

HOW MUCH DOES TOKENIZATION AFFECT NEURAL MACHINE TRANSLATION?

Authors: Miguel Domingo¹, Mercedes García-Martínez², Alexandre Helle², Francisco Casacuberta¹, Manuel Herranz²

¹ Pattern Recognition and Human Language Technology Research Center

Universitat Politècnica de València

{midobal, fcn}@prhlt.upv.es

² Pangeanic / B.I Europa

PangeaMT Technologies Division

{m.garcia, a.helle, m.herranz}@pangeanic.com



UNIVERSITAT
POLITÈCNICA
DE VALÈNCIA



INTRODUCTION

- In this work, we use tokenization referring to separating punctuation and splitting tokens into words or subwords.
- Tokenizing words has proven to be helpful to reduce vocabulary and increase the number of examples of each word.
- It is extremely important for languages in which there is no separation between words.
- In this study, we aim to find the impact of tokenization on the quality of the final translation produced using neural machine translation.

NEURAL MACHINE TRANSLATION

Given a source sentence $x_1^J = x_1, \dots, x_J$ of length J , NMT aims to find the best translated sentence $\hat{y}_1^{\hat{I}} = \hat{y}_1, \dots, \hat{y}_{\hat{I}}$ of length \hat{I} :

$$\hat{y}_1^{\hat{I}} = \arg \max_{I, y_1^I} Pr(y_1^I | x_1^J)$$

CORPORA

Languages: Japanese (Ja), English (En), Russian (Ru), Chinese (Zh), German (De) and Arabic (Ar).

Specific Domain	Language				
	Ja-En	Ru-En	Zh-En	De-En	Ar-En
Computer Software - Instructions for use		X			X
Medical Equipment and Supplies	X	X	X	X	X
Consumer Electronics	X		X	X	X
Industrial Electronics		X		X	
Stores and Retail Distribution	X	X	X		
Healthcare		X			

METRICS

- **BiLingual Evaluation Understudy (BLEU)**: geometric average of the modified n-grams precision, with a penalty factor for short sentences.
- **Translation Error Rate (TER)**: word edit operations normalized by the number of words in the final translation.

RESULTS

Language	SentencePiece		OpenNMT tokenizer		Moses tokenizer		Mecab		Stanford	
	BLEU	TER	BLEU	TER	BLEU	TER	BLEU	TER	BLEU	TER
Ja-En	32.0 ± 1.3	51.1 ± 1.5	29.1 ± 1.4	54.7 ± 1.4	36.3 ± 1.4	47.5 ± 1.3	36.0 ± 1.5	48.6 ± 1.4	-	-
En-Ja	26.5 ± 1.4	62.5 ± 1.9	25.0 ± 4.4	89.9 ± 4.1	33.6 ± 2.3	61.0 ± 2.5	45.8 ± 1.3	43.7 ± 1.3	-	-
Ru-En	12.9 ± 0.9	72.7 ± 1.1	11.9 ± 0.9	74.9 ± 1.3	15.3 ± 1.0	68.6 ± 1.2	-	-	-	-
En-Ru	12.2 ± 0.8	75.0 ± 1.0	11.3 ± 0.9	77.3 ± 1.1	16.3 ± 1.2	70.4 ± 1.6	-	-	-	-
Zh-En	20.5 ± 1.1	64.8 ± 1.2	23.1 ± 1.3	64.8 ± 1.3	27.5 ± 1.3	59.8 ± 1.2	-	-	26.0 ± 1.3	59.3 ± 1.2
En-Zh	17.1 ± 1.2	71.2 ± 1.2	10.4 ± 3.9	101.1 ± 3.1	21.4 ± 2.0	65.8 ± 1.7	-	-	29.9 ± 1.2	55.6 ± 1.2
De-En	21.4 ± 0.8	67.8 ± 2.1	29.6 ± 0.9	54.2 ± 0.9	30.3 ± 0.9	52.8 ± 0.9	-	-	-	-
En-De	16.1 ± 0.7	76.4 ± 2.3	22.5 ± 0.9	65.0 ± 1.5	23.6 ± 0.9	62.9 ± 1.0	-	-	-	-
Ar-En	17.9 ± 0.8	66.9 ± 1.3	14.8 ± 0.8	71.3 ± 1.1	19.1 ± 0.9	65.4 ± 1.9	-	-	-	-
En-Ar	10.1 ± 0.6	75.3 ± 1.3	9.2 ± 0.6	77.2 ± 0.9	12.4 ± 0.7	69.8 ± 0.9	-	-	-	-

