

Parallel Corpus Examples for Language Learning Applications

Johannes Graën

Thursday 16th May, 2019



GÖTEBORGS UNIVERSITET



Universitat
Pompeu Fabra
Barcelona



FONDS NATIONAL SUISSE
SCHWEIZERISCHER NATIONALFONDS
FONDO NAZIONALE SVIZZERO
SWISS NATIONAL SCIENCE FOUNDATION

From parallel corpora to multilingual exercises

Making use of large text collections and crowdsourcing techniques for innovative autonomous language learning applications

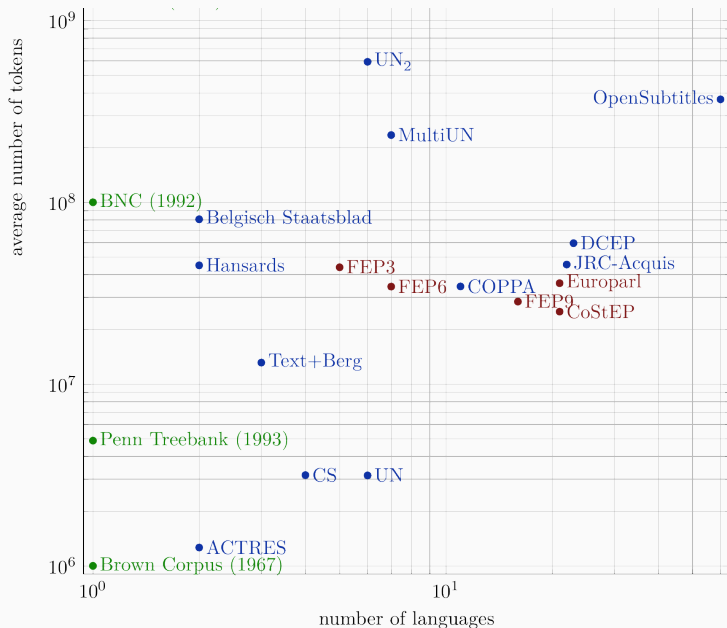
- Data-driven learning (DDL)
- Corpus-based learning is both effective and efficient (Cobb and Boulton 2015)
- Combination of NLP methods + CALL → ICALL applications

1. Parallel Corpora
2. Alignment
3. Applications Directed to Language Learners

Parallel Corpora

- Collection of translated texts (or speech)
- Wikipedia \neq parallel
- We have used Europarl for corpus linguistic studies

Parallel corpora



Comparing languages: article use

- **de**: "In unseren einzelnen Mitgliedstaat und gemeinsam **als Europäische Union** müssen wir [...]"
- **en**: "In our individual Member States, and collectively **as the European Union**, we must [...]"
- **es**: "En nuestros respectivos Estados miembros y, de manera colectiva, **en la Unión Europea** debemos [...]"
- **it**: "Sia nei singoli Stati membri che collettivamente, **come Unione europea**, dobbiamo esercitare [...]"
- **pt**: "Em cada um dos nossos Estados-Membros, e colectivamente **enquanto União Europeia**, temos que [...]"
- **sv**: "I våra enskilda medlemsstater, och samfällt **som Europeiska unionen**, måste vi [...]"

⇒ variable article use, zero articles (Callegaro 2017)

Translation direction

Wenn ihre Katze Bier trinkt , ist dies vielleicht der Grund , warum sie krank ist .

If her cat is drinking beer , then that is probably what is making the cat ill .

Si su gato bebe cerveza , probablemente sea eso lo que enferma al gato .

Jos hänen kissansa juo olutta , kissa tulee sairaaksi .

Si son chat boit de la bière , c' est probablement cela qui le rend malade .

Se il suo gatto beve birra , probabilmente è per quello che sta male .

Als haar kat bier drinkt , wordt hij daar waarschijnlijk ziek van .

Se o gato da senhora deputada bebe cerveja , provavelmente é isso que o traz doente .

Om hennes katt dricker öl så är det troligen det som gör katten sjuk .

Alignment

- **Correspondence relation** between units in parallel corpora (documents, sentences, words/tokens, ...)
- Also: process of finding those correspondences
- Initial purpose: training statistical machine translation models
- Pure statistics → no inherent linguistic properties

Multiparallel word alignment: linguistic perspective

Wir möchten nicht die Katze im Sack kaufen .

