

# LONG SHORT-TERM MEMORY AS A DYNAMICALLY COMPUTED ELEMENT-WISE WEIGHTED SUM

Anonymous authors  
Paper under double-blind review

## ABSTRACT

Long short-term memory networks (LSTMs) were introduced to combat vanishing gradients in simple recurrent neural networks (S-RNNs) by augmenting them with additive recurrent connections controlled by gates. We present an alternate view to explain the success of LSTMs: the gates themselves are powerful recurrent models that provide more representational power than previously appreciated. We do this by showing that the LSTM’s gates can be decoupled from the embedded S-RNN, producing a restricted class of RNNs where the main recurrence computes an element-wise weighted sum of context-independent functions of the inputs. Experiments on a range of challenging NLP problems demonstrate that the simplified gate-based models work substantially better than S-RNNs, and often just as well as the original LSTMs, strongly suggesting that the gates are doing much more in practice than just alleviating vanishing gradients.

## 1 INTRODUCTION

Long short-term memory networks (LSTM) (Hochreiter & Schmidhuber, 1997) have become the de-facto recurrent neural network (RNN) for learning representations of sequences in many research areas, including natural language processing (NLP). Like simple recurrent neural networks (S-RNNs) (Elman, 1990), LSTMs are able to learn non-linear functions of arbitrary-length input sequences. However, they also introduce an additional memory cell to mitigate the vanishing gradient problem (Hochreiter, 1991; Bengio et al., 1994). This memory is controlled by a mechanism of gates, whose additive connections allow long-distance dependencies to be learned more easily during backpropagation. While this view is mathematically accurate, in this paper we argue that it does not provide a complete picture of why LSTMs work in practice.

We present an alternate view to explain the success of LSTMs: the gates themselves are powerful recurrent models that provide more representational power than previously appreciated. To demonstrate this, we first show that LSTMs can be seen as a combination of two recurrent models: (1) an S-RNN, and (2) an element-wise weighted sum of the S-RNN’s outputs over time, which is implicitly computed by the gates. We hypothesize that, for many practical NLP problems, the weighted sum serves as the main modeling component. The S-RNN, while theoretically expressive, is in practice only a minor contributor that clouds the mathematical clarity of the model. By replacing the S-RNN with a context-*independent* function of the input, we arrive at a much more restricted class of RNNs, where the main recurrence is via the element-wise weighted sums that the gates are computing.

We test our hypothesis on NLP problems, where LSTMs are wildly popular at least in part due to their ability to model crucial language phenomena such as word order (Adi et al., 2017), syntactic structure (Linzen et al., 2016), and even long-range semantic dependencies (He et al., 2017). We consider four challenging tasks: language modeling, question answering, dependency parsing, and machine translation. Experiments show that while removing the gates from an LSTM can severely hurt performance, replacing the S-RNN with a simple linear transformation of the input results in minimal or no loss in model performance. We further show that in many cases, LSTMs can be further simplified by removing the output gate, arriving at an even more transparent architecture, where the output is a context-*independent* function of the weighted sum. Together, these results suggest that the gates’ ability to compute an element-wise weighted sum, rather than the non-linear transition dynamics of S-RNNs, are the driving force behind LSTM’s success.

## 2 THE MEMORY CELL COMPUTES AN ELEMENT-WISE WEIGHTED SUM

LSTMs are typically motivated as an augmentation of simple RNNs (S-RNNs), defined as follows:

$$\mathbf{h}_t = \tanh(\mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{W}_{hx}\mathbf{x}_t + \mathbf{b}_h) \quad (1)$$

S-RNNs suffer from the vanishing gradient problem (Hochreiter, 1991; Bengio et al., 1994) due to compounding multiplicative updates of the hidden state. By introducing a memory cell and an output layer that are controlled by a set of gates, LSTMs enable shortcuts through which gradients can flow easily when learning with backpropagation. This mechanism enables learning of long-distance dependencies while preserving the expressive power of recurrent non-linear transformations provided by S-RNNs.

