

Machine Translation Using Artificial Neural Networks

Jonáš Holcner*



Abstract

The goal of this paper is to develop a neural machine translation system (NMT). That is a machine translation system based on neural networks. Specifically, the system is based on the state-of-the-art architecture, which is encoder-decoder architecture, created with recurrent neural networks enabling sequence to sequence translation.

The system is build with libraries Keras and Tensorflow and is tested against Moses statistical machine translation tool.

This work does not bring some concrete model with new state of the art results but shows some insight into the topic as well as provides an open source python library that can interested reader use to easily conduct his own experiments.

Keywords: neural machine translation — NMT — recurrent neural networks — RNN — LSTM — encoder-decoder architecture — sequence to sequence — seq2seq — keras — moses — bleu

Supplementary Material: [Downloadable Library](#)

*xholcn01@stud.fit.vutbr.cz, Faculty of Information Technology, Brno University of Technology

1. Introduction

In the recent years, there is a significant increase in usage of machine learning and artificial intelligence. This is because only lately, the capacity and performance of computers caught up with the available amount of data that is being produced every day as to build and train large enough neural networks. Now days, neural networks are widely capable of recognizing images, transcribing spoken language and most interestingly for this paper, they are quite capable in translating sequences from one language to another.

The biggest advantage of modern NMT approach is that it does not have some of the problems the traditional machine translation systems had. Instead of being composed of many different complex parts, NMT has the ability to learn and translate directly in and-to-

end fashion.

Goal of this work is to develop and try out such system and provide an out of the box usable library. Solution proposed in this paper make use of the state-of-the-art NMT architecture which is encoder-decoder. Each of these two components is one recurrent neural network together capable of directly translating whole sequences from one language to another.

The result is python package *nmt* built with Keras and Tensorflow. With this package were conducted experiments, evaluated with the standard BLEU score. Results were compared with the system produced by the Moses [1] statistical machine tool.

17

18

19

20

21

22

23

24

25

26

27

28

29

30	2. Previous Works	
31	First idea of recurrent neural networks comes from the	
32	nineties [2]. The vanilla RNN, however, had a problem	
33	with long term dependencies because of the vanishing	
34	and exploding gradient [3].	
35	Thus came improved variants of the RNN – long	
36	short term memory (LSTM) [4, 5] and its simpler ver-	
37	sion, gated recurrent unit (GRU) [6]. These units have	
38	a memory, that stores and changes information in it	
39	over time, enabling the network to remember long term	
40	dependencies.	
41	Works [7, 8, 9] shows that good performing lan-	
42	guage models are possible to build with recurrent	
43	neural networks. This lays foundation for the neu-	
44	ral machine translation as language models are the	
45	vital part. The advantage of neural language model is	
46	that it learns embeddings in a continuous space for the	
47	words, which provides the model with more context it	
48	can learn from. Different variants of learning the word	
49	embddings are shown here [10, 11, 12, 13]. Pre-	
50	trained word embeddings, for example on some very	
51	large data set, can be used to boost performance of a	
52	NMT system, which would have to otherwise learn	
53	those embeddings by itself.	
54	Encoder-decoder architecture was proposed in [14]	
55	and was used for rescoring hypotheses produced by	
56	a phrase-based system with successful improvement.	
57	[15] then shows how to use encoder-decoder architec-	
58	ture for direct sequence to sequence translation and	
59	comes with the best results at the time. Furthermore,	
60	they found out the importance of reversing order of the	
61	words in all source sentences (reverse encoder), that	
62	improves models performance, by introducing short	
63	term dependencies between the source and the target	
64	sentence.	
65	Upon this builds [16] which shows even better re-	
66	sults with bi-directional encoder. What is even more	
67	important, they address the problem of encoder-decoder	
68	approach, where the meaning of the translated sen-	
69	tence is captured in a fixed-length vector and that can	
70	be problematic for translating long sentences. The pro-	
71	posed remedy is so called <i>attention</i> mechanism which	
72	lets the model at the time of decoding, look at the	
73	most important words from the source sentence for	
74	the currently translated word, resulting in even better	
75	performance.	
76	As translation is an open-vocabulary problem, the	
77	NMT systems have to somehow handle the vocabular-	
78	ies. This was typically done by using out-of-vocabulary	
79	words and by using very large vocabularies, which	
80	cases the models to be very memory and performance	
81	demanding. [17, 18] shows that using sub-word units	
	can be more efficient, help with rare and unknown	82
	words and improve the results.	83
	Current state-of-the-art results are published by	84
	Google [19, 20], which uses all of the techniques de-	85
	scribed, showing that they can be successfully applied	86
	on large production data sets.	87
	Another thing Google shows, is that with no changes	88
	to the model architecture, one model can be used to	89
	learn to translate from and to more languages [21],	90
	even to produce translations between languages, that it	91
	was not explicitly trained on (zero-shot translations).	92
	3. Seq2seq translation with encoder-	
	decoder architecture	93
	For a deeper overview of NMT based systems, I would	94
	point the reader to [22].	95
	4. Implementation of the NMT system	96
	datasets preprocessing bucketing fit generator? bidirec-	97
	tional encoder subwords - BPE beam search shuffling	98
	main.py repo	99
	5. Experiments and evaluation	100
	6. Conclusions	101
	[Paper Summary] What was the paper about, then?	102
	What the reader needs to remember about it?	103
	[Highlights of Results] Exact numbers. Remind	104
	the reader that the paper matters.	105
	[Paper Contributions] What is the original con-	106
	tribution of this work? Two or three thoughts that one	107
	should definitely take home.	108
	[Future Work] How can other researchers / devel-	109
	opers make use of the results of this work? Do you	110
	have further plans with this work? Or anybody else?	111
	Acknowledgements	112
	I would like to thank my supervisor X. Y. for his help.	113
	References	114
	[1] Philipp Koehn, Hieu Hoang, Alexandra Birch,	115
	Chris Callison-Burch, Marcello Federico, Nicola	116
	Bertoldi, Brooke Cowan, Wade Shen, Chris-	117
	tine Moran, Richard Zens, Chris Dyer, Ondřej	118
	Bojar, Alexandra Constantin, and Evan Herbst.	119
	Moses: Open source toolkit for statistical ma-	120
	chine translation. In <i>Proceedings of the 45th</i>	121
	<i>Annual Meeting of the ACL on Interactive Poster</i>	122
	<i>and Demonstration Sessions</i> , ACL '07, pages	123
	177–180, Stroudsburg, PA, USA, 2007. Associa-	124
	tion for Computational Linguistics.	125

