

CUni Multilingual Matrix in the WMT 2013 Shared Task

Abstract

We describe our experiments with phrase-based machine translation for the WMT 2013 Shared Task. We trained one system for 18 translation directions between English or Czech on one side and English, Czech, German, Spanish, French or Russian on the other side. We describe a set of results with different training data sizes and subsets.

1 Introduction

With so many official languages, Europe is a paradise for machine translation research. One of the largest bodies of electronically available parallel texts is being nowadays generated by the European Union and its institutions. At the same time, the EU also provides motivation and boosts potential market for machine translation outcomes.

Most of the major European languages belong to one of three branches of the Indo-European language family: Germanic, Romance or Slavic. Such relatedness is responsible for many structural similarities in European languages, although significant differences still exist. Within the language portfolio selected for the WMT shared task, English, French and Spanish seem to be closer to each other than to the rest.

German, despite being genetically related to English, differs in many properties. Its word order rules, shifting verbs from one end of the sentence to the other, easily create long-distance dependencies. Long German compound words are notorious for increasing out-of-vocabulary rate, which has

led many researchers to devising unsupervised compound-splitting techniques. Also, uppercase/lowercase distinction is more important because all German nouns start with an uppercase letter by the rule.

Czech is a language with rich morphology (both inflectional and derivational) and relatively free word order. In fact, the predicate-argument structure, often encoded by fixed word order in English, is usually captured by inflection (especially the system of 7 grammatical cases) in Czech. While the free word order of Czech is a problem when translating to English (the text should be parsed first in order to determine the syntactic functions and the English word order), generating correct inflectional affixes is indeed a challenge for English-to-Czech systems. Furthermore, the multitude of possible Czech word forms (at least order of magnitude higher than in English) makes the data sparseness problem really severe, hindering both directions.

Most of the above characteristics of Czech also apply to Russian, another Slavic language. Similar issues have to be expected when translating between Russian and English. Still, there are also interesting divergences between Russian and Czech, especially on the syntactic level. Russian sentences typically omit copula and there is also no direct equivalent of the verb “to have”. Periphrastic constructions such as “there is XXX by him” are used instead. These differences make the Czech-Russian translation interesting as well.

Our goal is to run one system under as similar conditions as possible to all fourteen translation directions, to compare their translation accuracies and see why some directions are

easier than others. Future work will benefit from knowing what are the special processing needs for a given language pair. The current version of the system does not include really language-specific techniques: we neither split German compounds, nor do we address the peculiarities of Czech and Russian mentioned above.

2 The Translation System

Our translation system is built around Moses¹ (Koehn et al., 2007). Two-way word alignment was computed using GIZA++² (Och and Ney, 2003), and alignment symmetrization using the *grow-diag-final-and* heuristic (Koehn et al., 2003). Weights of the system were optimized using MERT (Och, 2003). No lexical reordering model was trained.

For language modeling we use the SRILM toolkit³ (Stolcke, 2002) with modified Kneser-Ney smoothing (Kneser and Ney, 1995; Chen and Goodman, 1998).

3 Data and Pre-processing Pipeline

We applied our system to all the ten official language pairs. In addition, we also experimented with translation between Czech on one side and German, Spanish, French or Russian on the other side. Training data for these additional language pairs were obtained by combining parallel corpora of the officially supported pairs. For instance, to create the Czech-German parallel corpus, we identified the intersection of the English sides of Czech-English and English-German corpora, respectively; then we combined the corresponding Czech and German sentences.

We took part in the constrained task. Unless explicitly stated otherwise, the translation model in our experiments was trained on the combined News-Commentary v8 and Europarl v7 corpora.⁴ Note that there is only News Commentary and no Europarl for Russian. We were also able to evaluate several combinations with large parallel corpora: the UN corpus (English, French and Spanish), the Giga French-English corpus and CzEng

(Czech-English). We did not use any large corpus for Russian-English. Tables 1 and 2 show the sizes of the training data.

Corpus	SentPairs	Tkns lng1	Tkns lng2
cs-en	786,929	18,196,080	21,184,881
de-en	2,098,430	55,791,641	58,403,756
es-en	2,140,175	62,444,507	59,811,355
fr-en	2,164,891	70,363,304	60,583,967
ru-en	150,217	3,889,215	4,100,148
de-cs	657,539	18,160,857	17,788,600
es-cs	697,898	19,577,329	18,926,839
fr-cs	693,093	19,717,885	18,849,244
ru-cs	103,931	2,642,772	2,319,611

Table 1: Number of sentence pairs and tokens for every language pair in the parallel training corpus. Languages are identified by their ISO 639 codes: cs = Czech, de = German, en = English, es = Spanish, fr = French, ru = Russian. Every line corresponds to the respective version of EuroParl + News Commentary.