Rather than viewing the gates as simply an auxiliary mechanism to address a *learning* problem, we present an alternate view that emphasizes their *modeling* strengths. We argue that the LSTM should be interpreted as a hybrid of two distinct recurrent architectures: (1) the S-RNN which provides multiplicative connections across timesteps, and (2) the memory cell which provides additive connections across timesteps. On top of these recurrences, an output layer is included that simply squashes and filters the memory cell at each step.

Throughout this paper, let  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  be the sequence of input vectors,  $\{\mathbf{h}_1, \dots, \mathbf{h}_n\}$  be the sequence of output vectors, and  $\{\mathbf{c}_1, \dots, \mathbf{c}_n\}$  be the memory cell's states. Then, given the basic LSTM definition below, we can formally identify three sub-components.

$$\tilde{\mathbf{c}}_t = \tanh(\mathbf{W}_{ch}\mathbf{h}_{t-1} + \mathbf{W}_{cx}\mathbf{x}_t + \mathbf{b}_c) \quad (2)$$

$$\mathbf{i}_t = \sigma(\mathbf{W}_{ih}\mathbf{h}_{t-1} + \mathbf{W}_{ix}\mathbf{x}_t + \mathbf{b}_i) \quad (3)$$

$$\mathbf{f}_t = \sigma(\mathbf{W}_{fh}\mathbf{h}_{t-1} + \mathbf{W}_{fx}\mathbf{x}_t + \mathbf{b}_f) \quad (4)$$

$$\mathbf{c}_t = \mathbf{i}_t \circ \tilde{\mathbf{c}}_t + \mathbf{f}_t \circ \mathbf{c}_{t-1} \quad (5)$$

$$\mathbf{o}_t = \sigma(\mathbf{W}_{oh}\mathbf{h}_{t-1} + \mathbf{W}_{ox}\mathbf{x}_t + \mathbf{b}_o) \quad (6)$$

$$\mathbf{h}_t = \mathbf{o}_t \circ \tanh(\mathbf{c}_t) \quad (7)$$

**Content Layer (Equation 2)** We refer to  $\tilde{\mathbf{c}}_t$  as the content layer, which is the output of an S-RNN. Evaluating the need for the multiplicative recurrent connections in this content layer is the focus of this work. The content layer is passed to the memory cell, which decides which parts of it to store.

**Memory Cell (Equations 3-5)** The memory cell  $\mathbf{c}_t$  is controlled by two gates. The input gate  $\mathbf{i}_t$  controls what part of the content ( $\tilde{\mathbf{c}}_t$ ) is written to the memory, while the forget gate  $\mathbf{f}_t$  controls what part of the memory is deleted by filtering the previous state of the memory ( $\mathbf{c}_{t-1}$ ). Writing to the memory is done by adding the filtered content ( $\mathbf{i}_t \circ \tilde{\mathbf{c}}_t$ ) to the retained memory ( $\mathbf{f}_t \circ \mathbf{c}_{t-1}$ ).

**Output Layer (Equations 6-7)** The output layer  $\mathbf{h}_t$  passes the memory cell through a tanh activation function and uses an output gate  $\mathbf{o}_t$  to read selectively from the squashed memory cell.

