

# Morphological Pre-Processing for Turkish to English Statistical Machine Translation

Arianna Bisazza    Marcello Federico

FBK - Fondazione Bruno Kessler  
Via Sommarive 18, 38100 Povo (TN), Italy

{bisazza, federico}@fbk.eu

## Abstract

We tried to cope with the complex morphology of Turkish by applying different schemes of morphological word segmentation to the training and test data of a phrase-based statistical machine translation system. These techniques allow for a considerable reduction of the training dictionary, and lower the out-of-vocabulary rate of the test set. By minimizing differences between lexical granularities of Turkish and English we can produce more refined alignments and a better modeling of the translation task. Morphological segmentation is highly language dependent and requires a fair amount of linguistic knowledge in its development phase. Yet it is fast and light-weight – does not involve syntax – and appears to benefit our IWSLT09 system: our best segmentation scheme associated to a simple lexical approximation technique achieved a 50% reduction of out-of-vocabulary rate and over 5 point BLEU improvement above the baseline.

## 1. Introduction

Morphology plays a fundamental role in any NLP application involving agglutinative languages, such as Turkish. This is particularly true for statistical machine translation (SMT) from Turkish into English, in which the mismatch between the word formation mechanisms of the two languages severely contributes to the difficulty of the task. High word granularity differences reflect indeed on much higher data sparseness on the Turkish side and on the impossibility to properly model alignments between English and Turkish words. We approached this problem through morphological segmentation of Turkish, by taking advantage of linguistic knowledge of both the source and target languages. In particular we focused on the comparison of different segmentation rule sets in order to find an effective preprocessing scheme for the Turkish-English task of IWSLT09. This work shows how a specific linguistic preprocessing can benefit a purely statistics-based, language-independent NLP application like SMT.

The paper is organized as follows: the linguistic features of Turkish that are relevant to SMT and motivate our research will be presented in Section 2. After a brief overview of related work (Section 3) we will describe in detail the prepro-

cessing technique and the different segmentation schemes we have implemented (Section 4). In Section 5 the experimental results will be introduced and commented following several axes: segmentation schemes, distortion limit and finally lexical approximation. This will be followed by a global discussion of our approach and by an exposition of future works (Section 6).

## 2. Turkish Morphology and MT

Several linguistic features of Turkish [1] can directly affect the performance of an SMT system: (i) agglutinative morphology, (ii) vowel harmony and other phoneme alternation phenomena, and (iii) word order. Whereas the first two features are situated at the word level, the third concerns syntax and the global structure of sentences. In this work, our analysis focuses on word-level preprocessing, which we expect will then open the way to the exploration of sentence-level techniques of reordering.

### 2.1. Agglutination

Agglutination implies that the vocabulary is built by a wide range of basic suffix combinations. A Turkish word can thus correspond to a single English word, up to phrases of various length, or even to a whole sentence as shown in Table 1. Differences in token range can be observed in the IWSLT09 training parallel corpus, whose Turkish side is composed of around 139,000 tokens as opposed to the 182,000 tokens of the English side.

<i>oda</i>	‘room’
<i>odam</i>	‘my room’
<i>odamda</i>	‘in my room’
<i>odamdayım</i>	‘I am in my room’

Table 1: *Example of Turkish suffixation*

Given this premise, it is easy to imagine how inter-language alignments and in general any modeling of the language based on the notion of token may suffer from data sparseness. That is why morphological segmentation is needed as preprocessing.

## 2.2. Phonology

On a phonological level vowel harmony and other phoneme alternation phenomena systematically lead stems and suffixes to have several surface forms – i.e. allomorphy. For example (see Table 2) the possessive suffix *-(Im)* ‘my’ can have four different surface forms depending on the last vowel of the word it attaches to (ex.1-4), plus one if it is attached to a word ending with vowel (ex.5).

1)	<i>saç</i>	+ ( <i>Im</i> )	→	<i>saç<u>ım</u></i>	‘my hair’
2)	<i>el</i>	+ ( <i>Im</i> )	→	<i>el<u>im</u></i>	‘my hand’
3)	<i>kol</i>	+ ( <i>Im</i> )	→	<i>kol<u>um</u></i>	‘my arm’
4)	<i>göz</i>	+ ( <i>Im</i> )	→	<i>göz<u>üm</u></i>	‘my eye’
5)	<i>kafa</i>	+ ( <i>Im</i> )	→	<i>kaf<u>am</u></i>	‘my head’

Table 2: *Different surface forms of possessive suffix -(Im).*

If we envisage treating suffixes as single tokens, we foresee an additional increase of data sparseness due to suffix allomorphy. To cope with this problem, we need to introduce an abstract notation that factorizes different surface forms of the same suffix into one single form.

