

Eric Kafe

Sense keys

Mapping

Results

Discussion

Summary

References

Mapping Wordnets on the fly with permanent sense keys

Eric Kafe

kafe@megadoc.net

MegaDoc

January 2023

EHU/UPV, San Sebastian

Spain

Sense Keys

Eric Kafe

Sense keys

Mapping

Results

Discussion

Summary

References

- **Sense keys represent one word sense:**
Permanent since PWN 1.5SC (1994)
OEWN since Edition 2021
- **Examples:**
sequoia%1:20:00::
obtrusive%5:00:00:noticeable:00
- **Sense key encoding** (PWN Manual , 2010, Senseidx):
lemma%ss_type:lex_filenum:lex_id:head_word:head_id

Almost all wordnets are mapped at the synset level
But the sense keys can reveal splits and merges

Merge

Eric Kafe

Sense keys

Mapping

Results

Discussion

Summary

References

Sense Key	PWN _{3.0}	CILI _{3.0}	CILI _{3.1}	PWN _{3.1}
baseball%1:04:00::	00471613-n:	i37881	i37882	00472688-n
baseball_game%1:04:00::	00471613-n:	i37881	i37882	00472688-n
ball%1:04:01::	00474568-n	i37882	i37882	00472688-n

Merges concern whole synsets, so should not be a problem

But ILI-to-offset mappings lose merged synsets (with *Wn*):

- `pwn30.synsets(ili="i37881")[0].translate("oewn")`
→ []
- `oewn.synsets(ili="i37882")[0].lemmas()`
→ [*ball*, *baseball*, *baseball_game*]
- `oewn.synsets(ili="i37882")[0].translate("omw-en")[0].lemmas()`
→ [*ball*]

Would need a mapping between ILI versions:

`i378813.0 ↔ i378823.1`

Split and Merge

Eric Kafe

Sense keys

Mapping

Results

Discussion

Summary

References

French	PWN _{3.0}	CILI _{3.0}	Sense Key	CILI _{3.1}	PWN _{3.1}
Aides	09570298-n	i86957	aides%1:18:00::	i86957	09593427-n
ἄϊ	09570298-n	i86957	aidoneus%1:18:00::	i86957	09593427-n
Hadès	09570298-n	i86957	hades%1:18:00::	i86957	09593427-n
Pluton	09570298-n	i86957	pluto%1:18:00::	i86958	09593643-n
Dis	09570522-n	i86958	dis%1:18:00::	i86958	09593643-n
Ἄϊ	09570522-n	i86958	orcus%1:18:00::	i86958	09593643-n
Ἄϊ	Ἄϊ	Ἄϊ	dis_pater%1:18:00::	i86958	09593643-n

- Pluto moved from Greek gods to Roman gods
- CILI mapping works only for English:
For ex. French *Pluton* doesn't move
- Would need a mapping for particular senses:
 $Pluton_{French:3.0:i86957} \leftrightarrow Pluton_{French:3.1:i86958}$

Eric Kafe

Sense keys

Mapping

Results

Discussion

Summary

References

In NLTK since v. 3.7 (2021):
Load OMW with any PWN or OEWN version

```

SENSE_INDEXsource ← { ∀ sense ∈ source : sensekey → synset_idsource }
SENSE_INDEXtarget ← { ∀ sense ∈ target : sensekey → synset_idtarget }
MAP_TO_MANY ← { ∀ synset_idsource ∈ values(SENSE_INDEXsource) : synset_idsource → ∅ }
for sensekey ∈ SENSE_INDEXsource ∩ SENSE_INDEXtarget do
    MAP_TO_MANY[synset_idsource].append(synset_idtarget)
end for
MAP_TO_ONE ← { ∀ synset_idsource ∈ MAP_TO_MANY : synset_idsource → argmax(count(synset_idtarget)) }
return MAP_TO_ONE

```

- Map to most frequent target synset
- Runs in linear time
- Works with any synset id. (f. ex. offsets or ILI)