Czeng	SentPairs	Tkns lng1	Tkns lng2
cs-en	14,833,358	204,837,216	235,177,231
UN			
es-en	11,196,913	368,154,702	328,840,003
fr-en	12,886,831	449,279,647	372,627,886
Giga			
fr-en	22,520,400	854,353,231	694,394,577

Table 2: Sizes of additional large parallel corpora.

The News Test 2010 (2489 sentences in each language) and 2012 (3003 sentences) data sets⁵ were used as development data for MERT. BLEU scores reported in this paper were computed on the News Test 2013 set (3000 sentences each language). We do not use the News Tests 2008, 2009 and 2011.

All parallel and monolingual corpora underwent the same preprocessing. They were tokenized and some characters normalized or cleaned. A set of language-dependent heuristics was applied in an attempt to restore and normalize the directed (opening/closing) quotation marks (i.e. "quoted" → “quoted”). The motivation is twofold here: First, we hope

¹<http://www.statmt.org/moses/>

²<http://code.google.com/p/giza-pp/>

³<http://www-speech.sri.com/projects/srilm/>

⁴<http://www.statmt.org/wmt13/translation-task.html#download>

⁵<http://www.statmt.org/wmt13/translation-task.html>

that paired quotation marks could occasionally work as brackets and better denote parallel phrases for Moses; second, if Moses learns to output directed quotation marks, subsequent detokenization will be easier.

The data are then tagged and lemmatized. We used the Featurama tagger for Czech and English lemmatization and TreeTagger for German, Spanish, French and Russian lemmatization. All these tools are embedded in the Treex analysis framework (Žabokrtský et al., 2008).

The lemmas are used later to compute word alignment. Besides, they are needed to apply “supervised truecasing” to the data: we cast the case of the lemma to the form, relying on our morphological analyzers and taggers to identify proper names, all other words are lowercased. Note that guessing of the true case is only needed for the sentence-initial token. Other words can typically be left in their original form, unless they are uppercased as a form of HIGHLIGHTING.

3.1 Quotation Marks

A broad range of characters is used to represent quotation marks in the training data: straight ASCII quotation mark; Unicode directed quotation marks (U+2018 to U+201F); acute and grave accents; math symbols such as prime and double prime (U+2032 to U+2037) etc. Spaces around quotes in the original untokenized text ought to provide hints as to the direction of the quotes (no space between the opening quote and the next word, and no space between the closing quote and the previous word) but unfortunately there are numerous cases where superfluous spaces are inserted or required spaces are missing.

Nested quoting is also possible, such as in

As the Wise Men ’ s Report also says , and I quote : ’ It is elementary ’ common sense ’ that the Commission should have supported the Parliament ’ s decision - making process .

We want all possible quotation marks converted to one pair of characters. We do not mind the distinction between single and double quotes but we want to keep (or restore) the distinction between opening and closing quotes. In addition, we want to identify the apostrophe acting as grapheme in some lan-

guages, and keep it (or normalize it, as it could also be mis-typed as acute accent or something else):

As the Wise Men ’ s Report also says , and I quote : “ It is elementary “ common sense ” that the Commission should have supported the Parliament ’ s decision - making process .

We attempt at solving the problem by a set of rules that consider mutual positions of quotation marks, spaces and other punctuation, and also some language-dependent rules (especially on the lexical apostrophe, e.g. in French *d’*, *l’*).

Our rules applied to 1.84 % of Spanish sentences, 2.47 % Czech, 2.77 % German, 4.33 % English and 16.9 % French (measured on Europarl data).

Our approach is different from the normalization script provided and applied by the organizers of the shared task, which merely converts all quotes to the undirected ASCII characters. We believe that such MT output is incorrect, so we submitted two versions of each system run: the *primary* version is intended for human evaluation and does not apply the “official” normalization of punctuation. In contrast, the *secondary* version is normalized, which naturally leads to higher scores in the automatic evaluation.

4 Experiments

In the following section we describe several different settings and corpora combinations we experimented with. BLEU scores have been computed by our system, comparing truecased tokenized hypothesis with truecased tokenized reference translation.

Such scores must differ from the official evaluation—see Section 4.4 for discussion of the final results.

