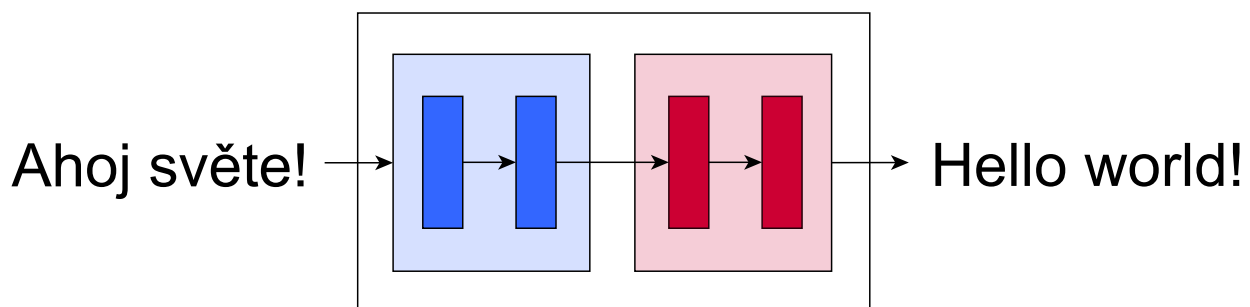


Machine Translation Using Artificial Neural Networks

Jonáš Holcner*



Abstract

The aim of this work 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 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 tradi-

tional machine translation systems had. Instead of being composed of many different complex parts, NMT has the ability to learn and translate directly in end-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 encoder-decoder architecture. 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

14
15
16
17
18
19
20
21
22
23
24
25
26

experiments, evaluated with the standard BLEU score. Results were compared with the system produced by the Moses [1] statistical machine tool.

2. Previous Works

First idea of recurrent neural networks comes from the nineties [2]. The vanilla RNN, however, had a problem with long term dependencies because of the vanishing and exploding gradient [3].

Thus came improved variants of the RNN – long short term memory (LSTM) [4, 5] and its simpler version, gated recurrent unit (GRU) [6]. These units have a memory, that stores and changes information in it over time, enabling the network to remember long term dependencies.

Works [7, 8, 9] shows that good performing language models are possible to build with recurrent neural networks. This lays foundation for the neural machine translation as language models are the vital part. The advantage of neural language model is that it learns embeddings in a continuous space for the words, which provides the model with more context it can learn from. Different variants of learning the word embeddings are shown here [10, 11, 12, 13]. Pre-trained word embeddings, for example on some very large data set, can be used to boost performance of a NMT system, which would have to otherwise learn those embeddings by itself.

Encoder-decoder architecture was proposed in [14] and was used for rescoring hypotheses produced by a phrase-based system with successful improvement. [15] then shows how to use encoder-decoder architecture for direct sequence to sequence translation and comes with the best results at the time. Furthermore, they found out the importance of reversing order of the words in all source sentences (reverse encoder), that improves models performance, by introducing short term dependencies between the source and the target sentence.

Upon this builds [16] which shows even better results with bi-directional encoder. What is even more important, they address the problem of encoder-decoder approach, where the meaning of the translated sentence is captured in a fixed-length vector and that can be problematic for translating long sentences. The proposed remedy is so called *attention* mechanism which lets the model, at the time of decoding, look at the most important words from the source sentence for the currently translated word, resulting in even better performance.

As translation is an open-vocabulary problem, the NMT systems have to somehow handle the words not

present at the time of training. This was typically done by using out-of-vocabulary tokens and by using very large vocabularies, which causes the models to be very memory and performance demanding. Solution for this can be sub-word units [17, 18], that are shown to be more efficient, help with rare and unknown words and improve the results.

Current state-of-the-art results are published by Google [19, 20], which uses all of the techniques described, showing that they can be successfully applied on large production data sets.

Another thing Google shows, is that with no changes to the model architecture, one model can be used to learn to translate from and to more languages [21], even to produce translations between languages, that it was not explicitly trained on (zero-shot translations).

3. Seq2seq translation with encoder-decoder architecture

Sequence to sequence (**seq2seq**) means that the process of translating takes one whole sentence at the input, in a source language, and its meaning translates into sequence in a target language. Its modeling probability of a target sentence based on source sentence as shown in figure 1.

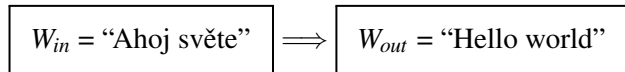


Figure 1. Seq2seq models a probability $P(W_{out}|W_{in})$. Meaning that the model learns to predict sentence W_{out} based on a sentence W_{in} , so it learns to translate.