No Splits

Eric Kafe

Sense keys

Mapping

Results

Discussion

Summary

References

Sense Key	PWN _{3.0}	CILI _{3.0}	CILI _{3.1}	PWN _{3.1}
aides%1:18:00::	09570298-n	i86957	i86957	09593427-n
aidoneus%1:18:00::	09570298-n	i86957	i86957	09593427-n
hades%1:18:00::	09570298-n	i86957	i86957	09593427-n
pluto%1:18:00::	09570298-n	i86957	i86957	09593427-n
dis%1:18:00::	09570522-n	i86958	i86958	09593643-n
orcus%1:18:00::	09570522-n	i86958	i86958	09593643-n
dis_pater%1:18:00::	‡	‡	i86958	09593643-n

- Synset mappings can't handle splits
- Synsets are not split (Pluto doesn't move)
- All senses accurate, except one false positive (*fp*)

Results

Eric Kafe

Sense keys

Mapping

Results

Discussion

Summary

References

Table: The 15 biggest Wordnets in OMW 1.4, mapped to OEWN 2021, using synset offsets vs. CILI 1.0

Language	Synsets	Map _{Offset}			Map _{CILI}		
	PWN 3.0	OEWN 2021	Lost	%	OEWN 2021	Lost	%
<i>English</i>	117659	117454	205	0.17	117427	232	0.20
<i>Finnish</i>	116763	116562	201	0.17	116535	228	0.20
<i>Thai</i>	73350	73240	110	0.15	73223	127	0.17
<i>French</i>	59091	59015	76	0.13	59005	86	0.15
<i>Japanese</i>	57184	57086	98	0.17	57080	104	0.18
<i>Romanian</i>	56026	55941	85	0.15	55931	95	0.17
<i>Catalan</i>	45826	45773	53	0.12	45769	57	0.12
<i>Portuguese</i>	43895	43844	51	0.12	43840	55	0.13
<i>Slovenian</i>	42583	42520	63	0.15	42513	70	0.16
<i>Mandarin Chinese</i>	42300	42249	51	0.12	42240	60	0.14
<i>Spanish</i>	38512	38431	81	0.21	38418	94	0.24
<i>Indonesian</i>	38085	38018	67	0.18	38011	74	0.19
<i>Standard Malay</i>	36911	36843	68	0.18	36836	75	0.20
<i>Italian</i>	35001	34964	37	0.11	34960	41	0.12
<i>Polish</i>	33826	33798	28	0.08	33794	32	0.09
<i>Average</i>	55800.80	55715.87	84.93	0.15	55705.47	95.33	0.16

- Synset offsets always better than CILI 1.0
- 99.85% true positives in average

Performance between $PWN_{3.0}$ and $OEWN_{2021}$

Eric Kafe

Sense keys

Mapping

Results

Discussion

Summary

References

<i>Synsets</i>	Mapped	Not Mapped
True	$PWN_{3.0} \cap OEWN_{2021}$ $tp = 117454$	\emptyset $tn = 0$
False	<i>Splits</i> $fp = 44$	$\begin{matrix} \subset PWN_{3.0} \\ \subset OEWN_{2021} \end{matrix}$ $fn = 205$

Almost perfect performance results:

$$precision = \frac{tp}{tp + fp} = 0.9996 \quad (1)$$

$$recall = \frac{tp}{tp + fn} = 0.9983 \quad (2)$$

$$f1 = \frac{2 * precision * recall}{precision + recall} = 0.9989 \quad (3)$$

Eric Kafe

Sense keys

Mapping

Results

Discussion

Summary

References

Strengths :

- map in linear time
- almost perfect accuracy
- no problem with merges

Limitations :

- needs stable sense keys
- cannot handle splits

Opportunities :

- very few problems, easy manual review
- link multilingual wordnets directly to OEWN, instead of going through PWN 3.0

Changed Sense Keys

Eric Kafe

Sense keys

Mapping

Results

Discussion

Summary

References

Sense Key	PWN 3.0	OEWN 2021
sequoia%1:20:00::	<i>either of two huge coniferous California trees that reach a height of 300 feet; sometimes placed in the Taxodiaceae</i>	‡
sequoia%1:20:01::	‡	<i>either of two huge coniferous California trees that reach a height of 300 feet; sometimes placed in the Taxodiaceae</i>

