

# Statistical Machine Translation of Croatian Weather Forecasts: How Much Data Do We Need?

---

Nikola Ljubešić, Petra Bago and Damir Boras

Department of Information Sciences, Faculty of Humanities and Social Sciences, University of Zagreb, Croatia

This research is the first step towards developing a system for translating Croatian weather forecasts into multiple languages. This step deals with the Croatian-English language pair. The parallel corpus consists of a one-year sample of the weather forecasts for the Adriatic, consisting of 7,893 sentence pairs. Evaluation is performed by the automatic evaluation measures BLUE, NIST and METEOR, as well as by manually evaluating a sample of 200 translations. We have shown that with a small-sized training set and the state-of-the-art Moses system, decoding can be done with 96% accuracy concerning adequacy and fluency. Additional improvement is expected by increasing the training set size. Finally, the correlation of the recorded evaluation measures is explored.

*Keywords:* statistical machine translation, automatic evaluation, manual evaluation, correlation between evaluation measures

## 1. Introduction

Machine translation (MT) has a long history dating back to the 1940s. Its progress has been motivated by advances in computer science and artificial intelligence. Traditional MT was rule-based and has relied on various levels of linguistic analysis on the source side and language generation on the target side.

In the late 1980s, the first statistical approach to machine translation (SMT) was pioneered by a group of researchers from IBM [2]. Since then, SMT has advanced from word-based to phrase-based models. At the beginning, SMT relied on the source-channel model consisting of a translation and a language model. The translation model ensures that the system produces target hypotheses corresponding to the source

sentence, while the language model ensures that the output is as grammatical and fluent as possible. Although early SMT models essentially ignored linguistic aspects, nowadays efforts are made to reintroduce linguistic information in both the translation and the language models [7].

Evaluating the output of an MT system is certainly not a simple task. Methods are usually divided into automatic and human.

Automatic measures rely on reference translations of source sentences and calculate the likeness of the system output and the reference translations. The best known automatic measures are BLEU, NIST and METEOR.

BLEU [14] is the geometric mean of modified n-gram precisions for different n-gram lengths (usually from one to four), multiplied by a factor (brevity penalty) that penalizes producing short sentences containing only highly reliable portions of the translation.

NIST [5] is the arithmetic mean of clipped n-gram precisions for different n-gram lengths multiplied by a brevity penalty. Also, when computing the NIST score, n-grams are weighted according to their frequency, so that less frequent (and thus more informative) n-grams are given more weight.

While BLEU and NIST are based on precision, METEOR [1] calculates both recall and precision on the unigram level, assigning in the harmonic mean more weight to recall than precision. Additionally, METEOR enables matching on the stem and synonym level. All three

measures correlate highly (around 0.9) with human judgments at the corpus level. METEOR is reported to correlate higher than BLEU and NIST [1]. Another advantage of METEOR is that it produces also scores on sentence level. However, the correlation with human judgment on sentence level is much lower than on corpus level (0.403 in [1]).

Human evaluation mostly consists of scoring every translation by adequacy and fluency on a scale from 0 to 5. Adequacy indicates the extent to which the information contained in the source is included in the translation, whereas fluency measures how grammatical and natural the translation sounds.

## 2. Experimental Design

### 2.1. The corpus

The corpus used in this research is a one-year sample of the weather forecasts for the Adriatic, published by the Croatian Meteorological and Hydrological Service [4]. The forecasts were published twice a day in four languages: Croatian, English, German and Italian.

This research deals only with the Croatian-English language pair. The pair consists of 720 documents and 2800 paragraphs (4 paragraphs/sections per document). Building the translation model and the decoding is performed with the Moses system [12]. The input for training the translation model with Moses is a sentence-aligned corpus. For this reason, the corpus is automatically sentence split and tokenized. The Croatian part consists of 8,409 sentences and the English part consists of 8,368 sentences. Furthermore, the sentences are aligned by the Gale & Church sentence alignment algorithm [6]. For the sake of simplicity, only those sentences that translate into one sentence are chosen. The resulting sentence-aligned corpus consists of 7,893 sentence pairs. Thereby, some 6% of the data is lost.