Best neural network architecture for this modeling is **encoder-decoder**. Both encoder and decoder are recurrent neural network, usually one or more layers of LSTM or GRU.

Encoder takes embeddings of the input tokens (words or subwords) and processes the input sequence in source language into fixed-length vector, so called “thought” vector – it captures the essence of the given sentence.

Decoders hidden state is initialized with this output from encoder. The decoding process is started with start token, which is given to decoder as first input. Decoder then generates sequence in target language until it reaches end token that tells it to stop. There is a difference between time when model is being trained and between time when it is used to predict. In the training time, decoder is fed the correct, expected output. This is called “teacher forcing” [22]. In the inference time, decoder is fed its own output each time step. Figure 2 shows the encoder-decoder architecture.

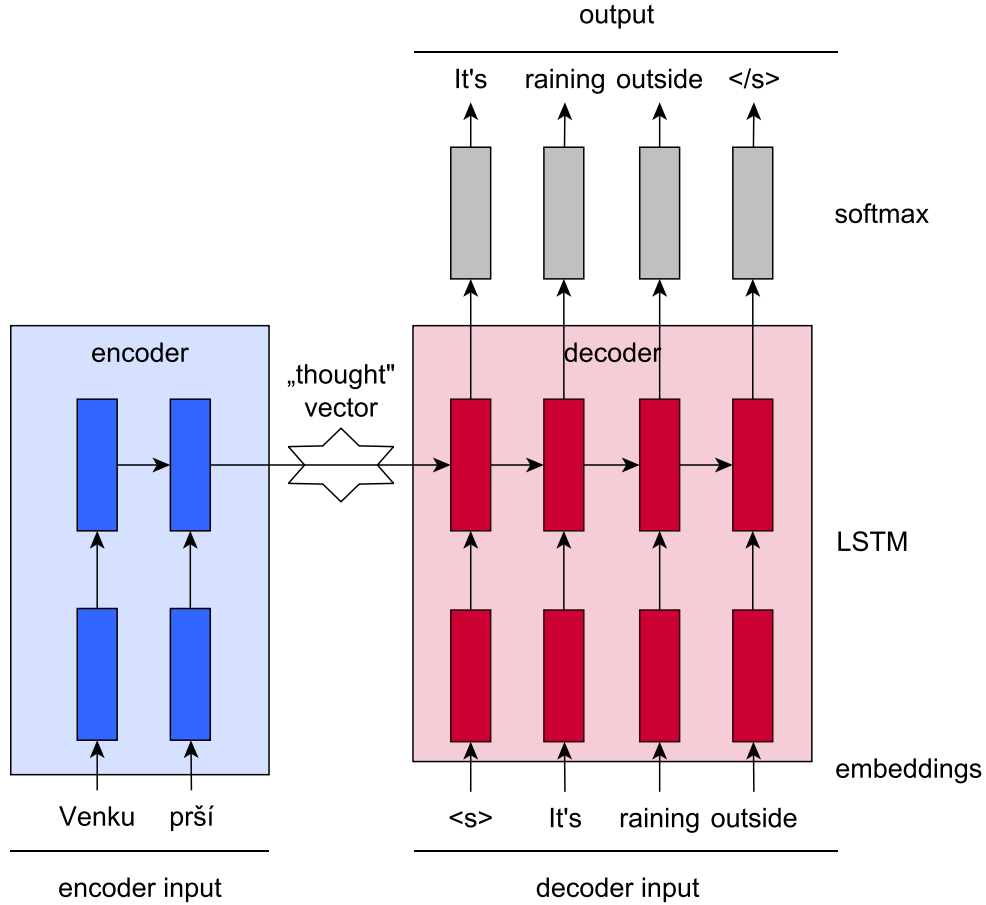


Figure 2. Encoder-decoder architecture shown on translation from Czech sentence to its English equivalent. Encoder processes embeddings of input sequence and produces fixed-length “thought” vector. This vector is used as initial state of decoder, it tells it from what context should it produce output in target language. Prediction is started with the $\langle s \rangle$ starting token. Then the decoder is fed either correct output tokens during training time or its own output, from time $t - 1$, during inference time, until it generates the ending $\langle /s \rangle$ token.

