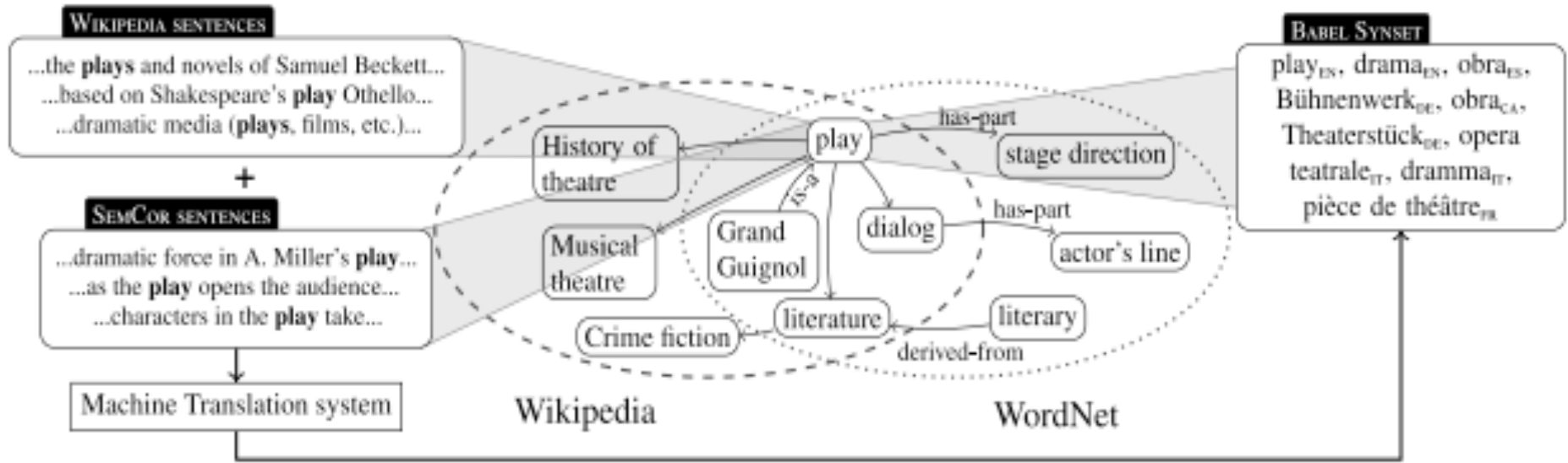


BabelNet: The automatic construction, evaluation and application of a wide-coverage multilingual semantic network

Roberto Navigli,
Simone Paolo Ponzetto

What is BabelNet

- ‘a very large, wide-coverage multilingual semantic network’
- Built from WordNet and Wikipedia
- 3 million concepts
- Mapping Accuracy of 82%
- Covers 52% of noun senses in WordNet
- Lemmas in 6 languages
 - 72.3% of synsets contain all 6



- WordNet
 - Synset(s) for each word
 - Labeled, Hierarchical relations Relations
 - Glosses

{play^{1_n}, drama^{1_N}, dramatic play^{1_N}}

play^{1_N} *is a* dramatic composition^{1_N}

Shakespeare^{1_N} *instance-of* dramatist^{1_N}

Stage direction^{1_N} *part-of* play^{1_N}

- WordNet
 - Synset(s) for each word
 - Labeled, Hierarchical relations Relations
 - Glosses

{play_N⁸, child's play_n²}

- Wikipedia
 - Page title
 - Optional label
 - Structured Relations
 - Redirect pages
 - disambiguation pages
 - internal links
 - inter-language links
 - categories

PLAY

PLAY (THEATRE)

Building BabelNet

1. Map Wikipedia to WordNet
2. Harvest multilingual lexicalizations from Wikipedia
3. Identify relations between synsets using WordNet and Wikipedia articles from all relevant languages