## 2.3. Word order

The typical structure of Turkish phrases is head-final. Sentences mostly belong to the subject-object-verb (SOV) kind, but word order is relatively free and discourse-related phrase movements are quite frequent. As a result alignments between Turkish and English are far from being monotonic, as shown by this example taken from the IWSLT09 corpus:

*Banyolu iki kişilik bir oda istiyorum.*  
 [with-bath] [two] [for-people] [a] [room] [I want]  
 ‘I’d like a twin room with a bath please.’

Although reordering rules seem hard to describe without using any syntactic information, we believe that morphological segmentation is a first necessary step to take in order to enable machine learning of refined alignments and complex word reordering patterns.

## 3. Related work

Morphological preprocessing of Turkish has been recently investigated by [2] in the context of an English to Turkish SMT system. The opposite translation direction, and the higher complexity of the language represented in the data made that task considerably different from the one we used for our experiments. As we were not concerned with the generation of morphologically complex words, we didn’t work on the target language model but focused on comparing the impact of different segmentation schemes applied to the source. For this purpose we referred to the methodology exposed by [3] on an Arabic-English task. Arabic is also morphologically rich but its segmentation schemes are

much simpler than those for Turkish, given that the number of involved clitics and suffixes is typically smaller<sup>1</sup>. As pointed by [2], Turkish employs about 30,000 root words and about 150 distinct suffixes. Although not all possible suffix combinations are grammatical, the number of potential inflected/derived forms of a given root word is still extremely high. This implies that linguistic knowledge becomes crucial to guide the investigation of meaningful segmentation schemes among all possible rule combinations. Another difference with respect to [3] is that in our work we consider not only splitting but also removing suffixes.

In previous editions of IWSLT, [4] and [5] tried to further decrease the out-of-vocabulary rate of the Arabic test set by a so-called lexical approximation approach. This idea consists in finding words of the training that are morphologically close to OOVs and introducing them into the translation process by various techniques – i.e. best replacer computation at run-time in [4] or phrase table expansion in [5]. This method was shown to have a positive effect on Arabic-English SMT systems. In this work, we developed for Turkish a technique of lexical approximation similar to [4] and tested it on the Turkish-English task of IWSLT09.

## 4. Morphological segmentation schemes

### 4.1. Preprocessing technique

Our preprocessing workflow starts with morphological analysis, which consists in running K. Oflazer’s [6] suffix combinatory FSTs to each entry of the corpus dictionary. This operation is carried out through the *lookup* command of the Xerox Finite-State Tool’s suite [7]. As more than one analysis is often possible (with differences in features but also in the lemma), disambiguation is performed on the words in context through the perceptron-based tool developed by [8] (see an example of disambiguation in Table 3).

‘Are there any tours of famous stars’ homes?’		
<i>Ünlü yıldızların evine turlar var mı ?</i>		
	<i>ev+Noun+A3sg+P2sg+Dat</i>	[to your house]
→	<b><i>ev+Noun+A3sg+P3sg+Dat</i></b>	<b>[to his/her/its house]</b>
	<i>evin+Noun+A3sg+Pnon+Dat</i>	[to the kernel]

Table 3: *Morphological disambiguation of a Turkish word in context.*

As a result of this process, each token is replaced by its lemma followed by a sequence of tags representing lexical features of the analyzed word. While some of these tags actually have a surface realization, some others simply encode morphological features (e.g. *Noun* and *Verb* indicate lexical category, *A3sg* stands for ‘singular’, and *Pnon* for ‘no possessive’). The use of feature tags provides a means to abstract

<sup>1</sup>As opposed to Arabic – a Semitic language – Turkish belongs to the Turkic language group, itself part of the larger Altaic family according to many linguists.

from suffix allomorphy. For example all the forms of possessive suffix *-(I)m* of Table 2 are replaced by the symbol *P3sg*. The advantage of using lexical features (e.g. *P3sg*, *Fut*, *Fut-Part*) in place of the suffixes themselves (eg. *sH*, *yAcAk*) is that features are less ambiguous and make our rules more readable. Hence, we can now define different segmentation rules on the tags by using simple regular expressions<sup>2</sup>.