The aligned corpus is described through some basic statistics in Table 1. The statistics show that, as expected, Croatian has a lower token count, but a higher type count due to its rich inflectional morphology. Furthermore, English sentences tend to consist of more characters and

more words. In general, the type-token ratio emphasizes the overall simplicity of the text with only 802 (Croatian), ie. 592 (English) types on approximately 100,000 tokens.

language	Croatian	English
mean(wps)	10.589	13.652
mean(cps)	69.136	72.237
count(token)	87681	111944
count(type)	802	592
type-token ratio	0.00915	0.00529

Table 1. Corpus statistics (wps – words per sentence, cps – characters per sentence).

### 2.2. Research questions

The basic questions this research deals with are:

1. How much data is necessary for training a good translation model for the text complexity level of weather forecasts?
2. What is the correlation between the four most popular measures used in the evaluation of statistical machine translation?

#### 2.2.1. Training

The first question will be answered through an experiment where the translation and language models are trained on ten different corpus sizes. In these ten steps, the corpus size, on which the models are trained and tested, grows from 789 to 7,890 sentence pairs.

To ensure a good estimate of the calculated measures, every step is repeated ten times. In each iteration, the sample is built from scratch by selecting sentence pairs randomly from the whole corpus. After building a sample, the sample is split into a 9 : 1 ratio – a training and a test set, respectively.

The training procedure consists of training the language model with the SRILM tool [15] on the English training instance using Kneser-Ney smoothing [9], which has often proven to achieve the best results [3]. The translation model is trained in Moses [8] with default settings. These are defined in the script `train-factored-phrase-model.perl` [10].

### 2.2.2. Evaluation

After training the translation model, Croatian sentences are decoded using the default Moses settings. The output of the decoding step is evaluated by three previously described methods – BLEU, NIST (implementations mteval-v13a [13]) and METEOR (implementation meteor-1.0 [11]).

Unknown words are also monitored and recorded through the unknown word rate (UWR) which is the percentage of words in the source that have no translation in the translation model. It is important to note that the UWR measure is not an evaluation measure, since it is constant for a specific training and test set and does not depend on the machine translation method.

The recorded measures are used for answering both research questions: calculating the progress rate as the corpus size increases, and observing the correlation between the four measures.

At the end of this research, human evaluation is performed on a sample of 200 sentences. The primary goal of this evaluation is to achieve a clear insight in the quality of the translations. Error types are also recorded giving additional information regarding the causes of the observed translation errors. Results of human evaluation are compared to METEOR and UWR, since only these measures are capable of calculating agreement on the sentence level.

## 3. Results

### 3.1. Automatic evaluation

As described in the experimental design, three automatic evaluation metrics and the unknown word rate (UWR) are recorded as the corpus size increases in ten steps. In each step ten iterations are undertaken. The results are normalized to the [0, 1] scale for easy visual and numerical comparison. Thereby, UWR changes its sign since it grows as the translation quality decreases. Figure 1 shows the mean of the ten iterations on the four measures as corpus size increases in ten steps. The mean of all four measures is shown with a full line (ALL).

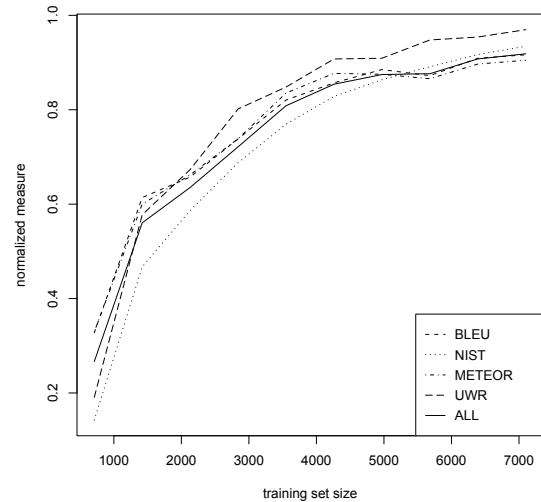


Figure 1. Normalized recorded measures as the training set size increases.