Mapping Wikipedia to WordNet

1. All available word senses and labeled relationships between them are pulled from WordNet
2. All encyclopedic entries and unspecified relationships between them are pulled from English Wikipedia
3. Combine them, and merge any intersecting senses
 1. If a sense exists in Wikipedia, but not in WordNet, the sense gets a null mapping
 2. If there is only one sense in Wikipedia and one sense in WordNet, create a one-to-one mapping
 3. For all remaining Wikipedia senses, map using the mapping algorithm:

$$\arg \max_{s \in \text{Senses}_{WN}(w)} p(s|w) = \arg \max_s p(s, w)$$
$$p(s, w) = \frac{\text{score}(s, w)}{\sum_{\substack{s' \in \text{Senses}_{WN}(w), \\ w' \in \text{Senses}_{Wiki}(w)}} \text{score}(s', w')}$$

Mapping Wikipedia to WordNet

$$p(s,w) = \frac{score(s,w)}{\sum_{\substack{s' \in Senses_{WN}(w), \\ w' \in Senses_{Wiki}(w)}} score(s',w')}$$

Bag-of-words

$$score(s,w) = |Ctx(s) \cap Ctx(w)| + 1$$

- Simple and fast

Graph-based

$$score(s,w) = \sum_{cw \in Ctx(w)} \sum_{s' \in Senses_{WN}(cw)} \sum_{p \in paths_{WN}(s,s')} e^{-(length(p)-1)}$$

- Exploits structural information available in WordNet and Wikipedia
- Provides mappings even when the intersection between WordNet and Wikipedia disambiguation contexts is empty.

Disambiguation Contexts

WordNet

- Synonymy
- Hypernymy/hyponymy
- Gloss (lemmatized, non-stop words)