Nous ne voulons pas acheter chat en poche .

Não estamos interessados em comprar gato por lebre .

We are not interested in buying a pig in a poke .

Vi är inte intresserade av att köpa grisen i säcken .

Simple word alignment



Die rumänische Gesellschaft ist erwachsen geworden.



Romanian society has grown up.



La sociedad rumana ha madurado.



La société roumaine a mûri.



De Roemeense samenleving is volwassen geworden.

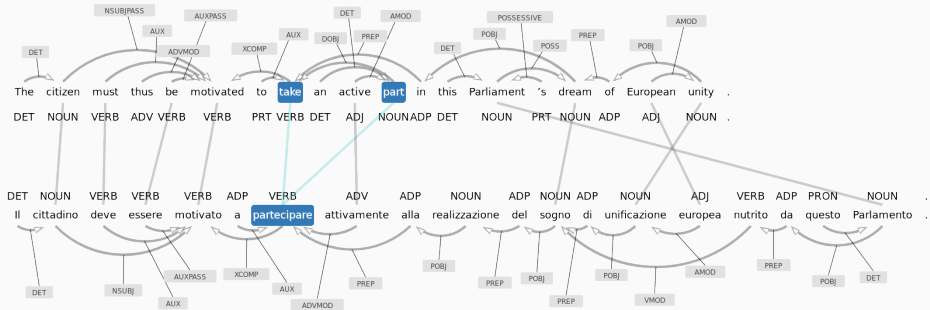


A sociedade romena amadureceu.



Det rumänska samhället har blivit vuxet.

Approach: alignment combined with annotation



- Language-wise annotation: part-of-speech tagging and parsing
- Alignment as interlingual (correspondence) annotation

Applications Directed to Language Learners

- CEFRLex project¹: lists of words & expressions extracted from graded textbook corpora for different languages (Dürlich and François 2018; François et al. 2016; Gala et al. 2013)
- Single words vs. multi-word expressions (e.g. pull + back)
- Lexical vs. grammatical proficiency

¹<http://cental.uclouvain.be/cefrlex/>

Good dictionary examples

- Following the idea of GDEX for lexica (Kilgarriff et al. 2008)
- Linguistic criteria (sentence length, vocabulary frequency, ...)
- Two dimensions: linguistic complexity and dependency on context

Sentence candidates for exercise generation

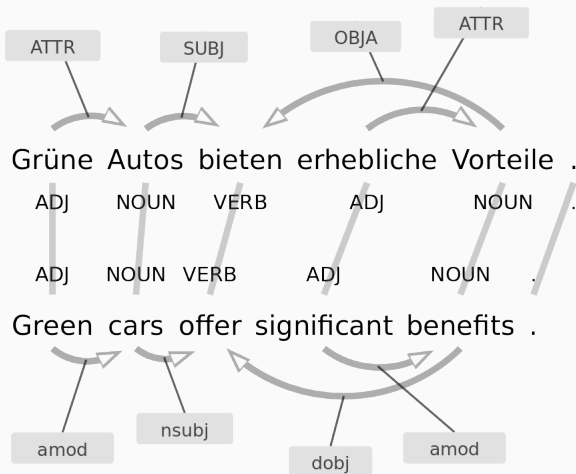
Nr	Criterion	Nr	Criterion
	Search term		Additional structural criteria
1	<i>Absence of search term</i>	13	Negative formulations
2	Number of matches	14	<i>Interrogative sentence</i>
3	<i>Position of search term</i>	15	<i>Direct speech</i>
	Well-formedness	16	<i>Answer to closed questions</i>
4	<i>Dependency root</i>	17	Modal verbs
5	Ellipsis	18	Sentence length
6	<i>Incompleteness</i>		Additional lexical criteria
7	Non-lemmatized tokens	19	Difficult vocabulary
8	Non-alphabetical tokens	20	Word frequency
	Context independence	21	Out-of-vocabulary words
9	<i>Structural connective in isolation</i>	22	Sensitive vocabulary
10	Pronominal anaphora	23	Typicality
11	Adverbial anaphora	24	Proper names
12	L2 complexity in CEFR level	25	Abbreviations

Criteria for HitEx framework (Pilán et al. 2016)

