

# Response to Bahdanau, Cho, Bengio (2016)

Matthew Evans

April 7, 2025

## Overview

Translating between languages is challenging because words do not align one-to-one, and their order often varies. Building on earlier work by Cho et. al (2014)[1], the Bahdanau et. al propose a new approach[2] that avoids summarizing an entire sentence with one fixed snapshot, and instead creates a flexible summary for each translated word. This design allows the model to scan the entire original sentence and focus on the most relevant parts when generating each word. In doing so, the approach aims to capture the true meaning of the input even when word order differs significantly between languages.

## Approach

### Prior Work

Earlier work in neural machine translation, particularly the RNN Encoder-Decoder model introduced by Cho et al. (2014), laid the foundation for generating translations from a source sentence. In this framework, an encoder RNN processes the input sentence and compresses it into a fixed-length context vector, denoted by  $\mathbf{c}$ , which is intended to capture all the necessary information. The decoder RNN then uses this single vector to generate the output sentence, word by word.

$$\mathbf{c} = \tanh(\mathbf{W}\mathbf{h}^{(T)})$$

Here,  $\mathbf{h}^{(T)}$  is the hidden state of the encoder RNN and  $\mathbf{W}$  is a learned weight matrix. Thus, the fixed-length context vector  $\mathbf{c}$  acts as a lossy compression of the entire input sequence, forcing the model into a “hard alignment” that obscures intricate, variable word-to-word correspondences, which becomes especially problematic for longer sentences with critical nuanced dependencies.

### Novelty

The novel approach introduced by the authors centers on an attention-based mechanism that computes a variable-length context for each target word, rather than compressing the entire source sentence into a single fixed-length vector. For instance, the context vector for the  $i$ -th target word is defined as

$$\mathbf{c}_i = \sum_{j=1}^{T_x} \alpha_{ij} \mathbf{h}_j,$$

where the alignment weights  $\alpha_{ij}$  are determined via a differentiable scoring function. This mechanism enables the model to focus on different terms of the source sentence (indexed by  $j = 1, \dots, T_x$ ) as needed while keeping that every layer of the network is fully differentiable and trainable.

Furthermore, the model employs a bidirectional RNN architecture that processes the input sentence in both forward and backward directions. By combining these two perspectives, the model overcomes the challenges posed by differing grammatical sequences in languages. Unlike earlier methods that used neural components as auxiliary features in statistical machine translation, this approach delivers a fully trainable translation system, promising more coherent and contextually accurate translations.

## Considerations

Despite the promising aspects of the proposed model, several potential risks and weaknesses deserve attention.

- The authors note that due to the tendency of RNNs to better represent recent inputs, the annotation  $\mathbf{h}_j$  is biased toward the words surrounding  $x_j$ . This reliance on local context may be problematic in capturing long-range dependencies, and its universal applicability remains uncertain.
- The model employs a vocabulary that is twice the size of that used in previous work by [1]. This expansion may have contributed significantly to its performance gains on a training set of merely 15,000 words, raising the question of whether the observed improvements are primarily due to architectural innovations or simply a result of increased lexical coverage.
- The authors acknowledge that handling unknown or rare words continues to be a challenge. While the model outperforms state-of-the-art systems such as Moses on sentences composed solely of known vocabulary, Moses maintains an edge when no such restrictions are applied. This indicates that further enhancements are necessary to address rare word occurrences, which is crucial for the model to achieve consistently high performance in more diverse and realistic translation contexts.

## Strengths

The proposed approach offers several notable advantages.

- It eliminates the need for a separate monolingual dataset, as required by systems like Moses, thereby streamlining the training process.
- By employing a variable-length context vector, the model avoids the lossy semantic compression inherent in fixed-length representations, allowing it to capture a more faithful summary of the source sentence.
- The integration of an attention mechanism enables soft-alignment, which permits the model to selectively focus on the most relevant input words for each output word, irrespective of their position in the sentence. This targeted focus not only enhances translation accuracy but also accommodates the natural variability in word order across different languages.

## Measures of Success

The authors validated their findings using quantitative evaluations based on BLEU scores. Their experiments demonstrated that the proposed model outperformed the earlier RNN Encoder-Decoder model described in [1] across several test sets. In particular, when translations were evaluated on sentences composed entirely of known vocabulary, the model not only surpassed its predecessor but also achieved performance levels comparable to, or even exceeding, those of the state-of-the-art phrase-based system, Moses. However, as mentioned, this success was contingent on filtering out sentences with unknown words, which highlights a limitation in handling rare or unseen vocabulary.

## Impact

These innovations revolutionized sequence modeling and NLP by removing the need for monolingual datasets and enabling variable-length context vectors, improving translation of long sentences and overall performance.

The self-attention mechanism of the Transformer architecture [3] was built on the idea of dynamically focusing on relevant parts of the input—a concept first popularized by Bahdanau et al.[2]

Beyond translation, the success of attention has paved the way for major pre-trained language models like BERT[4] and the GPT series[5].

In summary, these advancements have streamlined training processes, enhanced the capacity to capture complex linguistic structures, and significantly improved performance across numerous applications in NLP and beyond.

## 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] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate, 2016.
- [3] 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.
- [4] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. *CoRR*, abs/1810.04805, 2018.
- [5] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training, 2018. Technical Report, OpenAI.