#### 4.2. Segmentation schemes

The schemes presented below are different combinations of rules determining the splitting or removal of tags from the analyzed words. The approach is incremental since a scheme includes all or most of the rules of the previous one. In this work, we mainly focused on nominal suffixation and also defined a few rules for the segmentation of verb forms. In order to find an effective rule set we tested eleven morphological segmentation schemes named *MS[1..11]*, but only the most meaningful among them will be described in the following<sup>3</sup>.

**MS2: Cases.** Nominal cases that are expected to have an English counterpart are split off from words: these are namely dative, ablative, locative and instrumental, often aligning with the English prepositions ‘to’, ‘from’, ‘in’ and ‘with/by’, respectively. The remaining case tags – nominative, accusative and genitive – are instead removed from the words because they are not expected to have English counterparts.

**MS6: Cases & Poss.** After treating case tags we remove the tag meaning absence of possessive suffixes and split off from nouns the possessive tags of all persons except the 3rd singular (*P3sg*), which is indeed removed. In fact, the latter often aligns with nothing on the target side: namely when it functions as marker of a noun compound’s head (ex.1) or when the possessor is expressed as a noun in the genitive case (ex.2). Yet, if the possessor is implicit the same suffix indeed aligns with English ‘his/her/its’ (ex.3). In lack of syntactic information we cannot easily distinguish among such cases. However the improvement in translation performance yielded by the removal of *P3sg* suggests that the third case has a minor impact.

- |    |                        |                   |
|----|------------------------|-------------------|
| 1) | <i>sirt çantası</i>    | ‘backpack’        |
|    | [back] [bag]+P3sg      |                   |
| 2) | <i>bayanın çantası</i> | ‘the lady’s bag’  |
|    | [lady]+Gen [bag]+P3sg  |                   |
| 3) | <i>çantası</i>         | ‘his/her/its bag’ |
|    | [bag]+P3sg             |                   |

**MS7: Cases, Poss, Copula.** This rule splits off copula from words, in addition to MS6’s rules. Unless it is implicit

<sup>2</sup>Morphological segmentation is carried out by a script that will be released soon under <http://hlt.fbk.eu/people/bisazza>.

<sup>3</sup>The omitted schemes differ by little from the ones mentioned – e.g. MS3 splits genitive instead of removing it – or are different combinations of the rules included in the schemes described in this paper – e.g. MS9 is like MS11 plus relative suffix splitting.

(ex.1), Turkish copula is expressed by suffixation (ex.2-3):

- |    |                                |                     |
|----|--------------------------------|---------------------|
| 1) | <i>bayan yorgun</i>            | ‘the lady is tired’ |
|    | [lady] [tired]                 |                     |
| 2) | <i>yorgunum</i>                | ‘I am tired’        |
|    | [tired]+DB+Verb+Zero+Pres+A1sg |                     |
| 3) | <i>yorgundum</i>               | ‘I was tired’       |
|    | [tired]+DB+Verb+Zero+Past+A1sg |                     |

**MS8: Cases, Poss, Copula, Relative.** The relative suffix *-ki* is used to form either adjectival phrases (ex.1,3) or pronominal expressions (ex.2,4) and in both cases it is not easily alignable with English:

- |    |                              |                                |
|----|------------------------------|--------------------------------|
| 1) | <i>lobideki bayan</i>        | ‘the lady (who’s) in the hall’ |
|    | [hall]+Loc+Rel [lady]        |                                |
| 2) | <i>lobideki</i>              | ‘the one (who’s) in the hall’  |
|    | [hall]+Loc+Rel               |                                |
| 3) | <i>bu sabahki gazete</i>     | ‘this morning’s paper’         |
|    | [this] [morning]+Rel [paper] |                                |
| 4) | <i>bu sabahki</i>            | ‘this morning’s one’           |
|    | [this] [morning]+Rel         |                                |

In this scheme we isolate relative suffixes *-ki* occurring after the locative case, but this rule appears to slightly worsen translation performances, therefore it is not included in the following scheme. The treatment of this suffix probably requires a more refined strategy of segmentation that takes context into account.

