

English–Hindi Translation in 21 Days

Ondřej Bojar, Pavel Straňák, Daniel Zeman*

Univerzita Karlova v Praze, ÚFAL,
Malostranské nám. 25, 118 00, Praha, Czech Republic
{bojar, stranak, zeman}@ufal.mff.cuni.cz

Abstract

We present experiments with a Moses-based English-to-Hindi translation system. We evaluate the impact of additional out-of-domain training data, both parallel and Hindi-only, and experiment with three methods for improving word order: standard Moses reordering model, rule-based pre-processing and language-independent suffix identification.

1 Data

1.1 Parallel

We tried to obtain as much parallel data as possible:

EILMT+TIDES: The parallel data provided by the organizers; 7k and 50k sentences.¹

Emille: The parallel part of this corpus consists of 200k words of text in English and its accompanying translations in Hindi and other languages.

Agriculture domain parallel corpus:

Resource Center for Indian Language

⁰This work was supported by grants of Czech Academy of Sciences: 1ET201120505, 1ET101470416, European Union: FP6-IST-5-034291-STP, and Ministry of Education: MSM0021620838.

¹We acknowledge the Department of IT (DIT), Government of India and the English-to-Indian Languages (EILMT) consortium for making the EILMT tourism dataset available.

Technology Solutions English-Hindi-Marathi-UNL parallel corpus. It contains 17,105 English and 13,248 Hindi words. Daniel Pipes: D. Pipes' website:² a limited-domain articles about the Middle East. 322 texts are available in Hindi.

1.2 Hindi

For Hindi language model, we generally use the target side of the parallel data mentioned above. To extend the data, we could use e.g. the monolingual sections of the Emille corpus.

These data are still quite small compared to data regularly used for e.g. English language models, so we decided to create a big monolingual corpus of Hindi ourselves. Starting from Hindi news sites lists we downloaded 27 websites, mostly news portals. After clean-up this amounted for 18.1M sentences and 309M tokens.

Downloading the data itself was done very simply.³ We then proceeded to clean-up the HTML and classify the texts by languages. The language classification is based on a model comparing frequencies of three-character suffixes of word forms with known suffix frequencies for each language. If no clearly-winning language is found, we back off to three simple models counting 1-, 2-, and 3-grams of characters.

²<http://www.danielpipes.org/>

³`wget -r -user-agent="" -timestamping -no-parent "$1"`

Finally we filtered the texts classified as Hindi by deleting all the lines that did not contain any Devanagari.

1.3 Tokenisation

For the data supplied by the organizers, we stick to their sentence segmentation and tokenization. For the additional data, we use a trainable tokenizer by (Klyueva and Bojar, 2008) that can be easily adapted to a new language simply by providing a few instances of sentence and token breaks.

2 Translation Setup

2.1 Baseline and Evaluation

We use a variation of Moses (Koehn et al., 2007) standard pipeline. Word alignments are obtained using GIZA++ (grow-diag-final heuristic). To reduce data sparseness for word-alignment estimation, we consider only first four letters (lowercased) of each English and Hindi word form.

The baseline setup is single-factored (i.e. plain phrase-based) wordform-to-wordform translation with distance-based reordering model. As a basic contrast experiment, we use the standard orientation-bidirectional reordering model based on both source and target side tokens.

We evaluate the translation quality using our implementation of BLEU with 95% confidence intervals obtained using the bootstrapping method by (Koehn, 2004). Due to the built-in tokenisation rules in `mteval-v11b` (and not in our system), the scores differ slightly.

We use Moses standard MERT (Och, 2003) to tune weights of the individual models in our setup. The final scores reported are evaluated on the unseen official Test set. Unless stated otherwise, we use only the EILMT Development and Test sets.⁴

⁴500 sentences each, 1 reference translation per sen-

Our baseline results as well as the impact of additional parallel data are given in Table 1. Unfortunately, we were not able to replicate the baseline scores reported by the organizers⁵, although we use the same maximum phrase-length (7 tokens), hopefully the same data and the same `mteval-v11b`.

