

Neural Network Models For Translating Input Sequences To Output Sequences

HETA DESAI

Abstract

Deep Neural Networks (DNNs) have shown strong performance on complex learning tasks, especially when large labeled datasets are available. However, they struggle with mapping input sequences directly to output sequences. In this paper, the authors propose a general, end-to-end method for sequence learning that relies on minimal assumptions about the structure of sequences. Their approach uses a multi-layer Long Short-Term Memory (LSTM) network to encode an input sequence into a fixed-size vector, followed by another deep LSTM network that decodes this vector into an output sequence.

A key result of the study is that this LSTM-based model achieved a BLEU score of 34.8 on the full test set for English-to-French translation using the WMT'14 dataset, despite being penalized for generating out-of-vocabulary words. Notably, the model performed well even on long sentences. For comparison, a phrase-based Statistical Machine Translation (SMT) system scored 33.3 on the same dataset.

Further improvement was observed when the LSTM was used to rerank the top 1000 translation hypotheses from the SMT system, raising the BLEU score to 36.5—close to the best result reported at the time. The LSTM also learned meaningful representations of phrases and sentences that were sensitive to word order and robust to changes like converting between active and passive voice.

Interestingly, performance improved significantly when the order of words in the source sentences was reversed (while keeping target sentences in normal order). This created more short-term dependencies between source and target sequences, which helped make the training process more effective.

1. Introduction

Deep Neural Networks (DNNs) are highly effective machine learning models that have demonstrated strong results on complex tasks like speech and visual object recognition. Their strength lies in their ability to carry out complex parallel computations with a limited number of steps. For example, it has been shown that a neural network with just two hidden layers of quadratic size can sort N N -bit numbers—highlighting the expressive power of DNNs. Though DNNs are related to traditional statistical models, they are capable of learning complex computations. When trained with supervised backpropagation and provided with sufficient labeled data, DNNs can learn parameter settings that allow them to excel at tasks that are relatively easy for humans to perform.

However, a major limitation of DNNs is that they typically require input and output to be fixed-length vectors. This restricts their application to tasks where the structure can be compactly represented in this way. Many real-world problems—like machine translation, speech recognition, and question answering—involve sequences of varying lengths, making it difficult to apply traditional DNN architectures directly. As such, a general-purpose method that can learn to map sequences of arbitrary length to other sequences would be highly beneficial.

Sequences are challenging for standard DNNs because they demand a flexible input-output dimensionality. This paper presents a straightforward yet powerful approach to handling such tasks using the Long Short-Term Memory (LSTM) architecture. The approach involves using one LSTM network to read and encode the input sequence into a fixed-size vector, and another LSTM to decode the vector back into an output sequence. The second LSTM function is similar to a language model but is conditioned on the encoded input.

This method is especially suited for problems with long-ranged dependencies between inputs and outputs, thanks to LSTM's ability to remember information over extended time steps. Prior work has explored similar goals. For example, Kalchbrenner and Blunsom were among the first to encode an input sentence into a fixed vector, and Cho et al. applied similar ideas for re-ranking translation hypotheses. Additionally, Graves proposed a differentiable attention mechanism that inspired more advanced translation systems like that of Bahdanau et al. Another technique, Connectionist Sequence Classification (CTC), has also been used for sequence learning, though it relies on a monotonic alignment between input and output sequences.

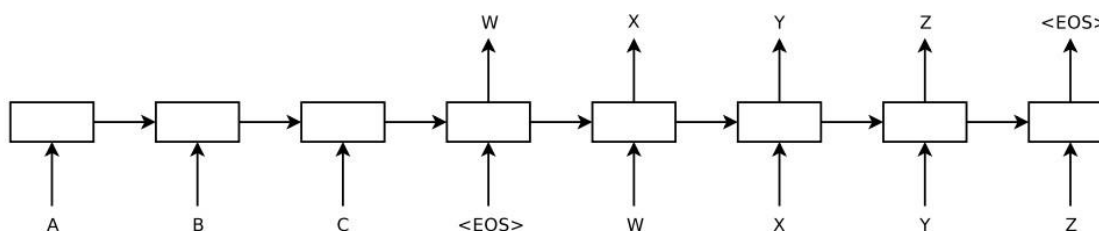


Figure 1: Our model reads an input sentence “ABC” and produces “WXYZ” as the output sentence. The model stops making predictions after outputting the end-of-sentence token. Note that the LSTM reads the input sentence in reverse, because doing so introduces many short term dependencies in the data that make the optimization problem much easier.

