

The AFRL-MITLL WMT15 System: There’s More than One Way to Decode It!

Jeremy Gwinnup[†], Timothy Anderson, Michael Kazi[‡], Elizabeth Salesky[‡],
Grant Erdmann, Katherine Young[†],
Christina May[†]
Brian Thompson[‡]

Air Force Research Laboratory
jeremy.gwinnup.ctr,timothy.anderson.20,
grant.erdmann,katherine.young.1.ctr,
christina.may.3.ctr@us.af.mil

MIT Lincoln Laboratory
michael.kazi,elizabeth.salesky,
brian.thompson@ll.mit.edu

Abstract

This paper describes the AFRL-MITLL statistical MT systems and the improvements that were developed during the WMT15 evaluation campaign. As part of these efforts we experimented with a number of extensions to the standard phrase-based model that improve performance on the Russian to English translation task creating three submission systems with different decoding strategies. Out of vocabulary words were addressed with named entity postprocessing.

1 Introduction

As part of the 2015 Workshop on Machine Translation (WMT15) shared translation task, the MITLL and AFRL human language technology teams participated in the Russian–English translation task. Our machine translation systems represent enhancements to both our systems from IWSLT2014 (Kazi et al., 2014) and WMT14 (Schwartz et al., 2014), the addition of hierarchical decoding systems (Hoang and Koehn, 2008), neural network joint models (Devlin et al., 2014) and the utilization of Drem (Erdmann and Gwinnup, 2015), a method of scaled derivative-free trust-region optimization, during the system tuning process.

2 System Description

We submitted systems for the Russian-to-English machine translation shared task. In all submitted systems, we used either phrase-based or hierarchical variants of the Moses decoder (Koehn et al.,

2007). As in previous years, our submitted systems used only the constrained data supplied when training.

2.1 Data Usage

In training our Russian–English systems we utilized the following corpora to train translation and language models: Yandex¹, Commoncrawl (Smith et al., 2013), LDC Gigaword English v5 (Parker et al., 2011) and News Commentary. The Wikipedia Headlines corpus² was reserved to train named entity recognizers.

2.2 Data Preprocessing

As with our WMT14 submission systems, preprocessing to address issues with the training data was required to ensure optimal system performance. Unicode characters in the private use, control character (C0, C1, zero-width, non-breaking, joiner, directionality and paragraph markers), and unallocated ranges were removed. Punctuation normalization and tokenization using Moses preprocessing scripts were then applied before lowercasing the data. The Commoncrawl corpus was further processed as in Schwartz et al. (2014) to exclude wrong-language text and to normalize mixed-alphabet spellings.

2.3 Factored Data Generation

We generated a class-factored version of the parallel Russian–English training data by using `mkcls` to produce 600 word classes for each side of the data. The factored data was then used to create a factored translation model and an in-domain class language model (Brown et al., 1993) for the English portion.

[†]This work is sponsored by the Air Force Research Laboratory under Air Force contract FA-8650-09-D-6939-029.

[‡]This work is sponsored by the Air Force Research Laboratory under Air Force contract FA-8721-05-C-0002.

¹<https://translate.yandex.ru/corpus?lang=en>

²<http://statmt.org/wmt15/wiki-titles.tgz>

2.4 Phrase and Rule Table Training

Phrase tables and rule tables were trained on the preprocessed data using scripts provided with the `moses` distribution. Both rule tables and phrase tables utilized Good-Turing discounting (Gale, 1995). Hierarchical lexicalized reordering models (Galley and Manning, 2008) were also trained for use in the phrase-based systems.

An additional phrase table was trained on the lemmatized forms of the Russian training data. These lemmatized forms were generated by the `mystem`³ tool.

2.5 Language Model Training

The English data sources listed in §2.1 were used to train a very large 6-gram language model (BigLM15). The English portion of the parallel data was processed into class form as outlined in §2.3 to generate an in-domain 600 class language model. `kenlm` (Heafield, 2011) was used to train these 6-gram models. These models were then binarized and stored on local solid-state disks for each machine in our cluster to improve load time and reduce fileservers traffic.

2.6 Operation Sequence Models

Using both the Russian and English data generated in §2.3, we trained order-5 Operation Sequence models (Durrani et al., 2011) for both the surface and class-factored forms of the data. These models improve translation quality by introducing information on the sequence of operations occurring at both the surface and class factor level. These models were then used in our factored phrase-based system.

