Translation-based Lexicalization Generation and Lexical Gap Detection: Application to Kinship Terms

Senyu Li, Bradley Hauer, Ning Shi, Grzegorz Kondrak

Alberta Machine Intelligence Institute, Dept of Computing Science University of Alberta, Edmonton, Canada



An Error Case: Google Translate

- 堂哥 "elder son of father's brother" => "cousin"
- 堂姐 "elder daughter of father's brother" => "cousin"
- Other powerful translators make similar errors. (DeepL, Baidu, etc.)

Detect language English Chinese (Simplified	\sim \leftarrow	→ Chinese (Simplified) English Spanish ∨	
我有一个堂哥,但是没有堂姐	×	I have a cousin, but no cousin. 🕁	
Wŏ yŏu yīgè táng gē, dànshì méiyŏu táng jiě.			
♀ ↓ 14 / 5,000	拼 -	()	>
		*Google Translate, February 15,	, 2024

Sample Output of ChatGPT



Given a word that means [father's younger brother] in Chinese is [叔叔], and a word that means [mother's brother] in Chinese is [舅舅]. Is there a word that means [elder brother] in [English]? If yes, give me that word. If no, say no.



ChatGPT

Yes, the word in English that means "elder brother" is "brother."

*GPT-3.5, February 15, 2024.

Outline

- **Problem:** How to identify concept lexicalizations and lexical gaps efficiently?
- Idea: If a concept is an exclusive disjunction of its hyponym concepts then all three concepts should have different lexicalizations.
- **Method:** Generate a candidate lexicalization for each concept by translating an unambiguous lexicalization into the target language in the context of the concept gloss. Then filter out incorrect translations using the above idea.
- **Results:** Empirical evaluations demonstrate that our approach yields higher accuracy than BabelNet and ChatGPT.

Concepts

- Concept: discrete word meaning
- Kinship concepts have clear definitions and hierarchical structure
 - Well-studied, good gold-standard dataset (Khishigsuren et al, 2022)



Lexicalizations and Lexical Gaps

- Lexicalization: a single word which can express (i.e. lexicalize) a concept.
- Lexical Gap: a concept that has no lexicalization in a given language.

Concepts	En	Es	Fr	Ja	Fa	Zh	PI
1	Sibling	Ø	fratrie	Ø	Ø	同胞	Ø
2	Ø	Ø	Ø	Ø	Ø	Ø	Ø
3	Brother	hermano	frère	Ø	برادر	兄弟	brat
4	Sister	hermana	sœur	Ø	خواهر	姐妹	siostra
5	Ø	Ø	Ø	Ø	Ø	Ø	Ø
6	Ø	Ø	Ø	兄さん	Ø	哥哥	Ø
7	Ø	Ø	Ø	姉ちゃん	Ø	姐姐	Ø
8	Ø	tato	Ø	おとうと	Ø	弟弟	Ø
9	Ø	Ø	Ø	いもうと	Ø	妹妹	Ø

Data from Using Linguistic Typology to Enrich Multilingual Lexicons: the Case of Lexical Gaps in Kinship (Khishigsuren et al, 2022)

Task definition: LexGen and LexGap

- LexGen: Lexicalization Generation
 - Input: language *L*, concept *s*
 - Output: word *w* in *L* s.t. *w* lexicalizes *s*,
 OR a special token GAP indicating that no such *w* exists
- LexGap: Lexical Gap Detection
 - Input: language *L*, concept *s*
 - Output: True if no word in *L* lexicalizes *s*,
 False otherwise.
- LexGen(L,s) = GAP <u>if and only if</u> LexGap(L,s) = True



Theoretical Basis

• **Proposition 1:** If a concept P is an exclusive disjunction of its hyponym concepts C1 and C2, expressing P and C1 with the same word w can result in a colloquial contradiction.



Theoretical Basis

- **Proposition 1:** If a concept P is an exclusive disjunction of its hyponym concepts C1 and C2, expressing P and C1 with the same word w can result in a colloquial contradiction.
- **Proposition 2:** If a concept P is an exclusive disjunction of its hyponym concepts C1 and C2, expressing C1 and C2 with the same word w can result in a colloquial contradiction.



Theoretical Basis

- **Proposition 1:** If a concept P is an exclusive disjunction of its hyponym concepts C1 and C2, expressing P and C1 with the same word w can result in a colloquial contradiction.
- **Proposition 2:** If a concept P is an exclusive disjunction of its hyponym concepts C1 and C2, expressing C1 and C2 with the same word w can result in a colloquial contradiction.
- **Corollary:** If a concept P is an exclusive disjunction of its hyponyms C1 and C2 then all their lexicalizations should be different.







Our Method

- Generate a candidate lexicalization *w* for each concept by translating a seed word.
- **Filter** using our four-step procedure:

1. Multi-word filter: If w is not a single word (e.g. "male cousin"), return GAP

- 2. Horizontal filter (Proposition 2): If w was also generated for a sibling node of s, return GAP
- 3. Back-translation filter: If back-translating w does not recover the seed word, return GAP

4. Vertical filter (Proposition 1): If *w* was also generated for a parent node of *s*, **and** another child of that parent node has already been tagged as a GAP, then **return GAP**

• If *w* makes it past the filters, return *w* for LexGen, **False** for LexGap

Experimental Setup

- Data: Database of Lexical Diversity in Kinship by Khishigsuren et al. (2022)
- **Translator:** Google Translate
- **Metrics:** Accuracy for LexGen, F1 score for LexGap
- **Comparison:** All-Gaps, BabelNet 5.1, and ChatGPT w/ GPT-3.5 Turbo
- Languages
 - Development languages: English, Mandarin, and Persian.
 - Test languages: Spanish, Russian, French, German, Polish, Arabic, Italian, Mongolian, Hungarian, and Hindi.

*GPT-3.5 Turbo and Google Translate were accessed on February 15, 2024

Results



Conclusion

- Novel translate-and-filter method for:
 - Generating lexicalizations
 - Detecting lexical gaps
- Grounded in linguistic theory, with clear definitions and and propositions
- Leverages translation and hypernym/hyponym relations
- Future work: Beyond kinship to other domains

github.com/UAlberta-NLP/KinshipAutoLex