



Concurrent multilingual Neural Machine Translation of several EU languages

David Dowey ddowey@stanford.edu

Stanford University

How can we build well-aligned high-quality parallel corpora for NMT in several languages from noisy multilingual websites?

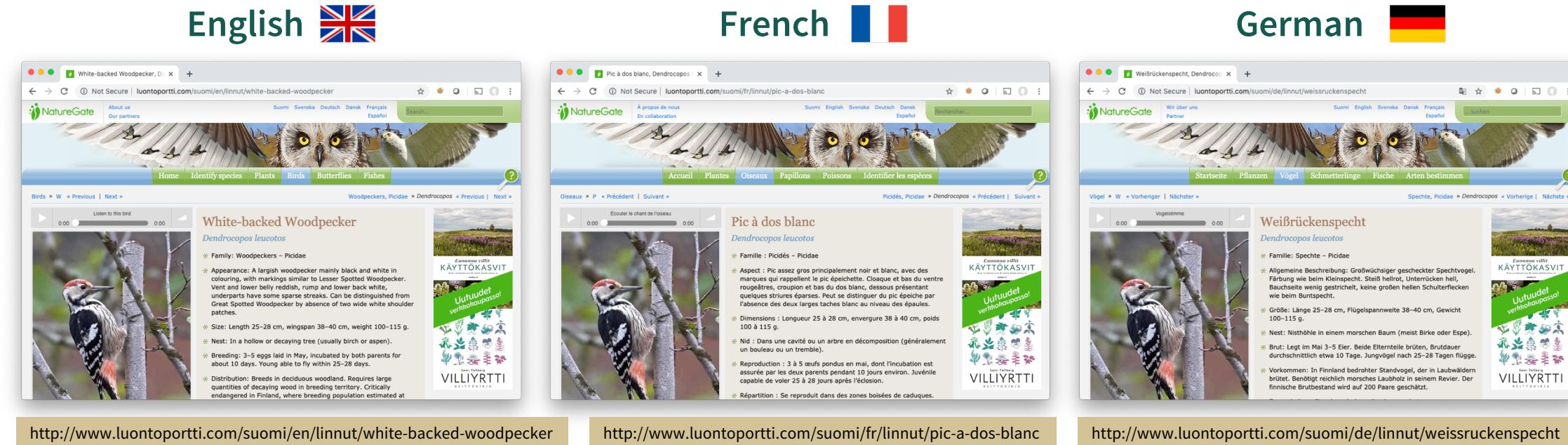
The problem addressed in this project is the messy but useful NLP task of aligning multilingual documents crawled from noisy website sources, for the purposes of generating high-quality parallel corpora of aligned documents across several languages. Automated NMT engines need parallel corpora that are both deep and broad, for example to improve translations tasks in specialized domains, such as for animal and plant species in the example.

The dataset

- The source data are web-crawled pages along with a language code (EN,FR,DE), the Mime type (text/html), the encoding charset (utf-8), URL, and the HTML and text full page content in Base64 encoding. The examples are scraped, processed and saved as lett files using the the bitextor scraper. The data come from the 2016 WMT shared task on document alignment, extended to multiple languages.
- The input data for the system are extracted from the source data – keeping the text of the webpage URLs and the language code. The URL may be similar across languages, or as in the example, quite different.
- The system outputs the best 1:1:1 match of URLs in EN,FR,DE as a triplet in that order. The URLs can then be used to extract the full page content from the source data to construct the parallel corpora.
- The testing triplets are hand chosen true matches (the oracle) used to measure the accuracy of the system

Domain	English pages	French pages	German pages	Testing triplets
luonoportti.com	3,645	1,796	2,146	70
vinci.com	3,564	3,374	3,005	50
dakar.com	17,420	14,582	1,655	50
krn.org	115	115	115	50
prohelvetia.ch	5,209	4,421	6,289	50
pawpeds.com	1,011	136	86	50
inst.at	3,203	543	17,076	50
schakportalen.nu	33	29	29	23
rehabcenter.lu	201	317	161	50
cyberspaceministry.org	1,534	958	173	40
galacticchannelings.com	4,231	1,283	983	50
	40,166	27,554	31,718	533

The task is not to align the from page content, but rather from the URLs - to use the techniques learned in class to determine which URLs from a scraped multilingual site lead to equivalent content in the different languages. In a sense, we seek to “translate” the URLs, but this is more difficult than it seems, since the URLs contain HTML code, redundant snippets and lots of non-letter characters that don’t translate.



An example is from the Finnish multilingual nature website luontoportti.com, (Nature Gate). Site used to identify species of plants, bird, butterflies and fish. Content is in English, French and German, but also Finnish, Swedish, Danish and Spanish. The site contains lots of media, web elements and links, but the content of interest is the center frame with information on species which is equivalent across multiple languages, which we can use to create domain-specific parallel corpora.

Processing pipeline