2.7 Neural Network Joint Models

Neural network joint models (Devlin et al., 2014) are neural network based language models with a source window context. We trained these models on the alignments produced by `mgiza` (Gao and Vogel, 2008) over the parallel training data and then used them to rescore n-best lists. As in (Devlin et al., 2014), we trained four different models. The standard model is “source-to-target, left-to-right,” (s2t, ltr) which evaluates $p(t_i|T, S)$ with target window $T = (t_{i-1}, t_{i-2}, \dots, t_{i-n})$ and $S = (s_{k-m}, \dots, s_k, \dots, s_{k+m})$, where s_k is word-aligned to t_i . The four permutations of this are defined by (a) whether to count upwards from i , in-

³<https://tech.yandex.ru/mystem>

stead of downwards (this is left-to-right vs right-to-left), and (b) whether to swap the sources and targets entirely (source-to-target vs target-to-source).

We experimented with NNJM decoding (via a simple feature function in Moses). We achieved some benefit (+0.48 BLEU) with this approach but rescoring a single NNJM source-to-target on 200-best lists produced better results in this case (+0.90 BLEU). This was on a single system tuned on `newstest2013`, tested on `newstest2014` (baseline 29.07 BLEU). In testing, 2-hidden layer rescoring models outperformed the 1-hidden layer decoding model.

The vocabulary for the NNJMs were created by using all words that appeared at least a certain number of times in the training data. We experimented with minimum counts of 20 and 25. Using 20, our vocabulary was approximately 80,000 Russian words and 40,000 English; with 25, it was 70,000 and 34,000, respectively. We compared rescoring with a single, standard model (s2t, l2r) to rescoring with all directions with results listed in Table 1.

| | Baseline | 1 NNJM | | 4 NNJMs | |
|------|----------|--------|-------|---------|-------|
| | | 20 | 25 | 20 | 25 |
| max | 27.71 | 27.90 | 28.05 | 27.90 | 28.07 |
| mean | 27.48 | 27.61 | 27.81 | 27.67 | 27.60 |

Table 1: NNJM Rescoring on `newstest2015`, optimizing on `newstest2014`, case-insensitive BLEU.

2.8 Processing of Unknown Words

In our submission systems, we allowed words unknown to the decoder to be passed through to the translated output. We developed three post-processing techniques to address unknown words: named entity (NE) tagging and translation (§2.8.2), permissive NE translation (§2.8.3), and selective transliteration of the remaining OOV words (§2.8.4). The first two techniques rely on our in-house transliteration mining of NE pairs, which is described in §2.8.1.

We applied all three post-processing steps to the output of our factored phrase-based submission system; due to time constraints, only the last two steps were applied to the output of our phrase-based and hierarchical submission systems.

Score improvements in uncased BLEU are reported in Table 2. We see that application of

permissive lookup and selective transliteration yielded an improvement of +0.48 BLEU versus a baseline system, while the application of named entity tagging and translation, permissive lookup and selective transliteration yielded a +0.57 BLEU gain.

2.8.1 Transliteration Mining

Both NE processing steps (§2.8.2 and §2.8.3) make use of a NE pairs list that we developed through transliteration mining of the Russian-English CommonCrawl. In transliteration mining (Kumaran et al., 2010; Zhang et al., 2012), we use transliteration as a tool to detect similar-sounding words in the parallel text that may correspond to names. Our process for detecting transliterated NE is generative and rule-based. We used `mystem` to tag NE in the Russian text, and then used capitalization and transliteration as clues to find matching NE in the parallel English sentences. English words were considered candidate matches if they were capitalized, but not sentence-initial; we excluded all-caps words, since acronyms often do not transliterate well. We also required the English candidate words to match the initial sound of the Russian NE.

We checked the initial sound match by transliterating the Russian words according to the textbook values of the Russian letters, and then checking for matches with the English spellings, allowing certain spelling variations. These variations include instances where Russian lacks an English sound, and substitutes a similar sound (e.g., English *h* written in Russian with the letters for *x* or *g*, and English *w* written with the Russian letters for *v* or *u*), as well as common English spelling alternations like *n/kn*, *s/c*, *c/k*, etc.

An iterative process of refining spelling alternations was applied by manual observation of known NE pairs that were not matched via existing rules; notably, this introduced spelling variations for words originating from a third language. For example, English *j* typically represents [dʒ] but may also indicate [h] in words of Spanish origin, so we need to allow the spelling alternation *x/j*. Similarly, the letters *gi* may represent [dʒ] in Italian names like *Giovanni*, so we need to allow transliterated Russian *dzh* to match English *gi*.