Our goal is to study how much each of these components contribute to the empirical performance of LSTMs. In particular, it is worth considering the memory cell in more detail to reveal why it could serve as a standalone powerful model of long-distance context. It is possible to show that it implicitly computes an *element-wise weighted sum* of all the previous content layers by expanding the recurrence relation in equation (5):

$$\begin{aligned} \mathbf{c}_t &= \mathbf{i}_t \circ \tilde{\mathbf{c}}_t + \mathbf{f}_t \circ \mathbf{c}_{t-1} \\ &= \sum_{j=0}^t \left( \mathbf{i}_j \circ \prod_{k=j+1}^t \mathbf{f}_k \right) \circ \tilde{\mathbf{c}}_j \\ &= \sum_{j=0}^t \mathbf{w}_j^t \circ \tilde{\mathbf{c}}_j \end{aligned} \quad (8)$$

Each weight  $\mathbf{w}_j^t$  is a product of the input gate  $\mathbf{i}_j$  (when its respective input  $\tilde{\mathbf{c}}_j$  was read) and every subsequent forget gate  $\mathbf{f}_k$ . An interesting property of these weights is that, like the gates, they are also soft element-wise binary filters.<sup>1</sup>

<sup>1</sup>GRUs also exhibit the same property of computing weighted sums over a content. Specifically, they compute weighted-averages, since the gates are coupled.

This sum is similar to recent architectures that rely on self-attention to learn context-dependent word representations (Cheng et al., 2016; Parikh et al., 2016; Vaswani et al., 2017). There are two major differences from self-attention: (1) instead of computing a weighted sum for each attention head, a separate weighted sum is computed for every dimension of the memory cell, (2) the weighted sum is accumulated with a dynamic program, enabling a linear rather than quadratic complexity in comparison to self-attention.

### 3 MEMORY CELLS ARE POWERFUL STANDALONE MODELS

The restricted space of element-wise weighted sums allows for easier mathematical analysis, visualization, and perhaps even learnability. However, constrained function spaces are also less expressive, and a natural question is whether these models will work well for NLP problems that need highly contextualized word representations. We hypothesize that the memory cell (which computes weighted sums) can function as a standalone contextualizer. To test this hypothesis, we present several simplifications of the LSTM’s architecture (Section 3.1), and show on a variety of NLP benchmarks that there is a qualitative performance difference between models that contain a memory cell and those that do not (Section 3.2). We conclude that the content and output layers are relatively minor contributors, and that the space of element-wise weighted sums is sufficiently powerful to compete with fully parameterized LSTMs (Section 3.3).

#### 3.1 SIMPLIFIED MODELS

The modeling power of LSTMs is commonly assumed to derive from the S-RNN in the content layer, with the rest of the model acting as a learning aid to bypass the vanishing gradient problem. We first isolate the S-RNN by ablating the gates (denoted as *LSTM – GATES* for consistency).

To test whether the memory cell has enough modeling power of its own, we take an LSTM and replace the S-RNN in the content layer from Equation 2 with a simple linear transformation, creating the *LSTM – S-RNN* model:

$$\begin{aligned}
 \tilde{\mathbf{c}}_t &= \mathbf{W}_{cx}\mathbf{x}_t \\
 \mathbf{i}_t &= \sigma(\mathbf{W}_{ih}\mathbf{h}_{t-1} + \mathbf{W}_{ix}\mathbf{x}_t + \mathbf{b}_i) \\
 \mathbf{f}_t &= \sigma(\mathbf{W}_{fh}\mathbf{h}_{t-1} + \mathbf{W}_{fx}\mathbf{x}_t + \mathbf{b}_f) \\
 \mathbf{c}_t &= \mathbf{i}_t \circ \tilde{\mathbf{c}}_t + \mathbf{f}_t \circ \mathbf{c}_{t-1} \\
 \mathbf{o}_t &= \sigma(\mathbf{W}_{oh}\mathbf{h}_{t-1} + \mathbf{W}_{ox}\mathbf{x}_t + \mathbf{b}_o) \\
 \mathbf{h}_t &= \mathbf{o}_t \circ \tanh(\mathbf{c}_t)
 \end{aligned} \tag{9}$$