Wikipedia

- Sense labels
- Links
- Redirections
- Categories

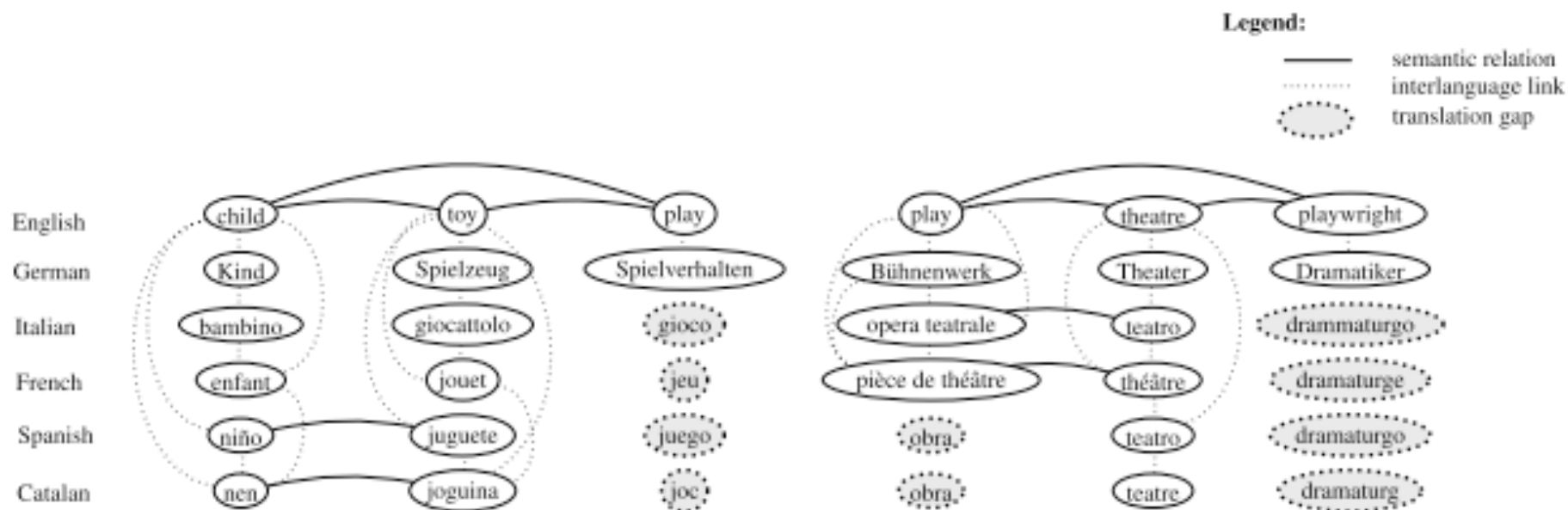
Translating Synsets

Synsets included

1. The word
2. The WordNet synset
3. Wikipedia redirections to the word
4. Other language Wikipedia articles
 1. Inter-language links
 2. Redirections to the inter-language links in the target language
5. Machine Translation using Google Translate
 1. WordNet contexts using SemCor contexts
 2. Wikipedia contexts from inter-wiki links to the page (BabelCor)

Translating Synsets

- Machine translation was not executed on named entities from Wikipedia
 - These were identified by titles containing at least 2 words with initial uppercase letters
 - 94% validation on a sample of 100 pages
- Contextless translation was performed on monosemous words from WordNet.



Finding Semantic Relations

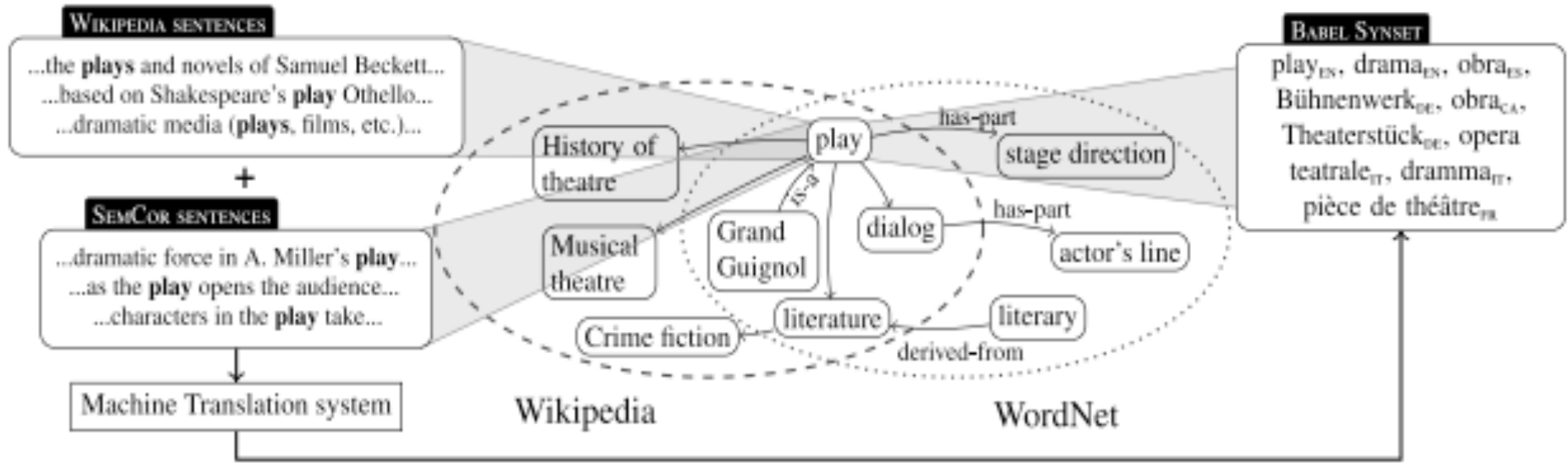
- Wikipedia:
 - Harvesting relations from hyperlinks within *all* Wikipedias in the relevant languages
 - Ex. Play (theatre) links to Acting through the German Wikipedia, but not the English Wikipedia
 - Evaluate the frequency of co-occurrence of word1 and word2 ($f_{w.w'}$) in a context.