**MS11: Cases, Poss, Copula, Verb person.** Besides applying rules of MS7, as a first attempt to reduce data sparseness due to verb inflection, we split off person suffixes from finite verb forms and copula. The following example shows an analyzed Turkish word before and after segmentation: the number of tokens increases from 1 to 5 as the word is split into noun, possessive, instrumental case, copula and verbal person:

*arkadaşımlayım* (‘I’m with my friend’):  
arkadaş+Noun+A3sg +P1sg +Ins ^DB+Verb+Zero+Pres +A1sg

The above segmentation rules were applied on the Turkish side of the parallel corpus used to train word alignment models for SMT. An example of the effect of word segmentation scheme MS11 on the resulting (symmetrized) word alignment is shown in Figure 1. In both sentences we observe a beneficial increase of 1-1 alignments, which permits on one side to correctly map the translation phrase-pair *kız arkadaş-‘girlfriend’* (left), and on the other to capture the complex word re-ordering of the phrase *bu yere* literally meaning ‘this to-place’ (right).

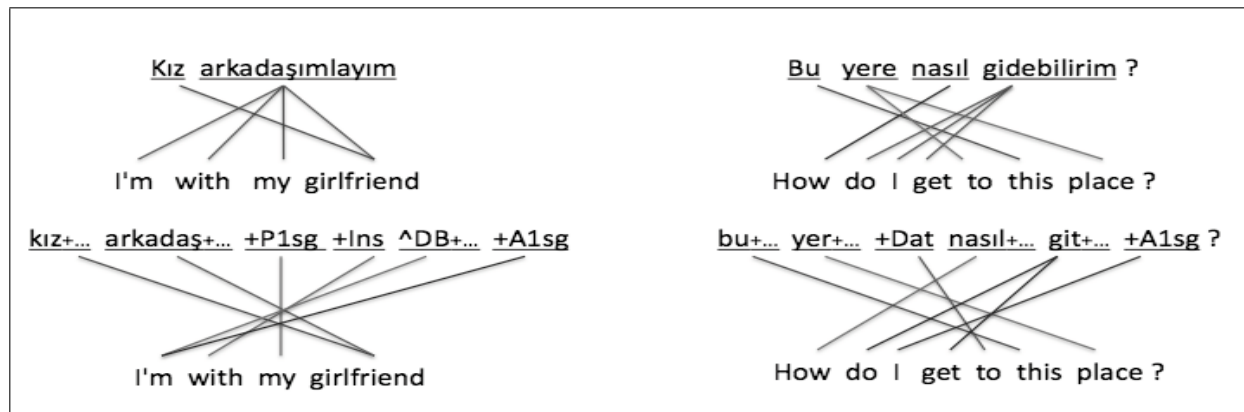


Figure 1: Two examples of sentence alignments before (up) and after (bottom) morphological segmentation MS11.

## 5. Experiments

### 5.1. Baseline

The baseline system is built upon the open-source MT toolkit Moses [9]. Phrase pairs are extracted from symmetrized word alignments generated by GIZA++ [10]. The phrase-based translation model provides direct and inverted frequency-based and lexical-based probabilities for each phrase pair included in a given phrase table. The decoder features a statistical log-linear model including a phrase-based translation model, a 5-gram language model, a lexicalized distortion model and word and phrase penalties. Distortion limit is 6 by default.

The weights of the log-linear combination are optimized by means of a minimum error training procedure [11] run on IWSLT09’s devset 1 using only the gold reference translation. Evaluation is performed on devset 2. The baseline preprocessing consists in simple tokenization and lowercasing of the source side data.

### 5.2. Morphological segmentation

Table 4 shows how morphological segmentation positively affects the training corpus dictionary size and the test OOV rate by reducing the differences in token granularity between Turkish and English: as the schemes become more complex, the number of words in the training corpus grows (from 6.9 to 8.4 words per sentence on average as opposed to 9.1 on the English side), whereas the number of different forms lowers. Thanks to our best segmentation scheme – MS11 – the OOV rate of the test set was reduced by more than half. We also observed a positive decrease of the test set’s cross-entropy<sup>4</sup>, estimated through a 5-gram language model trained on the

<sup>4</sup>Differently from perplexity, the computation of cross-entropy does not involve normalization on the number of tokens, but it gives us an estimate of the number of bits needed to encode the whole text. For this reason we chose it to compare language modeling across different segmentations schemes. A conventional dictionary upper bound size of  $10^7$  is assumed to make LMs with different OOV rates more comparable, although care must be taken in interpreting these figures.

source side of the parallel data. This can be seen as a further sign of the fact that the translation task is being better modeled.