EXAMPLE

Example 1	
Source	Revalidation of single-pilot single-engine class ratings
Reference	verlängerung von klassenberechtigungen für einmotorige flugzeuge mit einem piloten
SentencePiece	verlängerung der einzelantriebsklasse einmotorischer motorklasse
OpenNMT tokenizer	zur validierung der einmotorik-einzelmaschine mit einzelantrieb
Moses tokenizer	verlängerung von klassenberechtigungen für einmotorige flugzeuge mit einem piloten
Example 2	
Source	Cold drawing of wire
Reference	herstellung von kaltgezogenem draht
SentencePiece	kalt zeichnung des drahtes
OpenNMT tokenizer	kaltbildzeichnung
Moses tokenizer	herstellung von kaltgezogenem draht

TOKENIZERS

SentencePiece (Kudo and Richardson, 2018): an unsupervised text tokenizer and detokenizer mainly for Neural Network-based text generation systems.

OpenNMT tokenizer (Klein et al., 2017): the tokenizer included with the *OpenNMT* toolkit.

Moses tokenizer (Koehn et al., 2007): the tokenizer included with the *Moses* toolkit.

Mecab (Sim, 2013): an open source morphological analysis engine for Japanese, based on conditional random fields.

Stanford Word Segmenter (Tseng et al., 2005): a Chinese word segmenter based on conditional random fields.

TOKENIZATION

SentencePiece

Original: *In a browser window (Internet Explorer or Firefox) browse to www.dellconnect.com.*

Segmented: *_In _a _browser _window _ (Internet _Explorer _or _Firefox) _browse _to _www . dell connect . com .*

OpenNMT tokenizer

Original: *In a browser window (Internet Explorer or Firefox) browse to www.dellconnect.com.*

Segmented: *In a browser window (Internet Explorer or Firefox) browse to www . dellconnect . com .*

Moses tokenizer

Original: *In a browser window (Internet Explorer or Firefox) browse to www.dellconnect.com.*

Segmented: *In a browser window (Internet Explorer or Firefox) browse to www.dellconnect.com.*

Mecab

Original: ブラウザウィンドウ(Internet ExplorerまたはFirefox)で、www.dellconnect.comにアクセスします。

Segmented: ブラウザウィンドウ (Internet Explorer または Firefox) で 、 www . dellconnect . com に アクセス します 。

Stanford Word Segmenter

Original: 到 <http://www.kace.com/trial> , 然后“下 K1000 用版”, 将的 OVF (放虚化格式) 文件下到 vSphere 系。

Segmented: 到 <http://www.kace.com/trial> , 然后“下 K1000 用版”, 将的 OVF (放虚化格式) 文件下到 vSphere 系。

CONCLUSIONS

- Gains of up to 12 BLEU points and 15 TER points due to tokenization impact.
- Each tokenizer is more suitable for specific languages.
- *Moses tokenizer* seems to be the most suitable for European languages.

FUTURE WORK

- Evaluation of the corpus since impact using *SentencePiece*.
- Comparison of more segmentation strategies.
- Experimentation with generic data or bigger corpus size.

ACKNOWLEDGEMENTS

The research leading to these results has received funding from the Centro para el Desarrollo Tecnológico Industrial (CDTI) and the European Union through Programa Operativo de Crecimiento Inteligente (EXPEDIENT: IDI-20170964). We gratefully acknowledge the support of NVIDIA Corporation with the donation of a GPU used for part of this research.