At this point in the transliteration mining process, we have derived a list of capitalized English words that have initial spellings potentially matching the initial sound of the Russian NE word. If the

English sentence contains more than one such candidate, we select the word with the smallest edit distance from the Russian transliteration, using a length-normalized Levenshtein distance. For this calculation, any spelling variation counts as an edit distance change, so we penalize variations such as *k* for *c*.

For NE tagging and translation (§2.8.2), we return only the NE pairs with zero edit distance. For permissive NE translation, we allow some variation, as described in §2.8.3.

2.8.2 Named Entity Tagging and Translation

The named entity post-process uses Russian-English pairs in the combined names and titles lists from the Wikipedia Headlines corpus (the “Wiki pairs list”) and the transliteration-mined list (§2.8.1) to replace unknown words with English equivalents. We began by stemming each list to remove Russian noun and adjective endings. To the Wiki pairs list, we added additional pairs yielded by replacing word-internal punctuation marks in existing Wiki pairs with spaces. We used `giza++` (Och and Ney, 2003) to align Russian-English phrases from the Wiki list. We then used these alignments to start a generated list of pairs with only one Russian word and one English word in a pair. Of the aligned pairs, we only included pairs that were aligned with one another three or more times. Only one-to-one alignments would count toward the three alignment rule. We also removed entries where the English word in the pair occurred in a list of stop-words as well as where the English word consisted of only digits. To the generated list, we also added pairs directly from the Wiki list with both single Russian words and single English words. Finally, we also added the highest quality pairs from the transliteration-mined list.

Upon encountering a single word without word-internal punctuation, the system first searches through the generated list, and returns a list of found guesses. If no items are found in the generated list, the Wiki list is then searched. If still no guesses are found, then the transliteration-mined list is searched. The same process occurs for a word containing word-internal punctuation, but after a failed iteration of the search process, the punctuation is replaced with a space and the Wiki lists are searched. Finally if that iteration fails, then the search process occurs on each individual word and a concatenation of English definitions is added

to the guess list for every possible combination of guesses for each component word. An English language model is used to choose among the guesses.

2.8.3 Permissive Named Entity Translation

Permissive NE look-up is applied to translate OOV words that remain untranslated after NE tagging and translation (§2.8.2), or when the NE tagging and translation step is unavailable. In this second step, we expand the NE pairs list to include pairs with greater edit distance when they are validated by repeat occurrence.

While the NE tagging and translation step only uses transliteration-mined NE pairs which match exactly, the permissive step allows NE pairs that have some spelling variation. We apply two additional restrictions to ensure good quality matches, length disparity and instance ratio. We restrict the output to words which come from sentences that do not differ too much in length. A large length disparity suggests a sentence alignment error in the parallel text, which would make the NE match unreliable.

We also restrict the output to words which are fairly frequent among other matches for the same Russian words, calculating an instance ratio as the number of times we see this English word with this Russian word, divided by the total number of English matches we record for this Russian word. Rare instances may be mistakes or spelling variants that we would prefer to exclude. For example, we found the Russian name *Константин* matched with English *Constantine* 117 times, and matched with the spelling *Konstantine* only 1 time, so we do not want to collect *Константин/Konstantine* as a NE pair.

We keep the NE pairs if:

1. The length-normalized edit distance < 0.2
2. The length-normalized edit distance falls between 0.2 and 0.5, inclusive, and sentence length disparity < 2 and instance ratio > 0.01

With these restrictions, we derived 32,560 potential NE pairs.

Subsequently, an additional transliteration mining step was conducted, to collect NE pairs from any capitalized Russian words, not just the words tagged as NE by *system*. We excluded Russian acronyms, sentence-initial words, and personal pronouns (which are capitalized in some styles

of Russian writing). Applying the previously described restrictions for edit distance, instance ratio, and sentence length disparity, we derived an additional 22,370 capitalized-word NE pairs. The combined *system* tagged and capitalized-word NE pairs lists were used in the permissive translation of OOV words, considering both the original form of the Russian OOV word and its stemmed form.

For the phrase-based and hierarchical systems, which were processed without the NE tagging and translation step, the wiki pairs list was added to the mined NE pairs list for permissive OOV translation.