Preprocessing	Train		Test	
	Tokens	Dictionary	OOV%	H(bits)
baseline	139,514	17,619	6.16	59,435
MS2	151,410	14,343	4.35	58,382
MS6	156,390	12,009	3.49	57,628
MS7	157,927	11,519	3.18	57,462
MS8	158,950	11,296	3.07	57,432
MS11	168,135	10,450	2.54	57,379

Table 4: Effect of preprocessing on Turkish side’s training corpus size and dictionary, test OOV and cross-entropy.

The impact of morphological preprocessing on translation performances is shown in Table 5. In each system the same preprocessing is applied to the training, development and test data. Word-error rate variations are not very significant (except for MS11), while position independent word-error rate constantly decreases (except for MS8). This suggests that morphological segmentation is improving the system’s lexical choice much more than reordering.

Preprocessing	BLEU	BP	WER	PER
baseline	52.26	95.58	37.75	29.95
MS2	53.89	96.78	37.21	28.51
MS6	54.10	98.14	37.29	28.19
MS7	55.05	98.84	37.73	27.67
MS8	54.94	98.40	37.35	27.72
MS11	56.23	98.86	36.59	26.37

Table 5: BLEU scores, brevity penalties (BP), word-error rate (WER) and position independent word-error rate (PER) in percentages on the IWSLT09 Turkish-English task.

### 5.3. Distortion limit

Keeping in mind that lexical reordering is one of the most challenging problems of Turkish to English SMT, we investigated how the distortion limit (DL) affects translation performances. It seems fair to think that since the number of words has grown, the DL should also be raised consequently. Given the short average size of IWSLT corpora sentences we decided to test translation performances in unlimited distortion conditions. Results are presented in Table 6 (note that each system was run with limited and unlimited distortion by using the set of weights optimized with DL equal to 6). As expected we found out that the gain obtained by allowing unlimited reorderings is higher when Turkish text has been morphologically segmented (by a relative improvement of 3.0% against 1.3% in the baseline), since suffix splitting makes possible new movements inside the sentence. This may as well prove that segmentation helps to establish more refined alignments. It is also interesting to notice how word-error rate improves, which was not possible with DL being set to 6. In the case of MS11 the unlimited distortion made us gain nearly 3 points of WER, that is a reduction from 36.59% to 33.70%.

Preprocess.	DL	BLEU	$\Delta$	BP	WER	PER
baseline	6	52.26	1.3%	95.58	37.75	29.95
	$\infty$	52.96		95.65	37.18	29.71
MS6	6	54.10	1.4%	98.14	37.29	28.19
	$\infty$	54.87		98.16	36.69	28.35
MS11	6	56.23	3.0%	98.86	36.59	26.37
	$\infty$	57.91		99.22	33.70	25.69

Table 6: BLEU score relative improvement with no distortion limit (DL). Additional scores reported as in Table 5.

### 5.4. Lexical approximation of OOV words

The figures of Table 4 suggest that splitting all suffixes would make the OOV rate fall close to zero. However this would not benefit the translation task itself because we would go beyond English word granularity and force the system to translate morphemes instead of words, thus making the choice of translation options and reordering far too complex.

Supposing that we reached the threshold of positive segmentation through our best scheme MS11, we tried to further reduce the OOV rate by operating on the test set. The idea consists in replacing each OOV in the decoder input by the most similar word found in the training among the words sharing the same lemma – i.e. lexical approximation. We designed a simple similarity function that gives high priority to the words sharing a large number of contiguous tags, and penalizes candidates to replacement whose tag sequence differs more from that of the original OOV word<sup>5</sup>.

<sup>5</sup>More precisely: score = match  $\times$  20 – diff<sub>1</sub>  $\times$  2 – diff<sub>2</sub>  $\times$  5, where match, diff<sub>1</sub> and diff<sub>2</sub> are respectively the numbers of shared con-

Table 7 shows a subset of candidates to the replacement of OOV word *çıkışlar* (‘exits’, ‘checkouts’) as ranked by our similarity function. The best result of lexical approximation in this case is the singular form *çıkış* (‘exit’).