The results show an obvious consistency between different measures. Additional improvement is to be expected as the corpus size increases. Interestingly, the NIST measure has the smoothest curve showing the least sensitivity to different data.

In the next step, the relationship between the four measures is shown in a scatter plot in Figure 2.

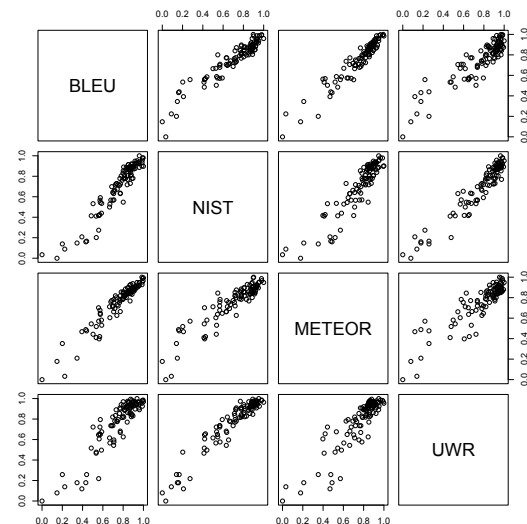


Figure 2. Scatter plot of the four recorded measures.

The correlation coefficients of the four variables are shown in Table 2. The data show that BLEU and METEOR have the most consistent results. BLEU and NIST, in theory the most similar measures, are second most consis-

tent, while METEOR and UWR have the lowest correlation coefficient.

	NIST	METEOR	UWR
BLEU	0.95179	0.95590	0.91322
NIST		0.91960	0.95803
METEOR			0.89478

Table 2. Pearson correlation coefficient between the four measures.

Since NIST has the smoothest curve as corpus size increases, its data is used to calculate the progress in specific steps. In the last three steps, the progress is 3.2%, 2.9% and 1.9%, respectively. These numbers show, as does Figure 1, the possibility of further improvement of results by increasing the corpus size.

Regarding the unknown word rate, it is on average 5.52% on the smallest corpus, whereas on the biggest corpus it drops to 0.86%. Percentages of decrease, as the corpus size increases, correspond to the previous numbers given for the NIST metric.

Additionally, standard deviation of a specific metric is calculated to examine the consistency of the results. Standard deviation is calculated on normalized results for the tenth and final step. The results for the specific measures, together with the non-normalized mean and standard deviation, are given in Table 3. The normalized standard deviation measure shows that unknown word rate deviates least from its central tendency. NIST is the metric with the lowest standard deviation among the evaluation metrics. This corresponds to the smoothness of its curve in Figure 1. BLEU is the metric with the highest deviation.

measure	mean	sd	sd norm
BLEU	0.718	0.008	0.053
NIST	9.111	0.083	0.037
METEOR	0.845	0.005	0.045
UWR	0.007	0.001	0.021

Table 3. Mean and standard deviation of results of the four non-normalized and normalized measures in the final step.

### 3.2. Human evaluation

At the end of this research, human evaluation of 200 translations is undertaken to get a clear pic-

ture of the quality of the automated translation. Out of ten experiments with the largest corpus size, sampling is done from the results that contain most of the medians of the four recorded measures. Out of these 7890 translations, 200 random translations are chosen and given to the human evaluator for manual evaluation.

The human evaluator is first given the target translation to evaluate its fluency and later the source to evaluate adequacy. Both adequacy and fluency are graded on a scale from 0 to 5.