The key finding of this study is that, on the WMT’14 English-to-French translation task, the model achieved a BLEU score of **34.81** using an ensemble of five deep LSTM networks (each with 384 million parameters and 8,000-dimensional states). Translations were generated using a straightforward left-to-right beam search decoder. This score represents the highest result for direct translation using large neural networks to date. For comparison, a standard phrase-based Statistical Machine Translation (SMT) system scored **33.30** on the same task. Notably, the LSTM model used a fixed vocabulary of 80,000 words, so any word not in this vocabulary led to a penalty in the BLEU score. Despite this limitation, the LSTM outperformed the SMT system, indicating that even a relatively basic neural architecture with room for improvement can surpass traditional systems.

Additionally, the LSTM was used to rescore the 1,000-best hypotheses generated by the SMT system, pushing the BLEU score up to **36.5**—an improvement of 3.2 points over the baseline and nearly matching the previous best published score of **37.0**.

Interestingly, the LSTM handled **long sentences** well, even though similar architectures often struggle with this. The success is largely attributed to a simple but effective strategy: **reversing the order of words** in the source sentences during training and testing (while keeping the target sentences unchanged). This introduced short-term dependencies that made learning easier and allowed the LSTM to optimize more effectively using stochastic gradient descent (SGD). This reversal trick is considered a key technical innovation of the paper.

Another strength of the LSTM model is its ability to represent variable-length input sentences as **fixed-size vectors**. Since translation requires capturing the **meaning** of a sentence, the model is encouraged to produce embeddings where semantically similar sentences are close together in this learned space. Qualitative analysis shows that the model is sensitive to **word order** and robust to variations in sentence structure, such as **active vs. passive voice**.

2. The model

The Recurrent Neural Network (RNN) [31, 28] is a natural generalization of feedforward neural networks to sequences. Given a sequence of inputs (x_1, \dots, x_T) , a standard RNN computes a sequence of outputs (y_1, \dots, y_T) by iterating the following equation:

$$\begin{aligned} h_t &= \text{sigm}(W^{hx}x_t + W^{hh}h_{t-1}) \\ y_t &= W^{yh}h_t \end{aligned}$$

Recurrent Neural Networks (RNNs) are well-suited for mapping sequences to sequences when the alignment between the input and output is already known. However, they become challenging to apply when the input and output sequences are of different lengths and the relationship between them is complex or non-monotonic.

A straightforward approach for handling general sequence-to-sequence problems is to use one RNN to encode the input sequence into a fixed-size vector, and then use a second RNN to decode that vector into the output sequence—an idea also explored by Cho et al. [5]. While this setup theoretically provides the necessary information for the model to perform the task, it is difficult to train due to the long-term dependencies involved (see figure 1) [14, 4, 16, 15].

However, Long Short-Term Memory (LSTM) networks [16] are specifically designed to handle long-range temporal dependencies, making them a promising solution for this problem.

The main objective of the LSTM in this context is to model the conditional probability $p(y_1, \dots, y_T | x_1, \dots, x_T) p(y_1, \dots, y_{T'} | x_1, \dots, x_T) p(y_1, \dots, y_T | x_1, \dots, x_T)$, where (x_1, \dots, x_T) is the input sequence and (y_1, \dots, y_T) is the target sequence, which may differ in length. The process involves two steps: first, the LSTM encodes the input sequence into a fixed-length vector v , derived from its final hidden state. Then, a second LSTM acts as a language model, using v as its initial hidden state to generate the probability distribution over the output sequence y_1, \dots, y_T .

