

Parallel corpus from patents in German, Spanish, French, Croatian, Norwegian, and Polish, paired with English



The EuroPat Corpus: A Parallel Corpus of European Patent Data

 Kenneth Heafield, **Elaine Farrow**, Jelmer van der Linde, Gema Ramírez-Sánchez, Dion Wiggins

INTRO

- Patents are a rich source of technical vocabulary and product names, complementing other data sources
- We mined parallel corpora by aggregating and aligning patents in six European languages paired with English

CORPUS SIZE

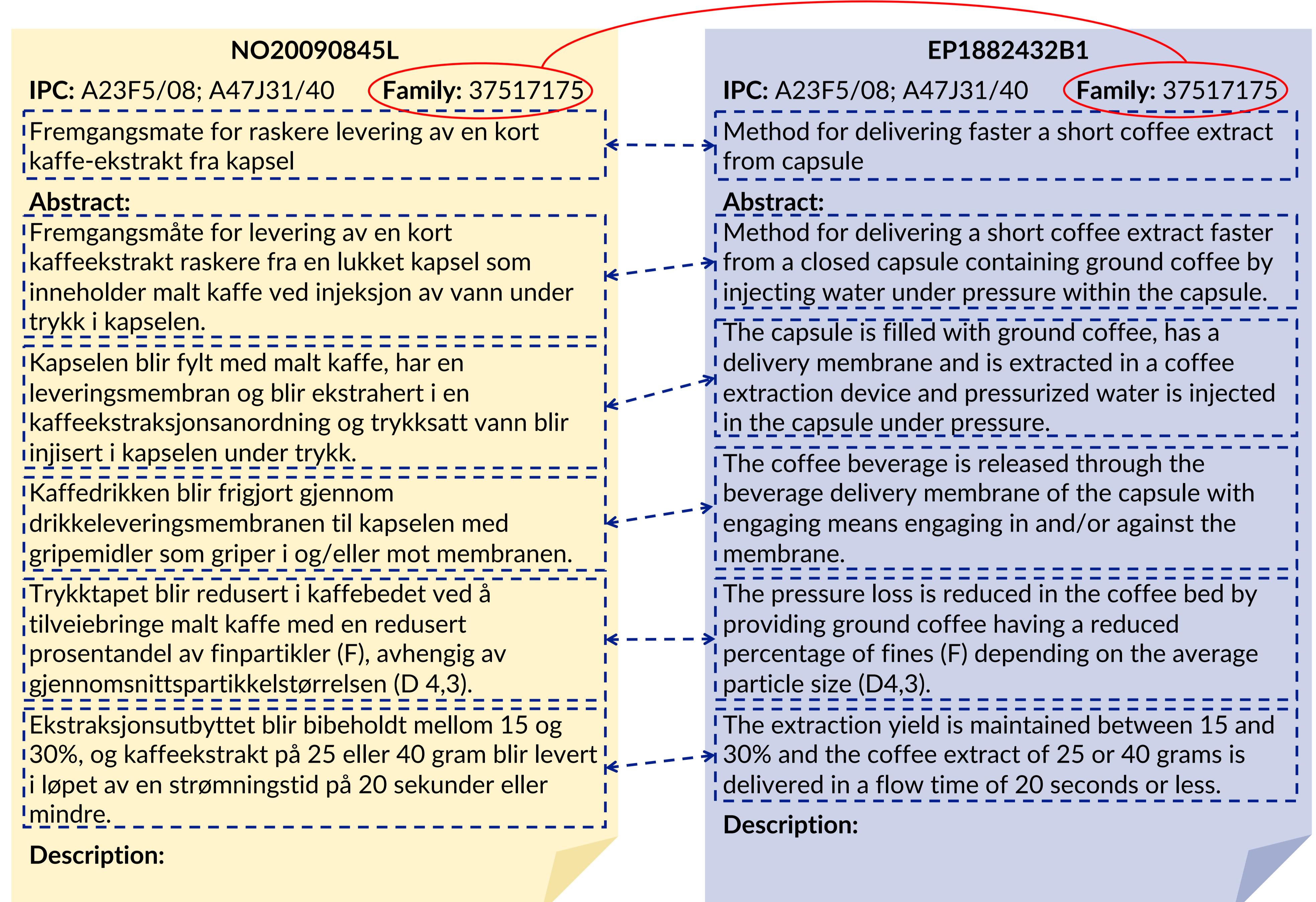
Language pair	Sentences (millions)
en-de	19.734
en-es	51.352
en-fr	11.098
en-hr	0.154
en-no	4.341
en-pl	0.332

CONCLUSION

- Creative Commons Zero, available at europat.net, OPUS, and ELRC-SHARE
- The EuroPat corpus is expected to be widely useful for training high-quality machine translation systems, particularly for those targeting technical documents such as patents and contracts

METHODS

1. Data acquisition: machine-readable text, scanned patent images
2. OCR and text extraction
3. Document alignment
4. Sentence alignment
5. Data cleaning
6. Domain grouping: International Patent Classification (IPC)
7. Corpus quality assessment: human and automatic
8. Data releases: RAW, TXT, and the translation memory standard TMX format



THE UNIVERSITY of EDINBURGH

prompsit

Omniscien™
TECHNOLOGIES



Co-financed by the European Union
Connecting Europe Facility

<https://europat.net>