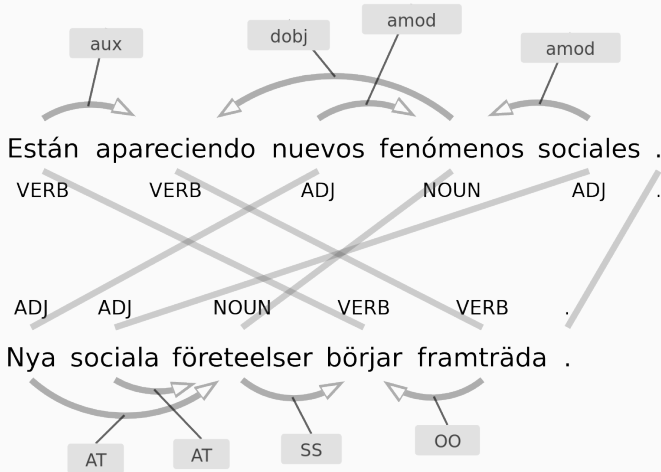
Candidates for parallel sentences

- Monolingual criteria plus translation characteristics
- Close (literal) or free (idiomatic) translation
- Correspondence on lexical level (single items vs. longer units)
- Similar part-of-speech and syntax or (systematic) difference between languages

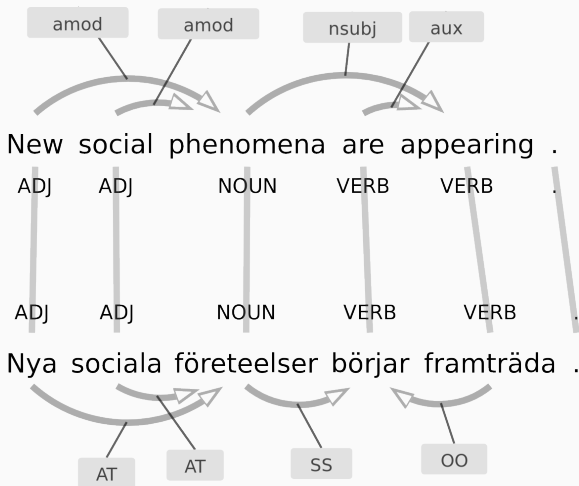
Straight correspondence \leftrightarrow literal translation



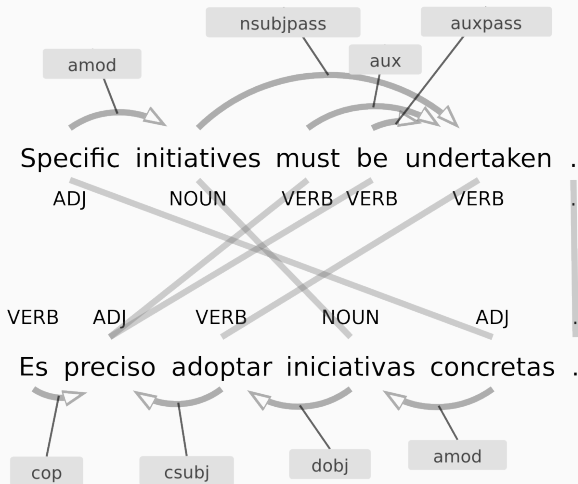
Diverging syntactical structure (es/sv)



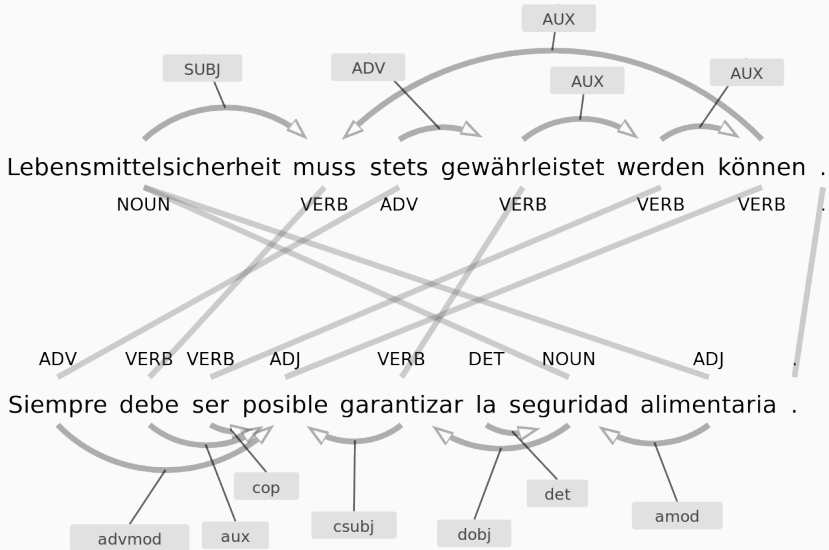
But similar structure in similar languages (en/sv)