- [2] Jeffrey L. Elman. Finding structure in time. *Cognitive Science*, 14(2):179–211, 1990.
- [3] Y. Bengio, P. Simard, and P. Frasconi. Learning long-term dependencies with gradient descent is difficult. *Trans. Neur. Netw.*, 5(2):157–166, March 1994.
- [4] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Comput.*, 9(8):1735–1780, November 1997.
- [5] Felix A. Gers, Jürgen A. Schmidhuber, and Fred A. Cummins. Learning to forget: Continual prediction with lstm. *Neural Comput.*, 12(10):2451–2471, October 2000.
- [6] Junyoung Chung, Çağlar Gülçehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *CoRR*, abs/1412.3555, 2014.
- [7] Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Janvin. A neural probabilistic language model. *J. Mach. Learn. Res.*, 3:1137–1155, March 2003.
- [8] Frederic Morin and Yoshua Bengio. Hierarchical probabilistic neural network language model. In Robert G. Cowell and Zoubin Ghahramani, editors, *Proceedings of the Tenth International Workshop on Artificial Intelligence and Statistics*, pages 246–252. Society for Artificial Intelligence and Statistics, 2005.
- [9] Tomáš Mikolov. *Statistické jazykové modely založené na neuronových sítích*. PhD thesis, Vysoké učení technické v Brně, Fakulta informačních technologií, 2012.
- [10] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *CoRR*, abs/1301.3781, 2013.
- [11] Tomas Mikolov, Wen-tau Yih, and Geoffrey Zweig. Linguistic regularities in continuous space word representations. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 746–751. Association for Computational Linguistics, 2013.
- [12] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. Glove: Global vectors for word representation. In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, 2014.
- [13] Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. Enriching word vectors with subword information. *CoRR*, abs/1607.04606, 2016.
- [14] Kyunghyun Cho, Bart van Merriënboer, Çağlar Gülçehre, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder-decoder for statistical machine translation. *CoRR*, abs/1406.1078, 2014.
- [15] Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. Sequence to sequence learning with neural networks. *CoRR*, abs/1409.3215, 2014.
- [16] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *CoRR*, abs/1409.0473, 2014.
- [17] Tomas Mikolov, Ilya Sutskever, Anoop Deoras, Hai-Son Le, Stefan Kombrink, and Jan Cernocký. Subword language modeling with neural networks. In *Subword Language Modeling with Neural Networks*, 2011.
- [18] Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with subword units. *CoRR*, abs/1508.07909, 2015.
- [19] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Lukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, and Jeffrey Dean. Google’s neural machine translation system: Bridging the gap between human and machine translation. *CoRR*, abs/1609.08144, 2016.
- [20] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *CoRR*, abs/1706.03762, 2017.
- [21] Melvin Johnson, Mike Schuster, Quoc V. Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat, Fernanda B. Viégas, Martin Wattenberg, Greg Corrado, Macduff Hughes, and Jeffrey Dean. Google’s multilingual neural machine translation system: Enabling zero-shot translation. *CoRR*, abs/1611.04558, 2016.
- [22] Graham Neubig. Neural machine translation and sequence-to-sequence models: A tutorial. *CoRR*, abs/1703.01619, 2017.