We further simplify the LSTM by removing the output gate from Equation 7, leaving only the activation function in the output layer (*LSTM – S-RNN – OUT*):

$$\begin{aligned}
 \tilde{\mathbf{c}}_t &= \mathbf{W}_{cx}\mathbf{x}_t \\
 \mathbf{i}_t &= \sigma(\mathbf{W}_{ih}\mathbf{h}_{t-1} + \mathbf{W}_{ix}\mathbf{x}_t + \mathbf{b}_i) \\
 \mathbf{f}_t &= \sigma(\mathbf{W}_{fh}\mathbf{h}_{t-1} + \mathbf{W}_{fx}\mathbf{x}_t + \mathbf{b}_f) \\
 \mathbf{c}_t &= \mathbf{i}_t \circ \tilde{\mathbf{c}}_t + \mathbf{f}_t \circ \mathbf{c}_{t-1} \\
 \mathbf{h}_t &= \tanh(\mathbf{c}_t)
 \end{aligned} \tag{10}$$

After removing the S-RNN and the output gate from the LSTM, the entire ablated model can be written in a modular, compact form:

$$\mathbf{h}_t = \text{OUTPUT}\left(\sum_{j=0}^t \mathbf{w}_j^t \circ \text{CONTENT}(\mathbf{x}_j)\right) \tag{11}$$

where the content layer  $\text{CONTENT}(\cdot)$  and the output layer  $\text{OUTPUT}(\cdot)$  are both context-independent functions, making the entire model highly constrained and interpretable. The complexity of modeling contextual information is needed only for computing the weights  $\mathbf{w}_j^t$ . As we will see in Section 3.2, both of these ablations perform on par with LSTMs on language modeling, question answering, dependency parsing, and machine translation.

There are many other models that can be expressed in the weighted-sum form (Equation 11). In this work, we focus on the closest variant of LSTM that satisfies this property; removing the S-RNN and the output gate is sufficient for the content and output functions to be context-independent. We leave more thorough investigations into the necessity of the remaining architecture as future work.

### 3.2 EXPERIMENTS

We compare model performance on four NLP tasks, with an experimental setup that is lenient towards LSTMs and harsh towards its simplifications. In each case, we use existing implementations and previously reported hyperparameter settings. Since these settings were tuned for LSTMs, any simplification that performs equally to (or better than) LSTMs under these LSTM-friendly settings provides strong evidence that the ablated component is not a contributing factor. For each task we also report the mean and standard deviation of 5 runs of the LSTM settings to demonstrate the typical variance observed due to training with different random initializations.<sup>2</sup> The code and settings to replicate these experiments are publicly available.<sup>3</sup>

#### 3.2.1 LANGUAGE MODELING

We evaluate on two language modeling datasets: the Penn Treebank (PTB) (Marcus et al., 1993), and Google’s billion-word benchmark (BWB) (Chelba et al., 2014). PTB contains approximately 1M tokens over a vocabulary of 10K words. We used the implementation of Zaremba et al. (2014) while replacing any invocation of LSTMs with simpler models. We tested two of their configurations: *medium*, which uses two layers of 650-dimension LSTMs, and *large*, which uses two layers of 1500-dimension LSTMs.

BWB is about a thousand times larger than PTB, and uses a more diverse vocabulary of 800K words. Using the implementation of Józefowicz et al. (2016), we tested their *LSTM-2048-512* configuration. Our experiments use exactly the same hyperparameters (dimensions, dropout, learning rates, etc) that were originally tuned for LSTMs (Józefowicz et al., 2016). Following their implementation, we project the hidden state at each time step down to 512 dimensions. Due to the enormous size of this dataset, we stopped training after 5 epochs.

Table 1 shows overall model performance. In all three cases, replacing the LSTM’s content layer with a linear transformation results in small differences in perplexity. The most important result is that the small fluctuations in performance between the various gated architectures are minuscule in comparison to the enormous gap between the S-RNN (*LSTM – GATES*) and the original LSTM. This striking difference strongly supports our hypothesis that the weighted sums computed by the gates – not the S-RNN – is the recurrent model that contributes mostly strongly to the final performance.