2.8.4 Selective Transliteration of Remaining Out-of-Vocabulary Words

As a final post-processing step, we transliterate some of the remaining OOV words. We attempt to distinguish OOV NE from common words, dropping common words and transliterating names. We hypothesize that retaining transliterated forms of NE will improve readability, even if the output is not a direct match to the English reference.

We attempt to distinguish NE from common words on the basis of capitalization in the Russian source file. Capitalized words that do not begin a sentence are assumed to be NE, and are transliterated. For example, transliteration is the source of the name *Kostenok* in first example sentence shown in Figure 1. Lowercased words, and capitalized words that begin a sentence, are assumed to be common words and are dropped from the output.

3 Results

We submitted three systems for evaluation, each employing a different decoding strategy: traditional phrasal-based, hierarchical, and factored phrasal-based. Each system is described below. Automatically scored results reported in BLEU (Papineni et al., 2002) for our submission systems can be found in Table 3.

Finally, as part of WMT15, the results of our submission systems listed in Tables 3 were ranked by monolingual human judges against the machine translation output of other WMT15 participants. These judgements are reported in WMT (2015).

3.1 Phrasal-Based

We used a standard phrase-based approach, using lowercased data. The lemma-based phrase table

| System | Process Applied | baseline BLEU | postproc BLEU | Δ BLEU |
|--------------|----------------------------------|---------------|---------------|---------------|
| phrase-based | PermLookup + SelTranslit | 27.72 | 28.20 | +0.48 |
| hiero | PermLookup + SelTranslit | 27.43 | 27.91 | +0.48 |
| pb-factored | NEProc + PermLookup+ SelTranslit | 27.18 | 27.75 | +0.57 |

Table 2: NE post-processing improvement measured in uncased BLEU

described in §2.4 was used as a backoff phrase table. We trained a hierarchical lexicalized reordering model, and used two separate class based (factored) language models; one using 600 classes on the in-domain target-side parallel data, and the other using the LDC Gigaword-English v5 NYT corpus. N-best lists from Moses were rescored with 4-way NNJMs, and the system weights were tuned with PRO (Hopkins and May, 2011). Selective transliteration as described in §2.8.4 was then applied to the decoder output.

3.2 Hierarchical

New for this year, we trained a hierarchical system using the same parallel data as our phrase-based systems. The rule table was created as outlined in §2.4 and then filtered to only contain rules relating to the Russian content of the `newstest` test set for years 2012–2015. This filtering was performed in order to reduce the size of the rule table for both system memory requirements and expediency. The incremental-search algorithm (Heafield et al., 2013) and BigLM15 were used to decode the dev (`newstest2014`) and test (`newstest2015`) data. Drem was employed to tune feature weights, optimizing the sum of the expected sentence-level BLEU and expected sentence-level Meteor (Denkowski and Lavie, 2014) metrics. Finally, selective transliteration was employed as described in §2.8.4.

3.3 Factored Phrase-Based

For our last system, we used a factored phrase-based approach (Koehn and Hoang, 2007) where the surface form of the training data was augmented with word classes. These classes were generated on the parallel training data outlined in §2.4 using `mkc1s` to group the words into 600 classes for both English and Russian portions of the parallel training corpus. A phrase table and hierarchical reordering model was then trained using the `moses` training process on both the surface form and the class factor. Order-5 operation sequence models were separately trained on the sur-

face forms and the class factors. An order-6 class-factor LM (Shen et al., 2006) was also trained on the English portion of the parallel training data to supplement the use of BigLM15. NNJMs as outlined in §2.7 were used to rescore the n-best lists from the decode. Following this rescoring, Drem was employed to tune feature weights, optimizing expected corpus-level BLEU (Smith and Eisner, 2006). After optimization and decoding of the test set, remaining unknown words were processed as described in §2.8.2 and §2.8.4.

| System | Cased BLEU | Uncased BLEU |
|--------------|------------|--------------|
| phrase-based | 27.0 | 28.2 |
| hiero | 26.7 | 27.9 |
| pb-factored | 26.4 | 27.8 |

Table 3: MT Submission Systems decoding `newstest2015`

4 Discussion

Our three submitted systems all scored similarly against the official test set. Manual examination of our systems’ output shows that there are significant differences in sentence structure and content.