$$p(y_1, \dots, y_{T'} | x_1, \dots, x_T) = \prod_{t=1}^{T'} p(y_t | v, y_1, \dots, y_{t-1})$$

In this approach, each probability $p(y_t | v, y_1, \dots, y_{t-1})$ is calculated using a softmax function applied over the entire vocabulary. The LSTM architecture follows the design from Graves [10]. Importantly, each sentence must end with a special end-of-sentence token (“<EOS>”), allowing the model to define a probability distribution over sequences of any length. Figure 1 illustrates this idea, where the LSTM processes the sequence “A”, “B”, “C”, “<EOS>” to generate a fixed representation, which is then used to predict the sequence “W”, “X”, “Y”, “Z”, “<EOS>”.

However, the actual implementation of our model differs from this basic setup in three key ways:

- 1. Separate LSTMs for Encoding and Decoding:** We use one LSTM for encoding the input sequence and a separate LSTM for decoding the output. This not only increases the number of model parameters without adding significant computational cost, but also facilitates training on multiple language pairs at the same time [18].
- 2. Deeper Architectures Perform Better:** We discovered that deep LSTMs outperform shallow ones. As a result, our model uses an LSTM with four layers.
- 3. Input Sentence Reversal:** A crucial improvement was reversing the word order of the input sentence. For instance, instead of training the model to map “a, b, c” to its translation “ α, β, γ ”, we train it to map “c, b, a” to “ α, β, γ ”. This reversal aligns the beginning of the input with the beginning of the output (e.g., “a” is close to “ α ”), which helps the model learn more efficiently. This simple change significantly boosted LSTM performance.

3. Experiments

We evaluated our approach on the WMT'14 English-to-French machine translation task using two main strategies. First, we used our model to translate input sentences directly, without relying on any traditional statistical machine translation (SMT) system. Second, we used our model to rescore the n-best candidate translations produced by a baseline SMT model. We report the performance of both methods, provide example translations, and explore the learned sentence representations.

3.1 Dataset Details

For training, we used the WMT'14 English-to-French dataset. Our training subset included 12 million sentence pairs, with approximately 348 million French words and 304 million English words. This dataset was chosen because it's a cleaned and preprocessed version made available by [29], which also provides tokenized training and test data along with 1000-best translation hypotheses from the SMT baseline.

Since neural language models rely on word embeddings, we used fixed-size vocabularies for each language. Specifically, we selected the 160,000 most frequent words for English (the source language) and the 80,000 most frequent words for French (the target language). Words not found in these vocabularies were replaced with a special "UNK" (unknown) token.

3.2 Decoding and Rescoring

The main focus of our experiments was training a large, deep LSTM using numerous sentence pairs. The training goal was to maximize the log-likelihood of producing the correct target sentence T given the source sentence S . In other words, the model was optimized to increase the probability of generating accurate translations.

$$\frac{1}{|S|} \sum_{(T,S) \in \mathcal{S}} \log p(T|S)$$

where S is the training set. Once training is complete, we produce translations by finding the most likely translation according to the LSTM:

$$\hat{T} = \arg \max_T p(T|S)$$

To find the most likely translation, we use a straightforward left-to-right beam search decoder. This decoder keeps track of a limited number BBB of partial hypotheses at each step, where each partial hypothesis is a prefix of a potential translation. At every timestep, each hypothesis in the beam is extended by appending every possible word from the vocabulary. Since this generates a large number of new hypotheses, we retain only the top BBB hypotheses based on their log probability scores. Once a hypothesis includes the special end-of-sentence token

$\langle \text{EOS} \rangle$, it is removed from the beam and added to the set of complete translations. Although this method is not exact, it is easy to implement. Notably, even with a beam size of 1, our model performs well, and most of the beam search's benefits are retained with a beam size of just 2 (as shown in Table 1).

We also applied our LSTM to rescore the 1000-best candidate translations produced by the baseline SMT system [29]. For each hypothesis in the list, we calculated its log-likelihood using our LSTM and averaged this score with the one from the SMT system.