The confidence interval for most of the scores lies between ± 0.5 and ± 0.6 BLEU % points.

4.1 Baseline Experiments

The set of baseline experiments were trained on the supervised truecased combination of News Commentary and Europarl. As we had lemmatizers for the languages, word alignment was computed on lemmas. (But our previous

experiments showed that there was little difference between using lemmas and lowercased 4-character “stems”). A hexagram language model was trained on the monolingual version of the News Commentary + Europarl corpus (typically a slightly larger superset of the target side of the parallel corpus).

4.2 Larger Monolingual Data

Besides the monolingual halves of the parallel corpora, additional monolingual data were provided / permitted:

- The Crawled News corpus from the years 2007 to 2012, various sizes for each language and year.
- The Gigaword corpora published by the Linguistic Data Consortium, available only for English (5th edition), Spanish (3rd) and French (3rd).

Our experiments in previous years clearly showed that the Crawled News corpus, in-domain and large, contributed significantly to better BLEU scores. This year we included it in our baseline experiments for all language pairs: translation model on News Commentary + Europarl, language model on monolingual part of the two, plus Crawled News (`news.all`).

Table 3 gives the sizes of the subsets available for our experiments and Table 4 compares BLEU scores with Gigaword against the baseline. Gigaword mainly contains texts from news agencies and as such should be also in-domain. Nevertheless, the crawled news are already so large that the improvement contributed by Gigaword is rarely significant.

4.3 Larger Parallel Data

Various combinations with larger parallel corpora were also tested. We do not have results for all combinations because these experiments needed a lot of time and resources and not all of them finished in time successfully.

In general the UN corpus seems to be of low quality or too much off-domain. It may help a little if used in combination with news-euro. If used separately, it always hurts the results.

The Giga French-English corpus gave the best results for English-French as expected,

Corpus	Segments	Tokens
newsc+euro.cs	830,904	18,862,626
newsc+euro.de	2,380,813	59,350,113
newsc+euro.en	2,466,167	67,033,745
newsc+euro.es	2,330,369	66,928,157
newsc+euro.fr	2,384,293	74,962,162
newsc.ru	183,083	4,340,275
news.all.cs	27,540,827	460,356,173
news.all.de	54,619,789	1,020,852,354
news.all.en	68,341,615	1,673,187,787
news.all.es	13,384,314	388,614,890
news.all.fr	21,195,476	557,431,929
news.all.ru	19,912,911	361,026,791
gigaword.en	117,905,755	4,418,360,239
gigaword.es	31,304,148	1,064,660,498
gigaword.fr	21,674,453	963,571,174

Table 3: Number of segments (paragraphs in Gigaword, sentences elsewhere) and tokens of additional monolingual training corpora. “newsc+euro” are the monolingual versions of the News Commentary and Europarl parallel corpora. “news.all” denotes all years of the Crawled News corpus for the given language.

even without the core news-euro data. However, training the model on data of this size is extremely demanding on memory and time.

Finally, Czeg undoubtedly improves Czech-English translation in both directions. The news-euro dataset is smaller for this language pair, which makes Czeg stand out even more. See Table 5 for details.

4.4 Final Results

Table 6 compares our BLEU scores with those computed at `matrix.statmt.org`.

BLEU (without flag) denotes BLEU score computed by our system, comparing truecased tokenized hypothesis with truecased tokenized reference translation.

The official evaluation by `matrix.statmt.org` gives typically lower numbers, reflecting the loss caused by detokenization and new (different) tokenization.

4.5 Efficiency

The baseline experiments were conducted mostly on 64bit AMD Opteron quad-core 2.8 GHz CPUs with 32 GB RAM (decoding run on 15 machines in parallel) and the whole

Direction	Baseline	Gigaword
en-cs	0.1632	
en-de	0.1833	
en-es	0.2808	0.2856
en-fr	0.2987	0.2988
en-ru	0.1582	
cs-en	0.2328	0.2367
de-en	0.2389	0.2436
es-en	0.2916	0.2975
fr-en	0.2887	
ru-en	0.1975	0.2003
cs-de	0.1595	
cs-es	0.2170	0.2220
cs-fr	0.2220	0.2196
cs-ru	0.1660	
de-cs	0.1488	
es-cs	0.1580	
fr-cs	0.1420	
ru-cs	0.1506	

Table 4: BLEU scores of the baseline experiments (left column) on News Test 2013 data, computed by the system on tokenized data, versus similar setup with Gigaword. The improvement, if any, was typically not significant.

pipeline typically required between a half and a whole day.