Less literal translation ↔ more advanced correspondences



Complex structures and correspondences → high proficiency





Language Acquisition Reusing **Korp**

Bundle gap component

Show solution

Mellan juli och december förra året _____ bara 130 miljoner av de 700 miljoner som regeringen lagt undan .

Sen undrar jag om det inte var bola bola som _____ för köttbullar , ska försöka få tag i dottern hon vet .

" Alex , jag lyckades till slut spåra det där telefonnumret i Knightsbridge , det som _____ för restaurangens bordsbeställningar .

I veckor och månader efteråt _____ den här intervjun som sanningsvittne på att stortinget skulle få reducerad makt .

Continue



Language Acquisition Reusing **Korp**

Vocabulary Multiple Choice

B1

Change level

Click to generate!

4 Among other things he ordered a big _____ of the philosopher Nietzsche.

guinea pig



3 Rör om och _____ väl .

guinea pig

cannabis

bacon

bottle

portrait



2 Den upplevelsen kommer jag att _____ när jag är död .



1 Han pratade om att FN måste förändras för att _____ bättre .

kasta





Language Acquisition Reusing **Korp**

MultiPaVerX - Multilingual multiple choice particle verb exercise

Select language

I want to practice English ▼ particle verbs

Select mother tongue

My mother tongue is Spanish ▼

I also speak ...

Select proficiency in **English**:

Beginner Intermediate Advanced

Start



Language Acquisition Reusing **Korp**

MultiPaVerX - Multilingual multiple choice particle verb exercise

+1 word

50/50

more context

| Hints left: 2

Auseinanderbrechen

bryta _____

upp

ut

av

ihop

sönder

ned

MultiPaVerX - Multilingual multiple choice particle verb exercise

+1 word

50/50

more context

| Hints left: 0

Auseinanderbrechen

brechen

zerstören

bryta _____

upp

ut




av


ihop

sönder

ned

References

-  Callegaro, E. (2017). “Parallel Corpora for the Investigation of (Variable) Article Use in English: A Construction Grammar Approach”. PhD thesis. University of Zurich.
-  Cobb, T. and A. Boulton (2015). “Classroom applications of corpus analysis”. In: *The Cambridge Handbook of English Corpus Linguistics*. Ed. by D. Biber and R. Reppen. Cambridge University Press, pp. 478–497.
-  Dürlich, L. and T. François (2018). “EFLLex: A Graded Lexical Resource for Learners of English as a Foreign Language”. In: *11th International Conference on Language Resources and Evaluation (LREC)*.

-  François, T., E. Volodina, I. Pilán, and A. Tack (2016). “SVALex: a CEFR-graded Lexical Resource for Swedish Foreign and Second Language Learners.”. In: *10th International Conference on Language Resources and Evaluation (LREC)*.
-  Gala, N., T. François, and C. Fairon (2013). “Towards a French lexicon with difficulty measures: NLP helping to bridge the gap between traditional dictionaries and specialized lexicons”. In: *E-lexicography in the 21st century: thinking outside the paper*.
-  Kilgarriff, A., M. Husák, K. McAdam, M. Rundell, and P. Rychlý (2008). “GDEX: Automatically Finding Good Dictionary Examples in a Corpus”. In: *Proceedings of the 13th EURALEX International Congress*. Ed. by J. D. Elisenda Bernal. Institut Universitari de Lingüística Aplicada, Universitat Pompeu Fabra.



Pilán, I., E. Volodina, and L. Borin (2016). “Candidate sentence selection for language learning exercises: from a comprehensive framework to an empirical evaluation”. In: *Traitement Automatique des Langues* 57.3, pp. 67–91.