#### 3.2.2 QUESTION ANSWERING

For question answering, we use two different QA systems on the Stanford question answering dataset (SQuAD) (Rajpurkar et al., 2016): the Bidirectional Attention Flow model (BiDAF) (Seo et al., 2016) and DrQA (Chen et al., 2017). BiDAF contains 3 LSTMs, which are referred to as the phrase layer, the modeling layer, and the span end encoder. Our experiments replace each of these LSTMs with their simplified counterparts. We directly use the implementation of BiDAF from AllenNLP (Gardner et al., 2017), and all experiments reuse the existing hyperparameters that were tuned for LSTMs. Likewise, we use an open-source implementation of DrQA<sup>4</sup> and replace only the LSTMs, while leaving everything else intact.

Table 2 shows that all the gated models do comparably. Most importantly, ablating the S-RNN from the LSTM has a minor effect in comparison to the drop in performance when ablating the gates.

#### 3.2.3 DEPENDENCY PARSING

For dependency parsing, we use the Deep Biaffine Dependency Parser (Dozat & Manning, 2016), which relies on stacked bidirectional LSTMs to learn context-sensitive word embeddings for de-

<sup>2</sup>Due to time constraints, we only include the reported LSTM results for BWB (Józefowicz et al., 2015).

<sup>3</sup><http://anonymous>

<sup>4</sup><https://github.com/hitvoice/DrQA>

Configuration	Model	Perplexity
PTB (Medium Model)	LSTM	$83.9 \pm 0.3$
	- GATES	140.9
	- S-RNN	80.5
	- S-RNN - OUT	81.6
PTB (Large Model)	LSTM	$78.8 \pm 0.2$
	- GATES	126.1
	- S-RNN	76.0
	- S-RNN - OUT	78.5
BWB	LSTM (Józefowicz et al., 2016)	47.5
	- GATES	82.2
	- S-RNN	45.4
	- S-RNN - OUT	47.9

Table 1: The performance of simplified LSTM architectures on language modeling benchmarks, measured by perplexity.

System	Model	EM	F1
BiDAF	LSTM	$67.9 \pm 0.3$	$77.5 \pm 0.2$
	- GATES	62.9	73.3
	- S-RNN	68.4	78.2
	- S-RNN - OUT	67.4	77.2
DrQA	LSTM	$68.8 \pm 0.2$	$78.2 \pm 0.2$
	- GATES	56.4	66.5
	- S-RNN	67.7	77.0
	- S-RNN - OUT	67.0	76.2

Table 2: The performance of simplified LSTM architectures on the question answering benchmark, SQuAD, measured by exact match (EM) and span overlap (F1).

Model	UAS	LAS
LSTM	$90.60 \pm 0.21$	$88.05 \pm 0.33$
- GATES	87.75	84.61
- S-RNN	90.77	88.49
- S-RNN - OUT	90.70	88.31

Table 3: The performance of simplified LSTM architectures on the universal dependencies parsing benchmark, measured by unlabeled attachment score (UAS) and labeled attachment score (LAS).

terminating arcs between a pair of words. We directly use their released implementation, which is evaluated on the Universal Dependencies English Web Treebank v1.3 (Silveira et al., 2014). In our experiments, we use the existing hyperparameters and only replace the LSTMs with the simplified architectures.

We observe the same pattern in the ablations for dependency parsing. The differences in performance between the gated models fall within the differences between multiple experiments with LSTMs. Consistent with ablation results from other tasks, removing the gating mechanisms causes a 3-4 point drop in performance.