3.3 Reversing Source Sentences

Even though LSTMs are designed to handle long-range dependencies, we observed

significantly better training results when we reversed the order of words in the source sentences (while keeping target sentences unchanged). With this reversal, test perplexity improved from 5.8 to 4.7, and the BLEU score rose from 25.9 to 30.6.

Although we don't have a full explanation for why this works, we suspect it's because reversing the source sequence creates more short-term dependencies in the data. Typically, words in the source sentence are positioned far from their translated counterparts in the target sentence, leading to a high "minimal time lag" [17]. Reversing the source order doesn't change the average word-to-word distance, but it reduces the time lag for the initial parts of the sentences. This makes it easier for the model to learn connections between source and target sequences during backpropagation, which likely accounts for the notable improvement in translation quality.

At first, we assumed that reversing the input sentences would only improve the model's confidence for the initial part of the translated output, possibly weakening its predictions toward the end. However, we found that LSTMs trained on reversed input sequences performed significantly better on long sentences than those trained on the original order (see Section 3.7). This suggests that reversing the input helps the LSTM make more efficient use of its memory.

3.4 Training Details

Training the LSTM model turned out to be relatively straightforward. We used deep LSTMs with 4 layers, each layer containing 1,000 LSTM units and 1,000-dimensional word embeddings. The input vocabulary had 160,000 words, while the output vocabulary had 80,000. This means each sentence was represented using 8,000 values. Deep LSTMs showed a clear advantage over shallow ones, with each additional layer reducing perplexity by nearly 10%—likely due to the increased capacity of the deeper models. A standard softmax was used over the 80,000-word output vocabulary. The entire model had 384 million parameters, including 64 million dedicated to recurrent connections—32 million for the encoder and 32 million for the decoder.

Key training settings included:

- Initializing all LSTM parameters using a uniform distribution between -0.08 and 0.08.
- Using stochastic gradient descent (SGD) without momentum, starting with a fixed learning rate of 0.7. After 5 epochs, we halved the learning rate every half-epoch, training for a total of 7.5 epochs.
- Training was done in batches of 128 sequences, and the gradient was scaled accordingly.
- Although LSTMs typically avoid vanishing gradients, they can suffer from exploding gradients. To manage this, we applied gradient clipping: if the L2 norm of the gradient $s = \|g\|_2$ exceeded 5, we scaled it as $s > 5$, we set $g = 5g/s$.
- Since sentence lengths vary—most are between 20 to 30 words, while some exceed 100—random minibatches tended to be inefficient. To solve this, we grouped sentences of similar length together in each batch, which doubled training speed.

3.5 Parallelization

Initially, our C++ implementation of the deep LSTM (with the configuration described above) could process around 1,700 words per second on a single GPU. This wasn't fast enough, so we parallelized the model across an 8-GPU system. Each of the four LSTM layers was run on a separate GPU, and activations were passed to the next layer immediately after computation.

The remaining four GPUs handled the softmax computation, with each GPU responsible for multiplying a $1000 \times 20,000$ matrix. This optimized setup achieved a processing speed of 6,300 words per second (including both English and French words) with a minibatch size of 128.

Overall, training the model took about ten days using this parallelized setup.

3.6 Experimental Results

We assessed the translation quality using the **cased BLEU score** [24]. This was calculated with the `multi-bleu.pl` script on tokenized predicted translations and their reference counterparts. This evaluation method aligns with those used in [5] and [2], and accurately reproduces the 33.3 BLEU score reported in [29]. Interestingly, when we apply the same method to evaluate the best-performing system from WMT'14 [9] (whose output is available at statmt.org/matrix), we obtain a BLEU score of **37.0**, which is higher than the **35.8** reported on the website.

The main results are shown in Tables 1 and 2. The best performance came from an **ensemble of LSTM models**—each initialized differently and trained with different random minibatch orderings. Although this ensemble's direct translation results do not surpass the top WMT'14 system, it marks a significant achievement: this is the **first time a fully neural machine translation system has outperformed a phrase-based SMT baseline** on a large-scale machine translation task by a notable margin. This was achieved **despite limitations**, such as not being able to translate out-of-vocabulary (OOV) words.