More Changed Sense Keys

Eric Kafe

Sense keys

Mapping

Results

Discussion

Summary

References

Sense Key	PWN 3.0	OEWN 2021
stub_out%2:30:00::	<i>extinguish by crush- ing</i>	⚡
stub_out%2:35:01::	⚡	<i>extinguish by crush- ing</i>
obtrusive%3:00:00::	<i>undesirably notice- able</i>	⚡
obtrusive%5:00:00- :noticeable:00	⚡	<i>undesirably notice- able</i>
newfangled%5:00:00- :original:00	<i>(of a new kind or fashion) gratuitously new</i>	⚡
newfangled%5:00:00- :new:00	⚡	<i>(of a new kind or fashion) gratuitously new</i>

Summary

Eric Kafe

Sense keys

Mapping

Results

Discussion

Summary

References

- Wordnets need mapping, even with ILI
- Baseline algorithm extremely fast and accurate
- Only review a few splits and losses
- Perspective: link directly to OEWN

Eric Kafe

Sense keys

Mapping

Results

Discussion

Summary

References

- NLTK <https://www.nltk.org>
<https://github.com/nltk/nltk>
- Wn (Goodman and Bond, 2021)
<https://github.com/goodmami/wn>

Databases

- OEWN <https://github.com/globalwordnet/english-wordnet>
- *wndb* <https://github.com/x-englishwordnet/wndb>
- Collaborative Interlingual index (CILI)
<https://github.com/globalwordnet/cili>
- Sense Key Index (SKI) <https://github.com/ekaf/ski>

References I

Eric Kafe

Sense keys

Mapping

Results

Discussion

Summary

References

Francis Bond, Christiane Fellbaum, Shu-Kai Hsieh, Chu-Ren Huang, Adam Pease, and Piek Vossen. 2014. A multilingual lexico-semantic database and ontology. In *Towards the Multilingual Semantic Web*, pages 243–258. Springer.

Francis Bond, Luis Morgado da Costa, Michael Wayne Goodman, John Philip McCrae, and Ahti Lohk. 2020. Some issues with building a multilingual Wordnet. In *Proceedings of the 12th Language Resources and Evaluation Conference*, pages 3189–3197, Marseille, France. European Language Resources Association.

Francis Bond, Piek Vossen, John McCrae, and Christiane Fellbaum. 2016. CILI: the collaborative interlingual index. In *Proceedings of the 8th Global WordNet Conference (GWC)*, pages 50–57, Bucharest, Romania. Global Wordnet Association.

Eric Kafe

Sense keys

Mapping

Results

Discussion

Summary

References

Christiane Fellbaum. 1998. *WordNet, An Electronic Lexical Database*. MIT Press, Cambridge.

Michael Wayne Goodman and Francis Bond. 2021. Intrinsically interlingual: The wn python library for wordnets. In *Proceedings of the 11th Global Wordnet Conference*, pages 100–107, University of South Africa (UNISA). Global Wordnet Association.

Eric Kafe. 2018. Persistent semantic identity in wordnet. *Cognitive Studies | Études cognitives*, 18.

References III

Eric Kafe

Sense keys

Mapping

Results

Discussion

Summary

References

John Philip McCrae, Alexandre Rademaker, Ewa Rudnicka, and Francis Bond. 2020. English WordNet 2020: Improving and extending a WordNet for English using an open-source methodology. In *Proceedings of the LREC 2020 Workshop on Multimodal Wordnets (MMW2020)*, pages 14–19, Marseille, France. The European Language Resources Association (ELRA).

Piek Vossen, Francis Bond, and John McCrae. 2016. Toward a truly multilingual GlobalWordnet grid. In *Proceedings of the 8th Global WordNet Conference (GWC)*, pages 424–431, Bucharest, Romania. Global Wordnet Association.

WordNet-team. 2010. Wordnet 3.0 reference manual. In *WordNet Documentation*. Princeton University, <https://wordnet.princeton.edu/documentation>.