Equations for the encoder-decoder:

$$m_t^{(f)} = M_{f_t}^{(f)} \quad (1)$$

$$h_t^f = \begin{cases} RNN^{(f)}(m_t^{(f)}, h_{t-1}^{(f)}) & \text{if } t \geq 1, \\ 0 & \text{otherwise.} \end{cases} \quad (2)$$

$$m_t^{(e)} = M_{e_{t-1}}^{(e)} \quad (3)$$

$$h_t^e = \begin{cases} RNN^{(e)}(m_t^{(e)}, h_{t-1}^{(e)}) & \text{if } t \geq 1, \\ h_{|F|}^f & \text{otherwise.} \end{cases} \quad (4)$$

$$p_t^{(e)} = \text{softmax}(W_{hs}h_t^{(e)} + b_s) \quad (5)$$

- $RNN^{(f)}$ - encoder, $RNN^{(e)}$ - decoder
- $M_{f_t}^{(f)}$ - embedding for encoder input in time t
- $M_{e_t}^{(e)}$ - embedding for decoder input in time t
- $h_t^{(f)}$ - encoder hidden state in time t
- $h_t^{(e)}$ - decoder hidden state in time t
- W_{hs} - weights
- b_s - bias

Embedding is looked up for every word in time t (eq. 1). Then the hidden state of encoder is computed (eq. 2). After processing the whole input sequence, there should be enough context stored for decoder initialisation. Then, as for encoder, embedding is looked up for decoder input (eq. 3). Only now the word is from time $t - 1$ as decoder generates new word based on the last one. In the time t_0 decoder is fed starting symbol $\langle s \rangle$. Eq. 4 computes hidden state of decoder and shows that in the time t_0 , decoder is initialised with encoder final state. Final output probability is resolved with *softmax* function (eq. 5).

For a deeper overview of NMT based systems, I would point the reader to [23].

4. Implementation of the NMT system

datasets preprocessing bucketing fit generator? bidirectional encoder subwords - BPE beam search shuffling main.py repo

147 5. Experiments and evaluation

148 6. Conclusions

149 **[Paper Summary]** What was the paper about, then?
150 What the reader needs to remember about it?

151 **[Highlights of Results]** Exact numbers. Remind
152 the reader that the paper matters.

153 **[Paper Contributions]** What is the original con-
154 tribution of this work? Two or three thoughts that one
155 should definitely take home.

156 **[Future Work]** How can other researchers / devel-
157 opers make use of the results of this work? Do you
158 have further plans with this work? Or anybody else?

159 Acknowledgements

160 I would like to thank my supervisor Ing. Igor Szőke, Ph.D.
161 for his help.

162 References

- 163 [1] Philipp Koehn, Hieu Hoang, Alexandra Birch,
164 Chris Callison-Burch, Marcello Federico, Nicola
165 Bertoldi, Brooke Cowan, Wade Shen, Chris-
166 tine Moran, Richard Zens, Chris Dyer, Ondřej
167 Bojar, Alexandra Constantin, and Evan Herbst.
168 Moses: Open source toolkit for statistical ma-
169 chine translation. In *Proceedings of the 45th*
170 *Annual Meeting of the ACL on Interactive Poster*
171 *and Demonstration Sessions*, ACL '07, pages
172 177–180, Stroudsburg, PA, USA, 2007. Associa-
173 tion for Computational Linguistics.
- 174 [2] Jeffrey L. Elman. Finding structure in time. *Cog-*
175 *nitive Science*, 14(2):179–211, 1990.
- 176 [3] Y. Bengio, P. Simard, and P. Frasconi. Learning
177 long-term dependencies with gradient descent
178 is difficult. *Trans. Neur. Netw.*, 5(2):157–166,
179 March 1994.
- 180 [4] Sepp Hochreiter and Jürgen Schmidhuber. Long
181 short-term memory. *Neural Comput.*, 9(8):1735–
182 1780, November 1997.
- 183 [5] Felix A. Gers, Jürgen A. Schmidhuber, and
184 Fred A. Cummins. Learning to forget: Con-
185 tinual prediction with lstm. *Neural Comput.*,
186 12(10):2451–2471, October 2000.
- 187 [6] Junyoung Chung, Çağlar Gülçehre, KyungHyun
188 Cho, and Yoshua Bengio. Empirical evaluation
189 of gated recurrent neural networks on sequence
190 modeling. *CoRR*, abs/1412.3555, 2014.
- 191 [7] Yoshua Bengio, Réjean Ducharme, Pascal Vin-
192 cent, and Christian Janvin. A neural probabilistic