$$\frac{2 \times f_{w.w'}}{f_w + f_{w'}}$$

Finding Semantic Relations

- WordNet:
 - Gather lemmatized, non-stop word bags of words from:
 - Synonyms, gloss words, and directly linked synsets
 - Weighting relations using the Dice coefficient to determine their strength

$$\frac{2 \times |S \cap S'|}{|S| + |S'|}$$



Evaluation: Mapping

- Gold standard of 1,000 sample wikipages
 - 505 non-empty mappings
- Experienced annotator provided the correct WordNet sense for each page
 - 2nd annotator sense tagged subset of 200 pages to evaluate annotation accuracy
 - Inter annotator agreement of 0.9

Evaluation: Mapping

- Bag-of-words:
 - $F_1 = 65.1$ (just gloss)
 - $F_1 = 75.0$ (using taxonomy and gloss)
- Graph-based:
 - $F_1 = 77.7$ (depth 2, using taxonomy and gloss)
- Baseline:
 - $F_1 = 33.5$ (most frequent sense)
 - $F_1 = 31.9$ (random selection)

Evaluation: Translation

- Automatic
 - Compared coverage of BabelNet with gold standard corpora
 - Catalan & Spanish: Multilingual Central Repository
 - French: Wordnet Libre du Français
 - German: GermaNet subset of EuroWordNet for German
 - Italian: MultiWordNet

Evaluation: Translation

$$\text{SynsetCov}(B, F) = \frac{\sum_{S_B \in F} \delta(S_B, S_F)}{|\{S_F \in F\}|}$$

$$\text{WordCov}(B, F) = \frac{\sum_{S_B \in F} \sum_{s_F \in S_F} \delta'(S_B, s_F)}{|\{S_F \in S_F : S_F \in F\}|}$$

$$\text{WordExtraCov}(B, F) = \frac{\sum_{S_\varepsilon \in \varepsilon \setminus F} \delta(S_B, S_\varepsilon)}{|\{S_F \in F\}|}$$

$$\text{WordExtraCov}(B, F) = \frac{\sum_{S_\varepsilon \in \varepsilon \setminus F} \sum_{s_\varepsilon \in S_\varepsilon} \delta'(S_B, s_\varepsilon)}{|\{S_F \in S_F : S_F \in F\}|}$$

Evaluation: Automatic Translation

Table 6

Coverage against gold-standard wordnets (percentages).

Resource	WordCov (SENSES)				SynsetCov (SYNSETS)			
	WIKI		WORDNET	BABELNET	WIKI		WORDNET	BABELNET
Method	Links	Transl.	Transl.	All	Links	Transl.	Transl.	All
Catalan	20.3	46.9	25.0	64.0	25.2	54.1	29.6	73.3
French	70.0	69.6	16.3	86.0	72.4	79.6	19.4	92.9
German	39.6	42.6	21.0	57.6	50.7	58.2	28.6	73.4
Italian	28.1	39.9	19.7	52.9	40.0	58.0	28.7	73.7
Spanish	34.4	47.9	25.2	66.4	40.7	56.1	30.0	76.6

Table 7

Extra coverage against gold-standard wordnets (percentages).

Resource	WordExtraCov (SENSES)				SynsetExtraCov (SYNSETS)			
	WIKI		WORDNET	BABELNET	WIKI		WORDNET	D..BABELNET
Method	Links	Transl.	Transl.	All	Links	Transl.	Transl.	All
Catalan	100	204	71	340	35	105	42	142
French	255	223	92	514	63	102	67	159
German	1349	940	367	2298	506	668	303	902
Italian	160	234	83	419	87	153	68	213
Spanish	214	158	56	384	48	74	30	102

Evaluation: Manual Translation

- Manual Evaluation
 - Evaluates the novel synsets
 - Random set of 600 Babel synsets
 - Manually evaluated by expert annotators

- Wordnet only synsets: 72%
- Wikipedia only synsets: 95%
- WordNet and Wikipedia synsets: 80%

How this relates to our project

- It really doesn't anymore
- We're now working on topic identification and classification of journal articles, using attributes derived from the article context.
- BabelNet can be used for topic identification
 - ... so they're kind of related