Method	test BLEU score (ntst14)
Bahdanau et al. [2]	28.45
Baseline System [29]	33.30
Single forward LSTM, beam size 12	26.17
Single reversed LSTM, beam size 12	30.59
Ensemble of 5 reversed LSTMs, beam size 1	33.00
Ensemble of 2 reversed LSTMs, beam size 12	33.27
Ensemble of 5 reversed LSTMs, beam size 2	34.50
Ensemble of 5 reversed LSTMs, beam size 12	34.81

Table 1: The performance of the LSTM on WMT'14 English to French test set (ntst14). Note that an ensemble of 5 LSTMs with a beam of size 2 is cheaper than of a single LSTM with a beam of size 12.

Method	test BLEU score (ntst14)
Baseline System [29]	33.30
Cho et al. [5]	34.54
Best WMT'14 result [9]	37.0
Rescoring the baseline 1000-best with a single forward LSTM	35.61
Rescoring the baseline 1000-best with a single reversed LSTM	35.85
Rescoring the baseline 1000-best with an ensemble of 5 reversed LSTMs	36.5
Oracle Rescoring of the Baseline 1000-best lists	~45

Table 2: Methods that use neural networks together with an SMT system on the WMT'14 English to French test set (ntst14).

When the LSTM model is used to **rescore** the 1000-best list generated by the SMT baseline, it achieves a BLEU score within **0.5 points** of the top WMT'14 system—highlighting its strong performance.

3.7 Performance on Long Sentences

To our surprise, the LSTM model performed well even on longer sentences, as demonstrated by the quantitative results in Figure 3. Table 3 highlights some examples of long input sentences along with their corresponding translations.

3.8 Model Analysis

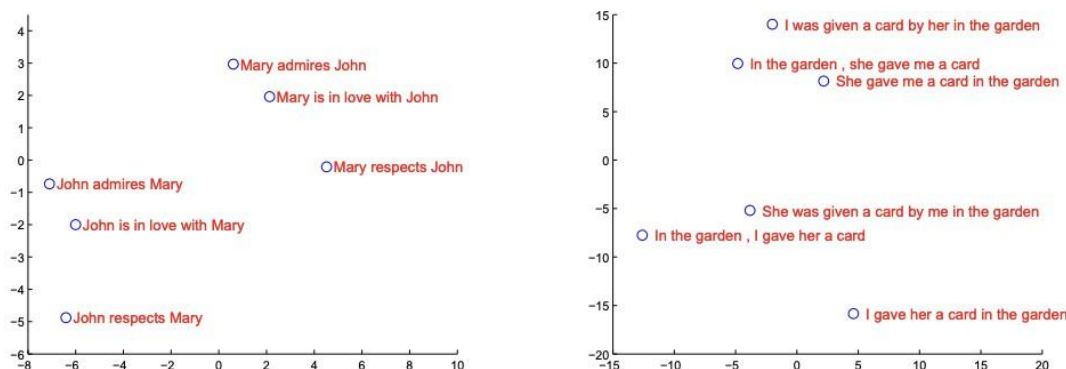


Figure 2: The figure shows a 2-dimensional PCA projection of the LSTM hidden states that are obtained after processing the phrases in the figures. The phrases are clustered by meaning, which in these examples is primarily a function of word order, which would be difficult to capture with a bag-of-words model. Notice that both clusters have similar internal structure.

A notable strength of our model is its capability to convert a sequence of words into a fixed-size vector. Figure 2 illustrates several of these learned representations. The visualization reveals that the model captures the importance of word order, yet remains relatively unaffected by switching between active and passive voice. These two-dimensional views were generated using Principal Component Analysis (PCA).