2.2 Web LM

Table 2 summarizes the impact of using a 4-gram language model (LM) based on 309M Hindi tokens (see Section 1.2), called `webLM`, in addition to the 3-gram model based on the Hindi side of EILMT. We use a separate weight for `webLM` in the MERT training to give Moses a chance of a weak domain adaptation.

Unfortunately, our results indicate that for the given domain of EILMT Test data the additional LM brings no improvement. This could be caused by three reasons: 1) the domain of EILMT is very specific and different from `webLM`, 2) `webLM` contains too much noise despite our attempts to keep only real Hindi sentences, 3) the tokenisation of `webLM` is different from EILMT data. The tokenizer we use can be easily trained to mimic any tokenisation style given *both* the original non-tokenized text and the intended tokenized form, but we had no access to the non-tokenized version of EILMT.

2.3 Unsupervised Stem-Suffix Segmentation

To lessen the impact of data sparseness, we wanted to use Moses' ability to work with factors. Splitting each word into the lexical stem and the grammatical suffix (thus defining two separate factors) is the obvious option here. However, being unfamiliar with Hindi morphological analysers and taggers, we decided

tence.

⁵<http://ltrc.iiit.ac.in/nlptools2008/resources2.php>

Parallel data	Distance Reordering		Reordering Using en+hi Forms	
	mteval-v11b	BLEU	mteval-v11b	BLEU
EILMT (7k sents)	.1399	.1382±.146	.1463	.1449±.144
EILMT+Tides+Our Additions (61k sents)	.1447	.1449±.135	.1497	.1500±.144
EILMT+Tides (57k sents)	.1448	.1450±.145	.1504	.1506±.146

Table 1: Impact of reordering model and parallel data size.

Parallel Data	Language Models	Distance Reordering	Reordering Using en+hi Forms
EILMT	EILMT+webLM	.1378±.144	.1437±.145
EILMT	EILMT	.1382±.146	.1449±.144
EILMT+Tides+Additions	EILMT+webLM	.1449±.135	.1500±.144
EILMT+Tides+Additions+webLM	EILMT	.1426±.137	.1506±.144

Table 2: Impact of additional language model.

to employ a tool for unsupervised segmentation of words into morphemes. The tool was originally published in the context of information retrieval, at the Morpho Challenge 2007 shared task (Zeman, 2008).

The tool has been trained on the word types of the Hindi side of the Tides corpus. For every word the algorithm searches for positions where it can be cut in two parts: the stem and the suffix. Then it tries to filter the stem and suffix candidates so that real stems and suffixes remain. The core idea is that real stems occur with multiple suffixes and real suffixes occur with multiple stems. For the purpose of filtering, a collection of stem and suffix candidates that have been observed together is called a *paradigm*.

Various techniques are applied to filter out spurious candidates:

1. If there are more suffixes than stems in a paradigm, the paradigm is removed.
2. If all suffixes in a paradigm begin with the same letter, there is another paradigm where the letter is part of the stem. The former paradigm is removed. Example:
suffixes: त, ति, तियां, तियों
stems: अनुभू, अभिव्यक्, अभिव्यक्, सस्कृ
3. If the suffixes of paradigm B form a subset of suffixes of paradigm A ($A \subset B$)

and there is no C, different from A, such that B is also subset of C: $\forall C \neq A : (B \subset C)$, we add the stems of B to the stems of A, and remove B. A subset paradigm is merged with its superset, as long as there is only one superset candidate.