However, we used machines with up to 500 GB RAM to train the large language models and translation models. Aligning the UN corpora with Giza++ took around 5 days. Giga French-English corpus was even worse and required several weeks to complete. Using such a large corpus without pruning is not practical.

5 Conclusion

We have described the Moses-based SMT system we used for the WMT 2013 shared task. We discussed experiments with large data for many language pairs from the point of view of both the translation accuracy and efficiency.

Acknowledgements

The work on this project was supported by the grant P406/11/1499 of the Czech Science Foundation (GAČR).

Dir	Parallel	Mono	BLEU
en-es	news-euro	+gigaword	0.2856
en-es	news-euro-un	+gigaword	0.2844
en-es	un	un+gigaw.	0.2016
en-fr	giga	+gigaword	0.3106
en-fr	giga	+newsall	0.3037
en-fr	news-euro-un	+gigaword	0.3010
en-fr	news-euro	+gigaword	0.2988
en-fr	un	un	0.2933
es-en	news-euro	+gigaword	0.2975
es-en	news-euro-un	baseline	0.2845
es-en	un	un+news	0.2067
fr-en	news-euro-un	+gigaword	0.2914
fr-en	news-euro	baseline	0.2887
fr-en	un	un+news	0.2737

Table 5: BLEU scores with different parallel corpora.

References

- Stanley F. Chen and Joshua Goodman. 1998. An empirical study of smoothing techniques for language modeling. In *Technical report TR-10-98, Computer Science Group*, Harvard, MA, USA, August. Harvard University.
- Reinhard Kneser and Hermann Ney. 1995. Improved backing-off for m-gram language modeling. In *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 181–184, Los Alamitos, California, USA. IEEE Computer Society Press.
- Philipp Koehn, Franz Josef Och, and Daniel Marcu. 2003. Statistical phrase-based translation. In *NAACL '03: Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology*, pages 48–54, Morristown, NJ, USA. Association for Computational Linguistics.
- Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondřej Bojar, Alexandra Constantin, and Evan Herbst. 2007. Moses: Open Source Toolkit for Statistical Machine Translation. In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics Companion Volume Proceedings of the Demo and Poster Sessions*, pages 177–180, Praha, Czechia, June. Association for Computational Linguistics.
- Franz Josef Och and Hermann Ney. 2003. A systematic comparison of various statistical alignment models. *Computational Linguistics*, 29(1):19–51.

Direction	<i>BLEU</i>	<i>BLEU_l</i>	<i>BLEU_t</i>
en-cs	0.1786	0.180	0.170
en-de	0.1833	0.179	0.173
en-es	0.2856	0.288	0.271
en-fr	0.3010	0.270	0.259
en-ru	0.1582	0.142	0.142
cs-en	0.2527	0.259	0.244
de-en	0.2389	0.244	0.230
es-en	0.2856	0.288	0.271
fr-en	0.2887	0.294	0.280
ru-en	0.1975	0.203	0.191
cs-de	0.1595	0.159	0.151
cs-es	0.2220	0.225	0.210
cs-fr	0.2220	0.191	0.181
cs-ru	0.1660	0.150	0.149
de-cs	0.1488	0.151	0.142
es-cs	0.1580	0.160	0.152
fr-cs	0.1420	0.145	0.137
ru-cs	0.1506	0.151	0.144

Table 6: BLEU scores with the large language models. *BLEU* is truecased computed by the system, *BLEU_l* is the official lowercased evaluation by `matrix.statmt.org`. *BLEU_t* is official truecased evaluation. Although lower official scores are expected, notice the larger gap in en-fr and cs-fr translation. There seems to be a problem in our French detokenization procedure.

Franz Josef Och. 2003. Minimum error rate training in statistical machine translation. In *ACL '03: Proceedings of the 41st Annual Meeting on Association for Computational Linguistics*, pages 160–167, Morristown, NJ, USA. Association for Computational Linguistics.

Andreas Stolcke. 2002. Srilm – an extensible language modeling toolkit. In *Proceedings of International Conference on Spoken Language Processing*, Denver, Colorado, USA.

Zdeněk Žabokrtský, Jan Ptáček, and Petr Pajas. 2008. TectoMT: Highly modular MT system with tectogrammatics used as transfer layer. In *ACL 2008 WMT: Proceedings of the Third Workshop on Statistical Machine Translation*, pages 167–170, Columbus, OH, USA. Association for Computational Linguistics.