

# Review of Cho et al. (2014)

Matthew Evans

April 5, 2025

## Overview

Semantic translation (rather than literal or “word-for-word” translation) between languages often yields differences in phrase length between the source and target language. Cho et al.[1] propose a novel approach to machine translation that leverages a recurrent neural network (RNN) to learn a mapping between variable-length input and output sequences.

## Approach

### RNN Encoder-Decoder

The authors propose an RNN encoder-decoder architecture designed to map variable-length input sequences to fixed-length context vectors, which are then used to produce variable-length output sequences. This model scores input-output pairs, such as source-target language phrases, and can generate target phrases from source phrases.

The model comprises two RNNs: an encoder that transforms a variable-length input into a fixed-length context vector  $c$ , and a decoder that generates a variable-length output conditioned on  $c$ . By modeling the conditional probability  $p(y|x)$ , the architecture accommodates sequences of differing lengths. The decoder’s predictions rely on the previously generated token, the current hidden state, and the context vector  $c$ . Both components are trained jointly using gradient-based optimization to maximize conditional likelihood, enabling the model to generate target sequences or evaluate input-output pairs based on learned probabilities.

### RNNs and LSTM

Recurrent neural networks learn sequential data using hidden units, a special type of node which receives its previous value as input (in addition to feature data input) at each time step. RNNs were improved upon by long short-term memory nodes (LSTM)[2] which introduced “gates” that control the degree to which a hidden state is updated on a given time-step, as well as the degree to which the hidden state influences the other nodes in the network.

### The GRU

The authors introduce the *gated recurrent unit* (i.e., GRU, a term later coined in [3]), a simplified hidden unit based on LSTM featuring *update* and *reset* gates. At each time-step, the hidden state is updated from two different sources: the hidden state at the previous time-step  $h_{t-1}$ , and an alternative “candidate” value,  $\tilde{h}_t$ .

It is the update gate which controls the proportional influence these two sources have on updating  $h_t$ . The reset gate is concerned with computing the value of the candidate  $\tilde{h}_t$ , specifically the degree to which it is influenced by the previous hidden state,  $h_{t-1}$ . Thus, for a given hidden unit, the update is given by

$$h_t = zh_{t-1} + (1 - z)\tilde{h}_t,$$

The update gate  $z$  is given by

$$z_t = \sigma(W_z x + U_z h_{t-1}).$$

The candidate, which is controlled by the reset gate  $r$  is given by

$$\tilde{h}_t = \tanh(Wx + U(r \odot h_{t-1})),$$

where the vector  $r \in (0,1)$  “throttles” the influence of the the previous hidden unit vector.  $W$ ,  $W_z$ ,  $U$ , and  $U_z$  are weight matrices,  $x$  are the input vector, and  $\odot$  is the element-wise product.

In this way, hidden units can learn to selectively filter the influence of input data received at future time steps.

## Enhancing SMT

*Statistical machine translation* (SMT) systems use a log-linear model to combine various features, including phrase translation probabilities and language model scores, to find the best translation. Neural networks enhance this process by providing more accurate scoring of phrase pairs or translation hypotheses based on learned representations, ultimately improving translation quality.

The authors applied their RNN Encoder-Decoder to rescore phrase pairs in an existing SMT system, focusing on capturing linguistic patterns rather than frequency-based statistics. By ignoring phrase frequencies during training, the model learned to distinguish plausible translations. The resulting phrase scores can thus be added as features to enhance SMT translation quality.

## Measures of Success

The authors present a number of experiments examining both the qualitative and quantitative performance of their architecture.

### Quantitative

The authors evaluated their proposed RNN Encoder-Decoder on the WMT’14<sup>1</sup> English-French translation task, using a filtered subset of a large bilingual corpus to build a phrase-based SMT baseline system. The authors then trained their RNN Encoder-Decoder model and a Continuous Space Language Model (CSLM) — a feedforward neural net language model. When used to score partial translations during decoding, the baseline SMT saw improved BLEU scores by both the CSLM and RNN models.

Furthermore, combining both the CSLM and the RNN Encoder-Decoder yielded the best BLEU scores, indicating their effects are complementary rather than redundant. Attempts to further improve performance by penalizing unknown words helped on the development set but not the test set. Overall, the study shows that integrating multiple neural models into SMT systems can yield additive performance gains.

### Qualitative

The authors compared the RNN Encoder-Decoder’s phrase scores with traditional translation probability driven models, which tend to favor frequent phrases due to their reliance on corpus statistics. In contrast, the RNN Encoder-Decoder, trained without frequency data, focuses on linguistic regularities, often producing more literal or appropriate translations, especially for long or rare source phrases.

The RNN Encoder-Decoder model showed preference for shorter target phrases and at times diverged significantly from traditional models, further highlighting its independence from frequency bias. Notably, the RNN can generate fluent, well-formed target phrases that are not in the phrase table, suggesting it could potentially replace or enhance parts of the traditional SMT phrase table in future systems.

## Considerations

The following considerations are observed regarding the authors’ approach.

- By using a fixed size context vector between the variable length input and output, the model may experience information loss, particularly when handling long or context rich sequences.

---

<sup>1</sup>Workshop on Statistical Machine Translation 2014, a task for evaluating machine translation systems.

- The model, as trained, only handles a 15,000 word vocabulary, mapping all other words to an “unknown” token. This limits its usefulness in domain-specific terminology or other rare words.
- While the model’s ability to translate sequences in instances of low word frequency is impressive, it also dismisses the very real importance that word frequency can encode, potentially producing suboptimal results.

## Impact

This paper had a significant impact on the development of machine translation and sequence modeling. It introduced the Gated Recurrent Unit, a simpler and computationally efficient alternative to LSTM that uses reset and update gates. It also pioneered the RNN Encoder-Decoder architecture, enabling mappings between variable-length input and output sequences and laying the foundation for future sequence-to-sequence models. By training on unique phrase pairs and ignoring frequency, the model demonstrated the ability to learn linguistic regularities beyond surface statistics (i.e., to learn semantics), influencing work on meaning-based and low-resource translation. The authors also showed that neural phrase scoring could enhance statistical machine translation systems when combined with traditional features, proving that neural and statistical models can be complementary. Collectively, these contributions helped motivate the shift toward fully end-to-end neural machine translation, directly influencing subsequent models like *seq2seq* with attention and the Transformer.

## References

- [1] Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation, 2014.
- [2] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.
- [3] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Y. Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. 12 2014.