4.1 Comparing Submitted Systems for Similarity

We scored one system output against another (as reference) with `mteval13a.pl` in both directions as BLEU scores are not symmetric. Results are listed in Table 4. Interestingly, the factored phrase-based and hierarchical systems were more similar to each other than to the traditional phrase-based system. This suggests that the addition of class factors serves a similar function to the use of hierarchical decoding.

4.2 A Closer Analysis of Performance between Submission Systems

We now examine two sentences translated with each of our submission systems and compare them with the supplied reference translation and a literal

| Test | Ref | BLEU |
|-------|-------|-------|
| PB | Hiero | 57.18 |
| PBFac | Hiero | 76.34 |
| Hiero | PB | 57.09 |
| PBFac | PB | 60.54 |
| PB | PBFac | 60.47 |
| Hiero | PBFac | 70.18 |

Table 4: Submission system similarity measured in uncased BLEU

translation. These comparisons are shown in Figure 1.

In the first sentence, the reference translation shows a reordering of the first clause to the end. The phrase-based system drops this clause. The pb-factored system has *informed* instead of *reported* which shifts the meaning; perhaps the translation was influenced by the fluent but different-meaning phrase *informed the Minister*. The hierarchical system follows the original order of the source sentence clauses; while missing *the*, it reads the best overall.

In the second sentence, *Учебный* “school” (adjective) is the probable source of *school*, *academic*, and *teach*. The phrase-based system handles this word best; the phrase-based factored system generates *academic* and *teach* but separates them; the hierarchical system generates *year to teach*. The hierarchical system does the best job with *no earlier than October*. The phrase-based factored system generates *no earlier* and *October* but reorders them (perhaps influenced by the common phrase, *in October*); and the phrase-based system creates *before October*, which reverses the meaning. The phrase-based system would have read best here, had it not neglected the negative particle.

5 Conclusion

In this paper, we present data preparation and processing techniques for our Russian–English submissions to the 2015 Workshop on Machine Translation (WMT15) shared translation task. Our submissions examine three different decoding strategies and the effectiveness of sophisticated handling of unknown words. While scoring similarly, each system produced markedly different output.

Opinions, interpretations, conclusions and recommendations are those of the authors and are not necessarily endorsed by the United States Government. Cleared for public release

References

- Peter Brown, Vincent Della Pietra, Stephen Della Pietra, and Robert Mercer. 1993. The mathematics of statistical machine translation: parameter estimation. *Computational Linguistics*, 19(2):263–311, June.
- Michael Denkowski and Alon Lavie. 2014. Meteor universal: Language specific translation evaluation for any target language. In *Proceedings of the Ninth Workshop on Statistical Machine Translation*, pages 376–380, Baltimore, Maryland, USA, June. Association for Computational Linguistics.
- Jacob Devlin, Rabih Zbib, Zhongqiang Huang, Thomas Lamar, Richard Schwartz, and John Makhoul. 2014. Fast and robust neural network joint models for statistical machine translation. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Proceedings of the ACL, Long Papers, pages 1370–1380, Baltimore, MD, USA.
- Nadir Durrani, Helmut Schmid, and Alexander Fraser. 2011. A joint sequence translation model with integrated reordering. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics (ACL ’11)*, pages 1045–1054, Portland, Oregon, June.
- Grant Erdmann and Jeremy Gwinnup. 2015. Drem: The AFRL submission to the WMT15 tuning task. In *Proceedings of the Tenth Workshop on Statistical Machine Translation (WMT’15)*, Lisbon, Portugal, September. To appear.
- William A. Gale. 1995. Good-turing smoothing without tears. *Journal of Quantitative Linguistics*, 2.
- Michel Galley and Christopher D. Manning. 2008. A simple and effective hierarchical phrase reordering model. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, EMNLP ’08, pages 848–856.
- Qin Gao and Stephan Vogel. 2008. Parallel implementations of word alignment tool. In *Software Engineering, Testing and Quality Assurance for Natural Language Processing*, pages 49–57, Columbus, Ohio, June.
- Kenneth Heafield, Philipp Koehn, and Alon Lavie. 2013. Grouping language model boundary words to speed k-best extraction from hypergraphs. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 958–968, Atlanta, Georgia, USA, June.
- Kenneth Heafield. 2011. KenLM: faster and smaller language model queries. In *Proceedings of the EMNLP 2011 Sixth Workshop on Statistical Machine*
on 11 Jun 2015. Originator reference number RH-15-114103. Case number 88ABW-2015-2973.