Word	Gloss	Preprocessed (MS11)	Score
<i>çıkışlar</i>	exits	çık+Verb+Pos`DB+Noun+Inf3+A3pl	
<i>çıkış</i>	exit	çık+Verb+Pos`DB+Noun+Inf3+A3sg	<b>93</b>
<i>çıkma</i>	going out	çık+Verb+Pos`DB+Noun+Inf2+A3sg	66
<i>çıkacak</i>	will go out	çık+Verb+Pos`DB+Noun+FutPart+A3sg	66
<i>çıkkan</i>	who goes out	çık+Verb+Pos`DB+Adj+PresPart	44
<i>çıkıyor</i>	is going out	çık+Verb+Pos+ProgI	27
<i>çıkılmıyor</i>	isn't going out	çık+Verb+Neg+ProgI	0
<i>çıkartır</i>	takes out	çık+Verb`DB+Verb+Caus+Pos+Aor	-15

Table 7: Example of lexical approximation.

Words whose lemma was never found in the training remain OOV. Another limit of the current implementation is that the best replacer is chosen in a deterministic fashion before decoding, which raises the chances of introducing noise in the text to translate. Although this technique still needs to be improved, it made us gain another 0.2 absolute points BLEU corresponding to a reduction of the OOV rate from 2.54% to 0.89% (see Table 8).

Preprocess.	DL	BLEU%	BP	WER%	PER%
MS11	$\infty$	57.91	0.9922	33.70	25.69
MS11 & lex. approx.	$\infty$	58.12	0.9945	33.87	25.36

Table 8: Effect of lexical approximation on the IWSLT09 Turkish-English task.

## 6. Discussion and future work

The experiments have shown that selectively splitting suffixes from morphologically analyzed and disambiguated Turkish text considerably improves the performance of an SMT system. In general it seems that the improvement increases with the complexity of the segmentation scheme – i.e. treating more classes of suffixes helps. Yet not all suffix splittings benefit the translation task.

We believe that a reasonable way to establish preprocessing rules is to formulate linguistic knowledge-based hypothesis and to validate them by an experimental phase involving retraining of the translation models. This is particularly true in the context of agglutinative languages, where suffix combinatory is too wide to allow for the testing of all possible rule combinations.

It is important to notice, though, that the IWSLT task is a particular one: its small training set allows for fast retraining and testing of many different preprocessing conditions, while the short size of its sentences lets tuning and decoding costs

iguous tags, different tags in the OOV word, different tags in the replacer candidate.

be affordable even in unlimited distortion conditions. This would probably not hold in other tasks.

Furthermore, it was shown by [3] that morphological preprocessing has a positive effect on Arabic-English translation performances in scarce-resource conditions, while it can even harm in very large-resource conditions.

For these reasons we plan to repeat our experiments on another Turkish-English task, with longer and more complex sentences, and in different data size conditions.

Despite the considerable improvement yielded by our best preprocessing scheme, finding the best rule set for Turkish morphological segmentation still remains an open problem. In this work we mainly focused on nominal suffixation, but the few rules we tested on verbs suggest that these should be preprocessed better. More research is needed in this direction, as verbal suffixation is more complex than nominal, and rules harder to define. In particular we would like to tackle non-finite verb forms, that are very difficult to align with English, for example those of relative clauses.

Concerning lexical approximation, we are aware that the scoring function used to pick the best candidate for OOV word replacement could be improved. It may also be helpful to feed Moses with multiple options of replacement (e.g. through XML markup or word lattice input) so that the translation and language models would contribute to the decision at decoding time.

Finally, we showed that morphological segmentation not only decreases data sparseness, but also positively impacts on word reordering, as reported by the experiments on limited vs unlimited distortion conditions. This outcome, to our view, opens interesting research perspectives on the machine learning of robust reordering rules between Turkish and English.

## 7. Acknowledgements

Thanks to all the authors of the linguistic resources mentioned in this paper, and thanks to Deniz Yuret for providing useful links to them<sup>6</sup>. This work was supported by the EuroMatrixPlus project (IST-231720), which is funded by the European Commission under the Seventh Framework Programme for Research and Technological Development.