language model. *J. Mach. Learn. Res.*, 3:1137–
1155, March 2003. 193 194

- [8] Frederic Morin and Yoshua Bengio. Hierarchi- 195
cal probabilistic neural network language model. 196
In Robert G. Cowell and Zoubin Ghahramani,
editors, *Proceedings of the Tenth International* 197
Workshop on Artificial Intelligence and Statistics, 198
pages 246–252. Society for Artificial Intelligence
and Statistics, 2005. 200 201
- [9] Tomáš Mikolov. *Statistické jazykové modely* 202
založené na neuronových sítích. PhD thesis,
Vysoké učení technické v Brně, Fakulta in- 203
formačních technologií, 2012. 204 205
- [10] Tomas Mikolov, Kai Chen, Greg Corrado, and 206
Jeffrey Dean. Efficient estimation of word repre- 207
sentations in vector space. *CoRR*, abs/1301.3781,
2013. 208 209
- [11] Tomas Mikolov, Wen-tau Yih, and Geoffrey 210
Zweig. Linguistic regularities in continuous
space word representations. In *Proceedings of the* 211
2013 Conference of the North American Chapter 212
of the Association for Computational Linguistics: 213
Human Language Technologies, pages 746–751. 214
Association for Computational Linguistics, 2013. 215 216
- [12] Jeffrey Pennington, Richard Socher, and Christo- 217
pher D. Manning. Glove: Global vectors for
word representation. In *Empirical Methods in* 218
Natural Language Processing (EMNLP), pages
1532–1543, 2014. 219 220 221
- [13] Piotr Bojanowski, Edouard Grave, Armand 222
Joulin, and Tomas Mikolov. Enriching word
vectors with subword information. *CoRR*,
abs/1607.04606, 2016. 223 224 225
- [14] Kyunghyun Cho, Bart van Merriënboer, Çağlar 226
Gülçehre, Fethi Bougares, Holger Schwenk, and
Yoshua Bengio. Learning phrase representations
using RNN encoder-decoder for statistical ma- 227
chine translation. *CoRR*, abs/1406.1078, 2014. 228 229 230
- [15] Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. 231
Sequence to sequence learning with neural net-
works. *CoRR*, abs/1409.3215, 2014. 232 233
- [16] Dzmitry Bahdanau, Kyunghyun Cho, and 234
Yoshua Bengio. Neural machine translation by
jointly learning to align and translate. *CoRR*,
abs/1409.0473, 2014. 235 236 237
- [17] Tomas Mikolov, Ilya Sutskever, Anoop Deoras, 238
Hai-Son Le, Stefan Kombrink, and Jan Cernocký.
Subword language modeling with neural net- 239
works. In *Subword Language Modeling with* 240
Neural Networks, 2011. 241 242

- 243 [18] Rico Sennrich, Barry Haddow, and Alexandra
244 Birch. Neural machine translation of rare words
245 with subword units. *CoRR*, abs/1508.07909,
246 2015.
- 247 [19] Yonghui Wu, Mike Schuster, Zhifeng Chen,
248 Quoc V. Le, Mohammad Norouzi, Wolfgang
249 Macherey, Maxim Krikun, Yuan Cao, Qin Gao,
250 Klaus Macherey, Jeff Klingner, Apurva Shah,
251 Melvin Johnson, Xiaobing Liu, Lukasz Kaiser,
252 Stephan Gouws, Yoshikiyo Kato, Taku Kudo,
253 Hideto Kazawa, Keith Stevens, George Kurian,
254 Nishant Patil, Wei Wang, Cliff Young, Jason
255 Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals,
256 Greg Corrado, Macduff Hughes, and Jeffrey
257 Dean. Google’s neural machine translation sys-
258 tem: Bridging the gap between human and ma-
259 chine translation. *CoRR*, abs/1609.08144, 2016.
- 260 [20] Ashish Vaswani, Noam Shazeer, Niki Parmar,
261 Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,
262 Lukasz Kaiser, and Illia Polosukhin. Attention is
263 all you need. *CoRR*, abs/1706.03762, 2017.
- 264 [21] Melvin Johnson, Mike Schuster, Quoc V. Le,
265 Maxim Krikun, Yonghui Wu, Zhifeng Chen,
266 Nikhil Thorat, Fernanda B. Viégas, Martin Wat-
267 tenberg, Greg Corrado, Macduff Hughes, and Jef-
268 frey Dean. Google’s multilingual neural machine
269 translation system: Enabling zero-shot transla-
270 tion. *CoRR*, abs/1611.04558, 2016.
- 271 [22] Anirudh Goyal, Alex Lamb, Ying Zhang,
272 Saizheng Zhang, Aaron Courville, and Yoshua
273 Bengio. Professor forcing: A new algorithm for
274 training recurrent networks. pages 4601–4609.
275 2016.
- 276 [23] Graham Neubig. Neural machine translation and
277 sequence-to-sequence models: A tutorial. *CoRR*,
278 abs/1703.01619, 2017.