### 3.2.4 MACHINE TRANSLATION

For machine translation, we used OpenNMT (Klein et al., 2017) to train English to German translation models on the multi-modal benchmarks from WMT 2016 (used in OpenNMT’s readme file). We use OpenNMT’s default model and hyperparameters, replacing the stacked bidirectional LSTM of its

Model	BLEU
LSTM	35.95
– GATES	12.22
– S-RNN	36.66
– S-RNN – OUT	36.39

Table 4: The performance of simplified LSTM architectures on the WMT 2016 multi-modal English to German translation benchmark, measured by BLEU.

encoder with the simplified architectures. Table 4 shows that while models containing memory cells perform more-or-less on par, removing the memory cell yields a substantial performance drop.

### 3.3 DISCUSSION

In the above experiments, we show three major ablations of the LSTM. In the S-RNN experiments (*LSTM – GATES*), we ablate the memory cell and the output layer. In the *LSTM – S-RNN* and *LSTM – S-RNN – OUT* experiments, we ablate the S-RNN. As consistent with previous literature, removing the memory cell degrades performance drastically. In contrast, removing the S-RNN makes little to no difference in the final performance, suggesting that the memory cell alone is largely responsible for the success of LSTMs in NLP. The results also confirm our hypothesis that weighted sums of context words is a powerful, yet more interpretable, model of contextual information.

## 4 WEIGHT VISUALIZATION

Given the empirical evidence that LSTMs are effectively learning weighted sums of the content layers, it is natural to investigate what weights the model learns in practice. Using the more mathematically transparent simplification of LSTMs, we can visualize the weights  $w_j^t$  that are placed on every input  $j$  at every timestep  $t$  (see Equation 11).

Unlike attention mechanisms, these weights are vectors rather than scalar values. Therefore, we can only provide a coarse-grained visualization of the weights by rendering their  $L^2$ -norm, as shown in Table 5. In the visualization, each column indicates the word represented by the weighted sum, and each row indicates the word over which the weighted sum is computed. Dark horizontal streaks indicate the duration for which a word was remembered. Unsurprisingly, the weights on the diagonal are always the largest since it indicates the weight of the current word. More interesting task-specific patterns emerge when inspecting the off-diagonals that represent the weight on the context words.

The first visualization uses the language model from BWB. Due to the language modeling setup, there are only non-zero weights on the current or previous words. We find that the common function words are quickly forgotten, while infrequent words that signal the topic are remembered over very long distances.

The second visualization uses the dependency parser. In this setting, since the recurrent architectures are bidirectional, there are non-zero weights on all words in the sentence. The top-right triangle indicates weights from the forward direction, and the bottom-left triangle indicates from the backward direction. For syntax, we see a significantly different pattern. Function words that are useful for determining syntax are more likely to be remembered. Weights on head words are also likely to persist until the end of a constituent.

This illustration provides only a glimpse into what the model is capturing, and perhaps future, more detailed visualizations that take the individual dimensions into account can provide further insight into what LSTMs are learning in practice.

## 5 RELATED WORK

Many variants of LSTMs (Hochreiter & Schmidhuber, 1997) have been previously explored. These typically consist of a different parameterization the gates, such as LSTMs with peephole connections (Gers & Schmidhuber, 2000), or a rewiring of the connections, such as GRUs (Cho et al., 2014).



However, these modifications invariably maintain the recurrent content layer. Even more systematic explorations of LSTM variants (Józefowicz et al., 2015; Greff et al., 2016; Zoph & Le, 2017) do not question the importance of the embedded S-RNN. This is the first study to provide apples-to-apples comparisons between LSTMs and LSTMs *without* the recurrent content layer.

Several other recent works have also reported promising results with recurrent models that are vastly simpler than LSTMs, such as quasi-recurrent neural networks (Bradbury et al., 2016), strongly-typed recurrent neural networks (Balduzzi & Ghifary, 2016), kernel neural networks (Lei et al., 2017), and simple recurrent units (Lei & Zhang, 2017), making it increasingly apparent that LSTMs are over-parameterized. While these works indicate an obvious trend, their focus is not to provide insight into what exactly LSTMs are learning. In our carefully controlled ablation studies, we propose and evaluate the minimal changes required to test our hypothesis that LSTMs are powerful because they dynamically compute element-wise weighted sums of content layers.