Type	Sentence
Our model	Ulrich UNK , membre du conseil d' administration du constructeur automobile Audi , affirme qu' il s' agit d' une pratique courante depuis des années pour que les téléphones portables puissent être collectés avant les réunions du conseil d' administration afin qu' ils ne soient pas utilisés comme appareils d' écoute à distance .
Truth	Ulrich Hackenberg , membre du conseil d' administration du constructeur automobile Audi , déclare que la collecte des téléphones portables avant les réunions du conseil , afin qu' ils ne puissent pas être utilisés comme appareils d' écoute à distance , est une pratique courante depuis des années .
Our model	“ Les téléphones cellulaires , qui sont vraiment une question , non seulement parce qu' ils pourraient potentiellement causer des interférences avec les appareils de navigation , mais nous savons , selon la FCC , qu' ils pourraient interférer avec les tours de téléphone cellulaire lorsqu' ils sont dans l' air ” , dit UNK .
Truth	“ Les téléphones portables sont véritablement un problème , non seulement parce qu' ils pourraient éventuellement créer des interférences avec les instruments de navigation , mais parce que nous savons , d' après la FCC , qu' ils pourraient perturber les antennes-relais de téléphonie mobile s' ils sont utilisés à bord ” , a déclaré Rosenker .
Our model	Avec la crémation , il y a un “ sentiment de violence contre le corps d' un être cher ” , qui sera “ réduit à une pile de cendres ” en très peu de temps au lieu d' un processus de décomposition “ qui accompagnera les étapes du deuil ” .
Truth	Il y a , avec la crémation , “ une violence faite au corps aimé ” , qui va être “ réduit à un tas de cendres ” en très peu de temps , et non après un processus de décomposition , qui “ accompagnerait les phases du deuil ” .

Table 3: A few examples of long translations produced by the LSTM alongside the ground truth translations. The reader can verify that the translations are sensible using Google translate.

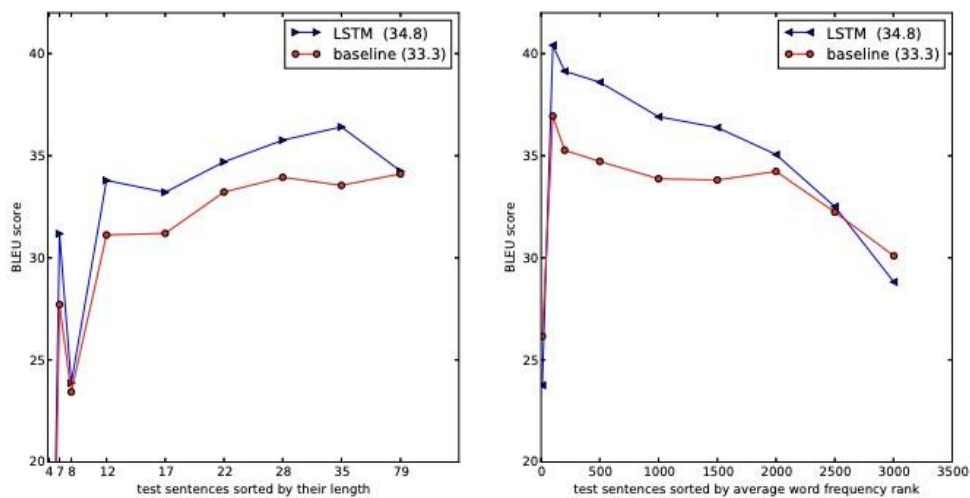


Figure 3: The left plot shows the performance of our system as a function of sentence length, where the x-axis corresponds to the test sentences sorted by their length and is marked by the actual sequence lengths. There is no degradation on sentences with less than 35 words, there is only a minor degradation on the longest sentences. The right plot shows the LSTM's performance on sentences with progressively more rare words, where the x-axis corresponds to the test sentences sorted by their "average word frequency rank".

4. Related Work

There has been substantial research on using neural networks for machine translation. Traditionally, the most straightforward and effective method has involved using RNN-based Language Models (RNNLMs) [23] or Feedforward Neural Network Language Models (NNLMs) [3] to score the then-best output of a strong statistical machine translation (SMT) system [22], which has consistently led to improvements in translation quality.

