

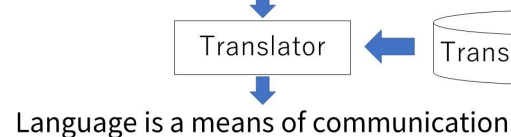
Abstract

Machine translation systems require many bilingual sentence pairs (translations of each other) as training data. We are researching **technology to build a parallel corpus (bilingual database)** by collecting bilingual text scattered on the Internet (Web) and in patent application archives. JParaCrawl, a web-based parallel corpus, was constructed by efficiently collecting many bilingual sentence pairs from the Web **using crowdsourcing**. JaParaPat, a patent parallel corpus, has improved the quality of its sentence pairs **by alternately extracting data and training models**. Both are the world's largest publicly available parallel corpus between Japanese and English. We will further enhance our technology to automatically build a high-quality parallel corpus in specific fields, such as medicine and finance, which are rich in specialized terminology, and in particular language pairs, such as Chinese and Japanese, **to implement a machine translation system customized to the needs of our customers**.

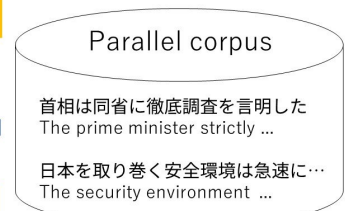
Machine translation and parallel data

The quantity and quality of parallel data determines machine translation accuracy

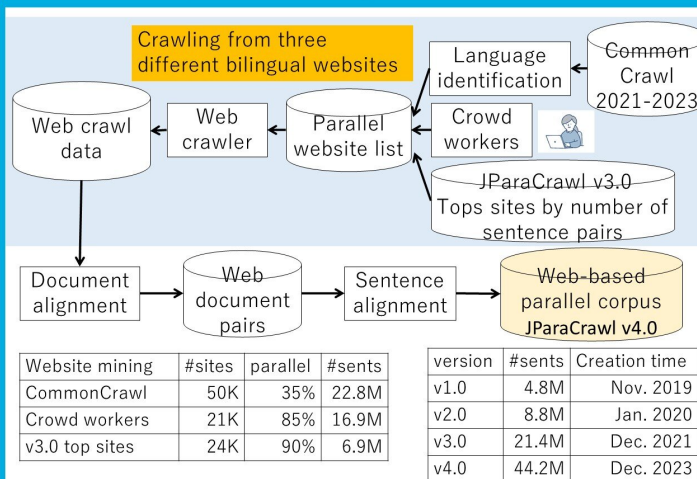
言語はコミュニケーションの手段である



Building a large and high-quality parallel corpus from the Internet (Web) and patents



Web-based parallel corpus: JParaCrawl v4.0

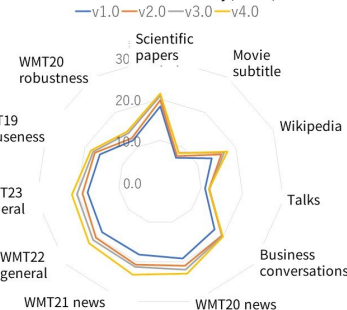


[JParaCrawl v4.0]

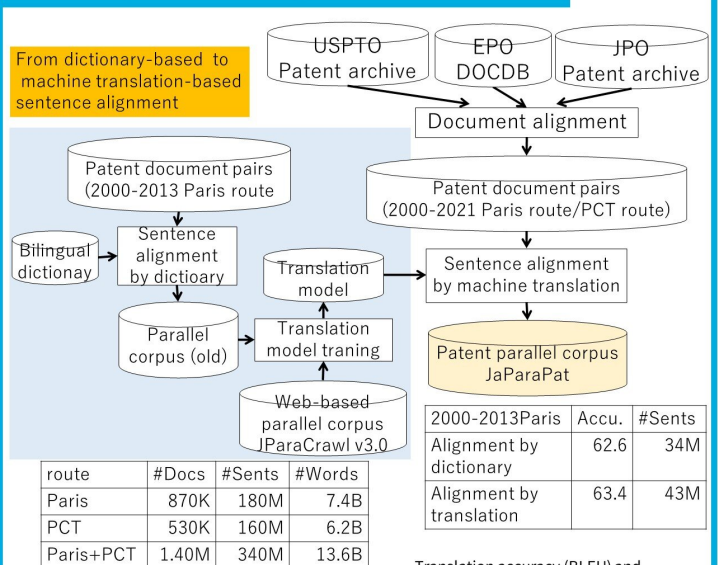
We built a web-based parallel corpus of more than 44M sentence pairs using three different bilingual website lists.

- Crowdsourced websites are more likely to be parallel, although each has fewer parallel sentences
- Improved translation accuracy in a variety of areas, including news, conversation, and scientific papers

JE Translation accuracy (BLEU)



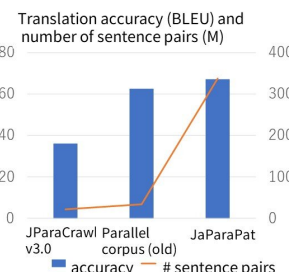
Patent parallel corpus: JaParaPat



[JaParaPat]

Based on patent families, we built a patent parallel corpus of more than 300M sentence pairs.

- Covers both the Paris route (priority claims) and the PCT route (international applications)
- Using machine translation for sentence alignment improves the quality and quantity



References

- [1] Makoto Morishita, Katsuki Chousa, Masaaki Nagata, "JParaCrawl v4.0: Building a Large Parallel Corpus with Crowdsourcing," *proceedings of the 30th annual meeting of Association for Natural Language Processing (NLP-2024)*, pp. 2330-2335, 2024 (in Japanese).
- [2] Masaaki Nagata, Makoto Morishita, Katsuki Chousa, Norihito Yasuda, "JaParaPat: A Large-scale Japanese-English Parallel Patent Application Corpus," *proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pp. 9452-9462, 2024.

Contact

Masaaki Nagata, Linguistic Intelligence Research Group, Innovative Communication Laboratory