As mentioned in Section 2, this weighted-sum view of LSTMs is highly related to neural attention (Bahdanau et al., 2015), which assigns a normalized scalar weight to each element as a function of its compatibility with an external element. The ability to inspect attention weights has driven the use of more interpretable neural models. Self-attention (Cheng et al., 2016; Parikh et al., 2016) extends this notion by computing intra-sequence attention. Vaswani et al. (2017) further showed that state-of-the-art machine translation can be achieved using only self-attention and without LSTMs. Recently, Arora et al. (2017) proposed a theory-driven approach to assign scalar weights to elements in a bag of words. The success of self-attention corroborates our findings that weighted sums are indeed a more effective method of learning context-sensitive representations than previously appreciated.

## 6 CONCLUSION

We presented an alternate view of LSTMs: they are a hybrid of S-RNNs and a gated model that dynamically computes weighted sums of the S-RNN outputs. Our experiments investigated whether the S-RNN is a necessary component of LSTMs. In other words, are the gates alone as powerful of a model as an LSTM? Results across four major NLP tasks (language modeling, question answering, dependency parsing, and machine translation) indicate that LSTMs suffer little to no performance loss when removing the S-RNN, but removing the gates can degrade performance substantially. This provides evidence that the gating mechanism is doing the heavy lifting in modeling context, and that element-wise weighted sums of context-independent functions of the inputs are often as effective as fully-parameterized LSTMs.

This work sheds light on the inner workings of the relatively opaque LSTM. By removing the S-RNN and the output gate, we also show that the resulting model is a far more mathematically transparent variant of LSTMs. This transparency enables a visualization of how the context affects the output of the model at every timestep, much like in attention-based models. We hope that this new outlook on LSTMs will foster better and more efficient models of contextualization.

## REFERENCES

- Yossi Adi, Einat Kermany, Yonatan Belinkov, Ofer Lavi, and Yoav Goldberg. Fine-grained analysis of sentence embeddings using auxiliary prediction tasks. In *ICLR*, 2017.
- Sanjeev Arora, Yingyu Liang, and Tengyu Ma. A simple but tough-to-beat baseline for sentence embeddings. In *ICLR*, 2017.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. In *ICLR*, 2015.
- David Balduzzi and Muhammad Ghifary. Strongly-typed recurrent neural networks. In *Proceedings of the 33rd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016*, pp. 1292–1300, 2016. URL <http://jmlr.org/proceedings/papers/v48/balduzzi16.html>.
- Yoshua Bengio, Patrice Y. Simard, and Paolo Frasconi. Learning long-term dependencies with gradient descent is difficult. *IEEE Transactions on Neural Networks*, 5(2):157–166, 1994.