More recently, efforts have been made to integrate source language information directly into these neural models. For instance, Auli et al. [1] enhanced an NNLM by combining it with a topic model of the source sentence, leading to better rescoring performance. Devlin et al. [8] also incorporated an NNLM into the MT decoder, using alignment data to identify and feed the most relevant source words into the neural model—an approach that yielded substantial performance gains.

Our approach closely aligns with the work of Kalchbrenner and Blunsom [18], who first proposed converting input sentences into fixed-length vectors and then decoding them back into output sentences. However, they used convolutional neural networks (CNNs), which tend to lose word order information. Similarly, Cho et al. [5] employed an LSTM-style RNN to perform sentence-to-vector and vector-to-sentence transformations, although their main focus was on enhancing SMT systems.

Bahdanau et al. [2] also pursued direct neural translation, introducing an attention mechanism to improve performance on longer sentences—an issue observed by Cho et al. [5]. Pouget-Abadie et al. [26] tackled the same challenge by translating sentence segments to produce smoother outputs, somewhat resembling a phrase-based model. We believe that training on reversed source sentences could provide them with similar benefits.

End-to-end training is also explored by Hermann et al. [12], who used feedforward networks to map inputs and outputs to nearby points in a vector space. However, their system cannot generate translations independently—it requires either a lookup from a pre-existing sentence database or the rescoring of candidate sentences.

5. Conclusion

In this study, we demonstrated that a large deep LSTM model—with a limited vocabulary and minimal assumptions about the structure of the task—can outperform a traditional SMT-based system, even one with an unlimited vocabulary, on a large-scale machine translation task. The strong performance of our simple LSTM-based methods suggests it could be effective for a wide range of sequence learning tasks, given sufficient training data.

Interestingly, we found that reversing the words in the source sentences significantly improved performance. This led us to conclude that encoding problems to increase the number of short-term dependencies can simplify the learning process. Notably, while we struggled to train a standard RNN on unreversed input (as shown in Figure 1), we believe that reversing the source sentences would make such training feasible, though we did not test this experimentally.

We were also surprised by the LSTM's capability to accurately translate long sentences. Initially, we assumed that the model's limited memory would hinder its performance on longer inputs—especially since similar models have shown poor results on such tasks [5, 2, 26]. However, our LSTM, when trained on the reversed dataset, handled long sentences effectively.

Most importantly, we showed that a straightforward, relatively unrefined approach could surpass an SMT system. This indicates that with further development, even better translation results are achievable. Overall, our findings suggest that this method holds strong potential for tackling other complex sequence-to-sequence learning problems.

6. Acknowledgments

We are grateful to Samy Bengio, Jeff Dean, Matthieu Devin, Geoffrey Hinton, Nal Kalchbrenner, Thang Luong, Wolfgang Macherey, Rajat Monga, Vincent Vanhoucke, Peng Xu, Wojciech Zaremba, and the Google Brain team for their valuable feedback and insightful discussions.