Example 1

| | |
|---------------|--|
| source: | Об этом сообщил министр образования и науки самопровозглашенной республики Игорь Костенок |
| literal: | Of this reported minister education and science self-proclaimed republic, Igor Kostenok |
| reference: | The Minister of Education and Science for the self-proclaimed republic, Igor Kostenok, reported this. |
| phrase-based: | The Minister of Education and Science of the self-declared republic, Igor Kostenok. |
| pb-factored: | This was informed the Minister of Education and Science of the self-declared republic, Igor Kostenok. |
| hierarchical: | This was announced by Minister of Education and Science of the self-proclaimed republic Igor Kostenok. |

Example 2

| | |
|---------------|---|
| source: | Учебный год в ДНР начнется не раньше октября. |
| literal: | School year in DNR begins not before October. |
| reference: | The academic school year in the DPR will begin no earlier than October. |
| phrase-based: | The school year in DNR will begin before October. |
| pb-factored: | Academic year in October, teach will begin. |
| hierarchical: | School year to teach will begin no earlier than October. |

Figure 1: Comparison of Submission System Translation Output

- Translation*, pages 187–197, Edinburgh, Scotland, United Kingdom, July.
- Hieu Hoang and Philipp Koehn. 2008. Design of the Moses decoder for statistical machine translation. In *Software Engineering, Testing, and Quality Assurance for Natural Language Processing*, SETQA-NLP '08, pages 58–65.
- Mark Hopkins and Jonathan May. 2011. Tuning as ranking. In *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing (EMNLP '11)*, pages 1352–1362, Edinburgh, Scotland, U.K.
- Michael Kazi, Elizabeth Salesky, Brian Thompson, Jessica Ray, Michael Coury, Tim Shen, Wade Anderson, Grant Erdmann, Jeremy Gwinnup, Katherine Young, Brian Ore, and Michael Hutt. 2014. The MIT-LL/AFRL IWSLT-2014 MT system. In *Proceedings of the 11th International Workshop on Spoken Language Translation (IWSLT'14)*, pages 65–72, Lake Tahoe, California, December.
- Philipp Koehn and Hieu Hoang. 2007. Factored translation models. In *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, pages 868–876.
- 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 ACL on Interactive Poster and Demonstration Sessions*, ACL '07, pages 177–180, Stroudsburg, PA, USA. Association for Computational Linguistics.
- A Kumaran, Mitesh M. Khapra, and Haizhou Li. 2010. Whitepaper of news 2010 shared task on transliteration mining. In *Proceedings of the 2010 Named Entities Workshop*, pages 29–38, Uppsala, Sweden, July. 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, March.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: A method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics (ACL '02)*, pages 311–318, Philadelphia, Pennsylvania, July.
- Robert Parker, David Graff, Junbo Kong, Ke Chen, and Kazuaki Maeda. 2011. English Gigaword Fifth Edition LDC2011T07. *Philadelphia: Linguistic Data Consortium*.
- Lane Schwartz, Timothy Anderson, Jeremy Gwinnup, and Katherine Young. 2014. Machine translation and monolingual postediting: The AFRL WMT-14 system. In *Proceedings of the Ninth Workshop on Statistical Machine Translation (WMT'14)*, pages 186–194, Baltimore, Maryland, USA, June.
- Wade Shen, Richard Zens, Nicola Bertoldi, and Marcello Federico. 2006. The JHU workshop 2006 IWSLT system. In *Proc. of the International Workshop on Spoken Language Translation*, Kyoto, Japan, November.

David A. Smith and Jason Eisner. 2006. Minimum risk annealing for training log-linear models. In *Proceedings of the COLING/ACL 2006 Main Conference Poster Sessions*, pages 787–794, Sydney, Australia.

Jason R. Smith, Herve Saint-Amand, Magdalena Plamada, Philipp Koehn, Chris Callison-Burch, and Adam Lopez. 2013. Dirt cheap web-scale parallel text from the common crawl. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (ACL '13)*, pages 1374–1383, Sofia, Bulgaria, August.

WMT. 2015. Findings of the 2015 Workshop on Statistical Machine Translation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation (WMT '15)*, Lisbon, Portugal, September.

Min Zhang, Haizhou Li, A. Kumaran, and Ming Liu, 2012. *Proceedings of the 4th Named Entity Workshop (NEWS) 2012*, chapter Report of NEWS 2012 Machine Transliteration Shared Task, pages 10–20. Association for Computational Linguistics.