In addition, if the grade is less than 5 on any of the criteria, the error type is also recorded. There are four error types with some examples (S – source, T – automatic translation, R – reference translation):

- all lexical items correct, but meaning changed by word order or punctuation  
T: In the rest of the Adriatic SE and SW wind 4-14, in the open sea up to 24 knots, **wind, diminishing**.  
R: In the rest of the Adriatic SE and SW wind 4-14, in the open sea up to 24 knots, also diminishing.
- lexical item translated incorrectly  
T: Sea 1-2, in the south Adriatic, **during and elsewhere** in the afternoon, 2-3.  
R: Sea 1-2, in the south, in the afternoon in the rest of the Adriatic as well, 2-3.
- unknown word in the source  
T: **Ujutro**, and again during the night a risk of fog patches.  
R: A risk of fog patches this morning and again during the night.
- typing error in the source  
S: **Tempertaure** zraka iste ili malo niže.  
T: **Tempertaure** temperatures with no change or dropping a little.  
R: Air temperatures with no change or a little lower.

Frequency of these error types is given in Table 4.

error type	absolute	relative
type 1	23	0.397
type 2	20	0.345
type 3	10	0.172
type 4	5	0.086

Table 4. Frequency of error types.

	UWR	METEOR
HE	0.52325	0.22604
UWR		0.24777

Table 5. Pearson's correlation coefficient between human evaluation (HE) and the UWR and METEOR measures.

From the grades given by the human evaluator, accuracy is calculated as the percentage of the assigned grades regarding the maximum grade. The accuracy given by the human evaluation is 96.15% on the sentence level. If the length of the sentence is taken into account, accuracy drops to 93.631%.

Since only a part of a sample was evaluated by humans, it is impossible to compare the result of human evaluation and automated evaluation on the corpus level. Out of four recorded measures, METEOR and UWR can also be calculated on the sentence level. The Pearson correlation coefficient between human evaluation (HE) and these two measures are given in Table 5. Correlation between human evaluation and unknown word rate is over 0.5, while METEOR and human evaluation correlate with only 0.22.

#### 4. Conclusion

In this research we have shown that with a small-sized sentence-aligned parallel corpus and the state-of-the-art Moses system, decoding can be done with 96% accuracy concerning adequacy and fluency. It is important to note that this domain-specific text is very simple, having only 600 types on almost 100,000 tokens.

As corpus size increases, automatic evaluation measures behave in a typical logarithmic way. With around 7,000 sentence pairs of training data, improvement falls down to 2%. Additional training data could further improve the results.

The relationship between the automatic evaluation measures BLEU, NIST and METEOR is also explored. All these metrics correlate very highly with each other as well as with the negated UWR.

Exploring the correlation of METEOR and UWR with the human evaluation on sentence level

shows a good correlation with UWR, but, as expected, a low correlation with METEOR.

The behavior of these automatic evaluation measures is still rather unknown, and we believe that this research has shed some light on it.

#### References

- [1] S. BANERJEE, A. LAVIE, METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, (2005) Ann Arbor, Michigan. Association for Computational Linguistics, pp. 65–72.
- [2] P. F. BROWN, J. COCKE, S. A. D. PIETRA, V. J. D. PIETRA, F. JELINEK, J. D. LAFFERTY, R. L. MERCER, P. S. ROOSSIN, A statistical approach to machine translation. *Computational Linguistics*, 16(2) (1990), 79–85.
- [3] S. CHEN, J. GOODMAN, An empirical study of smoothing techniques for language modeling. In *Proceedings of the 34th Annual Meeting of the ACL*, (1996), pp. 310–318.
- [4] PROGNOZA VREMENA ZA JADRAN, [http://prognoza.hr/prognoze.php?id=jadran\\_h](http://prognoza.hr/prognoze.php?id=jadran_h), 2010.
- [5] G. DODDINGTON, Automatic evaluation of machine translation quality using n-gram co-occurrence statistics. In *Proceedings of the second international conference on Human Language Technology Research*, (2002) San Francisco, CA, USA, pp. 138–145. Morgan Kaufmann Publishers Inc.
- [6] W. A. GALE, K. W. CHURCH, A program for aligning sentences in bilingual corpora. *Computational Linguistics*, 19(1) (1993), 75–102.
- [7] C. GOUTTE, N. CANCEDDA, M. DYMETMAN, G. FOSTER, *Learning Machine Translation*. The MIT Press, 2009.
- [8] H. HOANG, A. BIRCH, C. CALLISON-BURCH, R. ZENS, R. AACHEN, A. CONSTANTIN, M. FEDERICO, N. BERTOLDI, C. DYER, B. COWAN, W. SHEN, C. MORAN, O. BOJAR, Moses: Open source toolkit for statistical machine translation. In *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics*, (2007), pp. 177–180.
- [9] R. KNESER, H. NEY, Improved backing-off for m-gram language modeling. In *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing*, (1995), pp. 181–184.
- [10] P. KOEHN, *Moses – Statistical Machine Translation System*. University of Edinburgh, 2009.
- [11] THE METEOR AUTOMATIC MACHINE TRANSLATION EVALUATION SYSTEM, <http://www.cs.cmu.edu/~alavie/meteor/>.