References

- [1] M. Auli, M. Galley, C. Quirk, and G. Zweig. Joint language and translation modeling with recurrent neural networks. In EMNLP, 2013.
- [2] D. Bahdanau, K. Cho, and Y. Bengio. Neural machine translation by jointly learning to align and translate. arXiv preprint arXiv:1409.0473, 2014.
- [3] Y. Bengio, R. Ducharme, P. Vincent, and C. Jauvin. A neural probabilistic language model. In Journal of Machine Learning Research, pages 1137–1155, 2003.
- [4] Y. Bengio, P. Simard, and P. Frasconi. Learning long-term dependencies with gradient descent is difficult. IEEE Transactions on Neural Networks, 5(2):157–166, 1994.
- [5] K. Cho, B. Merriënboer, C. Gulcehre, F. Bougares, H. Schwenk, and Y. Bengio. Learning phrase representations using RNN encoder-decoder for statistical machine translation. In Arxiv preprint arXiv:1406.1078, 2014.
- [6] D. Ciresan, U. Meier, and J. Schmidhuber. Multi-column deep neural networks for image classification. In CVPR, 2012.
- [7] G. E. Dahl, D. Yu, L. Deng, and A. Acero. Context-dependent pre-trained deep neural networks for large vocabulary speech recognition. IEEE Transactions on Audio, Speech, and Language Processing-Special Issue on Deep Learning for Speech and Language Processing, 2012.
- [8] J. Devlin, R. Zbib, Z. Huang, T. Lamar, R. Schwartz, and J. Makhoul. Fast and robust neural network joint models for statistical machine translation. In ACL, 2014.
- [9] Nadir Durrani, Barry Haddow, Philipp Koehn, and Kenneth Heafield. Edinburgh’s phrase-based machine translation systems for wmt-14. In WMT, 2014.
- [10] A. Graves. Generating sequences with recurrent neural networks. In Arxiv preprint arXiv:1308.0850, 2013.
- [11] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber. Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks. In ICML, 2006.
- [12] K. M. Hermann and P. Blunsom. Multilingual distributed representations without word alignment. In ICLR, 2014.
- [13] G. Hinton, L. Deng, D. Yu, G. Dahl, A. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, T. Sainath, and B. Kingsbury. Deep neural networks for acoustic modeling in speech recognition. IEEE Signal Processing Magazine, 2012.
- [14] S. Hochreiter. Untersuchungen zu dynamischen neuronalen netzen. Master’s thesis, Institut für Informatik, Technische Universität München, 1991.
- [15] S. Hochreiter, Y. Bengio, P. Frasconi, and J. Schmidhuber. Gradient flow in recurrent nets: the difficulty of learning long-term dependencies, 2001.
- [16] S. Hochreiter and J. Schmidhuber. Long short-term memory. Neural Computation, 1997.
- [17] S. Hochreiter and J. Schmidhuber. LSTM can solve hard long time lag problems. 1997.
- [18] N. Kalchbrenner and P. Blunsom. Recurrent continuous translation models. In EMNLP, 2013.
- [19] A. Krizhevsky, I. Sutskever, and G. E. Hinton. ImageNet classification with deep convolutional neural networks. In NIPS, 2012.
- [20] Q. V. Le, M. A. Ranzato, R. Monga, M. Devin, K. Chen, G. S. Corrado, J. Dean, and A. Y. Ng. Building high-level features using large scale unsupervised learning. In ICML, 2012.
- [21] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. Proceedings of the IEEE, 1998.
- [22] T. Mikolov. Statistical Language Models based on Neural Networks. PhD thesis, Brno University of Technology, 2012.
- [23] T. Mikolov, M. Karafiát, L. Burget, J. Cernocký, and S. Khudanpur. Recurrent neural network based language model. In INTERSPEECH, pages 1045–1048, 2010.
- [24] K. Papineni, S. Roukos, T. Ward, and W. J. Zhu. BLEU: a method for automatic evaluation of machine translation. In

ACL, 2002.

[25] R.Pascanu,T.Mikolov,andY.Bengio.Onthedifficultyoftrainingrecurrentneural networks. arXiv preprint arXiv:1211.5063, 2012.

[26] J. Pouget-Abadie, D. Bahdanau, B. van Merriënboer, K. Cho, and Y. Bengio. Overcoming the curse of sentence length for neural machine translation using automatic segmentation. arXiv preprint arXiv:1409.1257, 2014.

[27] A.Razborov.Onsmalldepththresholdcircuits.InProc.3rdScandinavianWorkshopon Algorithm Theory, 1992.

[28] D.Rumelhart,G.E.Hinton,andR.J.Williams.Learningrepresentationsby back-propagating errors. Nature, 323(6088):533–536, 1986.

[29] H.Schwenk.Universitylemans.http://www-lium.univ-lemans.fr/~schwenk/cslm_joint_paper/, 2014. [Online; accessed 03-September-2014].

[30] M.Sundermeyer,R.Schluter,andH.Ney.LSTMneuralnetworksforlanguage modeling.In INTERSPEECH, 2010.

[31] P.Werbos.Backpropagationthrough time: what it does and how to do it. Proceedings of IEEE, 1990.