## 8. References

- [1] A. Göksel and C. Kerslake, “Turkish. A Comprehensive Grammar”, London and New York, Routledge, 2005, pp. 68–97, 103–108, 195–196, 284–285.
- [2] K. Oflazer, and I. Durgar El-Kahlout, “Exploring Different Representational Units in English-to-Turkish Statistical Machine Translation”, in *Proc. of the 2nd Workshop on Statistical Machine Translation*, 2007, pp. 25–32.
- [3] N. Habash, F. Sadat, “Arabic Preprocessing Schemes for Statistical Machine Translation”, in *Proc. of the Human Language Technology Conference of the NAACL*, 2006, pp. 49–52.
- [4] C. Mermer, H. Kaya and M. U. Dogan, “The TUBITAK-UEKAE Statistical Machine Translation System for IWSLT 2007”, in *Proc. of IWSLT*, 2007, pp. 176–179.
- [5] W. Shen, B. Delaney, T. Anderson and R. Slyh, “The MIT-LL/AFRL IWSLT-2008 MT System”, in *Proc. of the IWSLT*, 2008, pp. 69–76.
- [6] K. Oflazer, “Two-level description of Turkish morphology”, in *Literary and Linguistic Computing*, 1994, vol. 9, no. 2, pp. 137–148.
- [7] K. R. Beesley and L. Karttunen, “Finite State Morphology”, Palo Alto, CA: CSLI Publications, 2003.
- [8] H. Sak, T. Güngör and M. Saraçlar, “Morphological Disambiguation of Turkish Text with Perceptron Algorithm”, in *Proc. of CICLing*, 2007, pp. 107–118.
- [9] P. Koehn, et al., “Moses: Open source toolkit for statistical machine translation”, in *Proc. of the 45th Annual Meeting of the Association for Computational Linguistics. Demo and Poster Sessions*, Prague, Czech Republic, 2007, pp. 177–180.
- [10] F. J. Och and H. Ney, “A systematic comparison of various statistical alignment models”, *Computational Linguistics*, 2003, vol. 29, no. 1, pp. 19–51.
- [11] F. J. Och, “Minimum error rate training in statistical machine translation”, in *Proc. of the 41st Annual Meeting of the Association for Computational Linguistics*, 2003, pp. 160–167.

<sup>6</sup>See webpage: <http://denizyuret.blogspot.com/2006/11/turkish-resources.html>

<u>Japon Büyükelçiliği ile irtibata geçmek istiyorum .</u>
<i>Ref: I'd like to contact the Japanese Embassy .</i>
base: I'd like to contact with Japanese büyükelçiliği .
MS11: I'd like to contact with Japanese embassy .
<u>Bu film rulolarını banyo ettirip basabilir miydiniz ?</u>
<i>Ref: Could you develop and print these rolls of film ?</i>
base: Could you reissue ettirip rulolarını this film developed ?
MS11: Could you reissue roll of film developed ?
<u>Santralden santrale bir arama yapmak istiyorum .</u>
<i>Ref: I'd like to place a station-to-station call .</i>
base: santrale santralden I'd like to make a call .
MS11: Operator . I'd like to make a call to the operator
<u>Yirmi dakikada bir kalkar .</u>
<i>Ref: It leaves every twenty minutes .</i>
base: It leaves twenty minutes .
MS11: It leaves every twenty minutes .
<u>Onu bulmaktan ümidi hemen hemen kestim .</u>
<i>Ref: I've just about given up finding it .</i>
base: bulmaktan ümidi cut it right away .
MS11: I cut almost hope from find it .
<u>Takma dişlerim tam oturmuyor .</u>
<i>Ref: My dentures don't fit right .</i>
base: denture false right false .
MS11: denture My false teeth don't fit right .
<u>İskoçya'dan geliyorum .</u>
<i>Ref: I come from Scotland .</i>
base: İskoçya'dan back .
MS11: I come from Scotland .
<u>Belki bir doktora görünmelisin .</u>
<i>Ref: Perhaps you should see a doctor .</i>
base: Maybe görünmelisin a doctor .
MS11: Maybe you must see a doctor .
<u>Sığır eti harikaydı .</u>
<i>Ref: The beef was great .</i>
base: beef harikaydı .
MS11: beef was great .
<u>Bilgisayarında isminizi bulamıyorum .</u>
<i>Ref: I can't find your name on my computer .</i>
base: bilgisayarında I can't find your name .
MS11: I can't find your name in computer .
<u>Balkondan iki yer alabilir miyim ?</u>
<i>Ref: May I have two balcony seats ?</i>
base: Can I have two balcony ?
MS11: Can I have two balcony seats ?
<u>Şimdi kirazların çiçek açma mevsimi .</u>
<i>Ref: It's cherry blossom season .</i>
base: kirazların buds mail seasons now .
MS11: cherry blossoms bloom season now .

Table 9: Examples of translation outputs compared: baseline vs MS11, both with unlimited distortion.