- [12] MOSES – STATISTICAL MACHINE TRANSLATION SYSTEM, <http://www.statmt.org/moses/>, 2010.
- [13] INFORMATION TECHNOLOGY LABORATORY TOOLS, <http://www.itl.nist.gov/iad/mig/tools/>.
- [14] K. PAPINENI, S. ROUKOS, T. WARD, W. JING ZHU, Bleu: A method for automatic evaluation of machine translation. In *ACL '02: Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*, (2002), pp. 311–318.
- [15] A. STOLCKE, Srilm - an extensible language modeling toolkit. In *Proceedings of International Conference on Spoken Language Processing*, (2002).

Received: June, 2010

Accepted: November, 2010

Contact addresses:

Nikola Ljubešić  
 Department of Information Sciences  
 Faculty of Humanities and Social Sciences  
 University of Zagreb  
 Ivana Lučića 3, 10000 Zagreb  
 Croatia  
 e-mail: nljubesi@ffzg.hr

Petra Bago  
 Department of Information Sciences  
 Faculty of Humanities and Social Sciences  
 University of Zagreb  
 Ivana Lučića 3, 10000 Zagreb  
 Croatia  
 e-mail: pbago@ffzg.hr

Damir Boras  
 Department of Information Sciences  
 Faculty of Humanities and Social Sciences  
 University of Zagreb  
 Ivana Lučića 3, 10000 Zagreb  
 Croatia  
 e-mail: dboras@ffzg.hr

---

NIKOLA LJUBEŠIĆ is a senior research assistant at the Department of Information Sciences, Faculty of Humanities and Social Sciences, University of Zagreb. He received his PhD degree in 2009 on the topic of event detection in newspaper texts. His research interests include statistical methods in natural language processing and building language resources from the web. He is a member of the Croatian and Slovenian Language Technologies Society.

---



---

PETRA BAGO received her diploma in information sciences from the Faculty of Humanities and Social Sciences, University of Zagreb, Croatia, July 2008. She is currently working on her PhD and is employed as a research assistant at the Department of Information Sciences. Her research fields are: digital humanities, computational lexicography and encyclopaedic science.

---



---

DAMIR BORAS is the chairman and founder of the Chair for Lexicography and Encyclopaedic Science at the Department of Information Sciences, Faculty of Humanities and Social Sciences, University of Zagreb, Croatia, 2004. He received his PhD degree in information sciences from the same department in 1998 for a dissertation on theory and rules of segmenting text written in Croatian. He also holds the engineer diploma in radio communications received from the Faculty of Electrical Engineering, University of Zagreb, Croatia, 1974. Currently, he is a full professor at the Department of Information Sciences, Faculty of Humanities and Social Sciences, University of Zagreb, Croatia. His research interests include: lexicography, encyclopaedic science, digital humanities and natural language processing.

---