4. Paradigms with only one suffix are removed.

The following is an example of a paradigm that survived the filtering:

- suffixes: 0, त, े, ी
- stems: अहात खांच घुटन चढ़ाव ज्हटक दायर बचन भाल मदरस मसाल मेमन सितार खर्च

Given the lists of stems and suffixes obtained during training, we want to find the stem-suffix boundary in a word of the same language. Theoretically, we could use the learned stem-suffix combinations to require that both stem and suffix be known. However, this approach proved too restrictive, so we ended up in using just the list of suffixes. If a word ends in a string equal to a known suffix, the morpheme boundary is placed at the beginning of that substring.

2.4 Rule-based Reordering

None of the authors speaks Hindi, however, we obtained some information about its word order. Based on that, we decided to experi-

	EILMT	TIDES
Baseline Moses, Distance Reordering	13.82±1.46	9.28±0.70
Baseline Moses, Reordering Using en+hi Forms	14.49±1.44	9.99±0.71
Rule-based Reordering + Suffix LM+Reord	14.97±1.46	9.78±0.66
Suffix LM+Reord	14.09±1.38	10.17±0.71

Table 3: Impact of rule-based reordering and suffix LM and reordering on EILMT and TIDES datasets.

ment with two reordering transformations of the English sentences, so that their word order gets closer to Hindi. We first tagged the English side using Morče (Votrubec, 2005), then parsed dependencies using the MST parser (McDonald et al., 2005), and finally moved parts of the sentence, based on POS tags and the dependency tree. The following transformations have been applied:

1. Move finite verb forms to the end of the sentence (not crossing punctuation, “that”, WH-words).
2. Transform prepositions to postpositions

Table 3 displays some of the possible Moses configurations using our rule-based reordering and the unsupervised stem-suffix segmentation evaluated both on EILMT domain (7k training, 500 dev, 500 test sentences) and TIDES (50k training, 1k dev, 1k test sentences). The rule-based reordering is used as a simple deterministic pre-processing step before Moses training and translation. The stem-suffix segmentation is applied to the Hindi side of the training data. We use two output factors and two decoding steps: 1) English word forms are translated to Hindi word forms, and 2) suffixes are generated from the hypothesized output word forms. We apply two language models: 3-grams of Hindi forms and 10-grams of suffixes. We also use the target-side suffix factor in the reordering

model. We see that a similar improvement can be achieved by various methods with no clear winning combination.

3 Conclusion

We achieved significant improvements over the baseline Moses as we run it, however the organisers report a baseline BLEU that is beyond our best results. Frankly, we do not have any explanation for this.

The combination of stem-suffix segmentation and rule-based reordering helps for EILMT data, but not for TIDES. Segmentation itself helps on TIDES data but not on EILMT.

References

- Natalia Klyueva and Ondřej Bojar. 2008. UMC 0.1: Czech-Russian-English Multilingual Corpus. In *Proc. of International Conference Corpus Linguistics*, pages 188–195, October.
- Philipp Koehn et al. 2007. Moses: Open Source Toolkit for Statistical Machine Translation. In *Proceedings of ACL 2007*, pages 177–180, Prague, Czech Republic, June.
- Philipp Koehn. 2004. Statistical Significance Tests for Machine Translation Evaluation. In *Proceedings of EMNLP 2004*, Barcelona, Spain.
- Ryan McDonald, Fernando Pereira, Kiril Ribarov, and Jan Hajič. 2005. Non-Projective Dependency Parsing using Spanning Tree Algorithms. In *Proceedings of HLT/EMNLP 2005*, October.
- Franz Josef Och. 2003. Minimum Error Rate Training in Statistical Machine Translation. In *Proc. of the Association for Computational Linguistics*, Sapporo, Japan, July 6-7.
- Jan Votrubec. 2005. Volba vhodné sady rysů pro morfologické značkování češtiny. MSc. thesis.
- Daniel Zeman. 2008. Unsupervised acquiring of morphological paradigms from tokenized text. In Carol Peters et al., editor, *Advances in Multilingual and Multimodal Information Retrieval, CLEF 2007*, volume 5152/2008, pages 892–899. Springer.