

---

# Feed-Forward Networks with Attention Can Solve Some Long-Term Memory Problems

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

Recently, recurrent neural networks (RNNs) have been augmented with “attention” mechanisms which compute a fixed-length representation of entire sequences. We propose a simplified model of attention which is applicable to feed-forward neural networks and demonstrate that it can solve some long-term memory problems (specifically, those where temporal order doesn’t matter). In fact, we show empirically that our model can solve these problems for sequence lengths which are both longer and more widely varying than has been shown for RNNs.

## 1 Models for Sequential Data

Many problems in machine learning are best formulated using sequential data, i.e. data where a given observation may be dependent on previous observations. Such problems can be coarsely classified as sequence transduction (producing a new sequence given an input sequence), sequence classification (producing a single label or vector from an entire sequence), or sequence generation (producing a sequence from no input) tasks. Appropriate models for these tasks must be able to capture temporal dependencies in sequences, potentially of arbitrary length.

### 1.1 Recurrent Neural Networks

One such class of models are recurrent neural networks (RNNs). Recent advances in optimization techniques [], initialization methods [], and model architectures [] have facilitated the success of RNNs in a wide variety of fields, including machine translation [], speech recognition [], handwriting recognition [], ... This variety of applications demonstrates their applicability to most machine learning tasks involving sequences. In their simplest form, RNNs can be considered as a function  $f$  whose output at time  $t$  depends on input  $x_t$ , the previous state  $h_{t-1}$ , and the model’s learnable parameters  $\theta$ :

$$h_t = f(x_t, h_{t-1}, \theta)$$

In the supervised setting, the parameters  $\theta$  are optimized with respect to a loss function which measures the performance of  $f$ . A common approach is to use backpropagation through time [], which “unrolls” the RNN over time steps to compute the gradient of  $\theta$  with respect to the loss.

Despite the success of RNNs, their recursive nature makes them hard to both optimize and parallelize. Because the same function  $f$  is applied repeatedly over time, the gradient with respect to  $\theta$  can easily explode or vanish []. The use of gating architectures [], sophisticated optimization techniques [], gradient clipping [], and/or careful initialization [] can help mitigate this issue. However, these approaches do not solve this problem, and as a result RNNs are in practice typically only applied in tasks where sequential dependencies span at most hundreds of time steps []. Very long sequences can also make training computationally inefficient due to the fact that RNNs must be evaluated sequentially and cannot be fully parallelized.

## 1.2 Attention

A recently proposed method for easier modeling of long-term dependencies is “attention”. Attention mechanisms allow for a more direct dependence between the state of the model at different points in time. Following the definition from [], given a model which produces a hidden state  $h_t$  at each time step, attention-based models first compute a “context” vector  $c_t$  as the weighted mean of the state sequence  $h$  by

$$c_t = \sum_{j=1}^T \alpha_{tj} h_j$$

where  $T$  is the total number of time steps in the input sequence and  $\alpha_{tj}$  is a weight computed at each time step  $t$  for each state  $h_j$ . These context vectors are then used to compute a new state sequence  $s$ , where  $s_t$  depends on  $s_{t-1}$ ,  $c_t$  and, for sequence prediction, the model’s output at  $t - 1$ . The weightings  $\alpha_{tj}$  are then computed by

$$e_{tj} = a(s_{t-1}, h_j), \alpha_{tj} = \frac{\exp(e_{tj})}{\sum_{k=1}^T \exp(e_{tk})}$$

where  $a$  is a learned function which can be thought of as computing a scalar importance value for  $h_j$  given the value of  $h_j$  and the previous state  $s_{t-1}$ . This formulation allows the new state sequence  $s$  to have more direct access to the entire state sequence  $h$ . Attention-based RNNs have proven effective in a variety of sequence transduction tasks []. RNNs with attention can be seen as analogous to the recently proposed Memory Network [] and Neural Turing Machine [] models.

## 1.3 Feed-Forward Attention

A straightforward simplification to the attention mechanism described above which would allow it to be applied to sequence classification tasks could be formulated as follows: Instead of a sequence of context vectors, we produce a single context vector  $c$  as

$$c = \sum_{t=1}^T \alpha_t h_t, e_t = a(h_j), \alpha_t = \frac{\exp(e_t)}{\sum_{k=1}^T \exp(e_k)}$$

As before,  $a$  is a learnable function, but it now only depends on  $h_j$ . In this formulation, attention can be seen as producing a fixed-length embedding  $c$  of the input sequence by computing an adaptive weighted average of the state sequence  $h$ .

A consequence of using an attention mechanism is that the context vector  $c$  can model temporal dependencies because attention performs integration over time. It follows that by using this simplified form of attention, a model could perform sequence classification even if the calculation of  $h_t$  was feed-forward, i.e.  $h_t = f(x_t)$ . While using a completely feed-forward model for sequential modeling tasks will sacrifice the ability to solve some problems, we show that for certain tasks, feed-forward networks with attention can perform arbitrary-length sequence classification more effectively than RNNs.

We note here that feed-forward models can be used for sequence classification when the sequence length  $T$  is fixed, but when  $T$  varies across sequences, some form of temporal integration is necessary. An obvious straightforward choice, which can be seen as an extreme oversimplification of attention, would be to compute  $c$  as the unweighted average of the state sequence  $h_t$ . We will also explore the effectiveness of this approach for sequence classification.

## 2 Toy Long-Term Memory Problems

Different tasks

Previous results

108	<b>2.1 Fixed-Length Experiment</b>
109	
110	<b>2.2 Variable-length Experiment</b>
111	
112	<b>3 Limitations</b>
113	
114	
115	
116	
117	
118	
119	
120	
121	
122	
123	
124	
125	
126	
127	
128	
129	
130	
131	
132	
133	
134	
135	
136	
137	
138	
139	
140	
141	
142	
143	
144	
145	
146	
147	
148	
149	
150	
151	
152	
153	
154	
155	
156	
157	
158	
159	
160	
161	