- James Bradbury, Stephen Merity, Caiming Xiong, and Richard Socher. Quasi-recurrent neural networks. *CoRR*, abs/1611.01576, 2016.
- Ciprian Chelba, Tomas Mikolov, Mike Schuster, Qi Ge, Thorsten Brants, and Phillipp Koehn. One billion word benchmark for measuring progress in statistical language modeling. In *INTERSPEECH*, 2014.
- Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. Reading Wikipedia to answer open-domain questions. In *Association for Computational Linguistics (ACL)*, 2017.
- Jianpeng Cheng, Li Dong, and Mirella Lapata. Long short-term memory-networks for machine reading. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp. 551–561, Austin, Texas, November 2016. Association for Computational Linguistics. URL <https://aclweb.org/anthology/D16-1053>.
- 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. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 1724–1734, Doha, Qatar, October 2014. Association for Computational Linguistics. URL <http://www.aclweb.org/anthology/D14-1179>.
- Timothy Dozat and Christopher D. Manning. Deep biaffine attention for neural dependency parsing. *CoRR*, abs/1611.01734, 2016.
- Jeffrey L. Elman. Finding structure in time. *Cognitive Science*, 14:179–211, 1990.
- Matt Gardner, Joel Grus, Mark Neumann, Oyvind Tafjord, Pradeep Dasigi, Nelson Liu, Matthew Peters, Michael Schmitz, and Luke Zettlemoyer. Allennlp: A deep semantic natural language processing platform, 2017. URL [http://allennlp.org/papers/AllenNLP\\_white\\_paper.pdf](http://allennlp.org/papers/AllenNLP_white_paper.pdf).
- Felix A. Gers and Jürgen Schmidhuber. Recurrent nets that time and count. In *IJCNN*, 2000.
- Klaus Greff, Rupesh K Srivastava, Jan Koutník, Bas R Steunebrink, and Jürgen Schmidhuber. Lstm: A search space odyssey. *IEEE Transactions on Neural Networks and Learning Systems*, 2016.
- Luheng He, Kenton Lee, Mike Lewis, and Luke Zettlemoyer. Deep semantic role labeling: What works and what’s next. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, 2017.
- Sepp Hochreiter. Untersuchungen zu dynamischen neuronalen netzen. *Diploma, Technische Universität München*, 91, 1991.
- Sepp Hochreiter and Jürgen Schmidhuber. Long Short-term Memory. *Neural computation*, 9(8): 1735–1780, 1997.
- Rafal Józefowicz, Wojciech Zaremba, and Ilya Sutskever. An empirical exploration of recurrent network architectures. In *ICML*, 2015.
- Rafal Józefowicz, Oriol Vinyals, Mike Schuster, Noam Shazeer, and Yonghui Wu. Exploring the limits of language modeling. *arXiv preprint arXiv:1602.02410*, 2016.
- Guillaume Klein, Yoon Kim, Yuntian Deng, Jean Senellart, and Alexander M. Rush. Opennmt: Open-source toolkit for neural machine translation. In *Proc. ACL*, 2017. doi: 10.18653/v1/P17-4012. URL <https://doi.org/10.18653/v1/P17-4012>.
- Tao Lei and Yu Zhang. Training rnns as fast as cnns. *arXiv preprint arXiv:1709.02755*, 2017.
- Tao Lei, Wengong Jin, Regina Barzilay, and Tommi Jaakkola. Deriving neural architectures from sequence and graph kernels. In *ICML*, 2017.
- Tal Linzen, Emmanuel Dupoux, and Yoav Goldberg. Assessing the ability of lstms to learn syntax-sensitive dependencies. *TACL*, 4:521–535, 2016.

- Mitchell P. Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. Building a large annotated corpus of english: The penn treebank. *Computational Linguistics*, 19:313–330, 1993.
- Ankur Parikh, Oscar Täckström, Dipanjan Das, and Jakob Uszkoreit. A decomposable attention model for natural language inference. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp. 2249–2255, Austin, Texas, November 2016. Association for Computational Linguistics. URL <https://aclweb.org/anthology/D16-1244>.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100, 000+ questions for machine comprehension of text. In *EMNLP*, 2016.
- Min Joon Seo, Aniruddha Kembhavi, Ali Farhadi, and Hannaneh Hajishirzi. Bidirectional attention flow for machine comprehension. *CoRR*, abs/1611.01603, 2016.
- Natalia Silveira, Timothy Dozat, Marie-Catherine de Marneffe, Samuel Bowman, Miriam Connor, John Bauer, and Christopher D. Manning. A gold standard dependency corpus for English. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC-2014)*, 2014.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *arXiv preprint arXiv:1706.03762*, 2017.
- Wojciech Zaremba, Ilya Sutskever, and Oriol Vinyals. Recurrent neural network regularization. *arXiv preprint arXiv:1409.2329*, 2014.
- Barret Zoph and Quoc V Le. Neural architecture search with reinforcement learning. In *ICLR*, 2017.