# Tree-Structured Composition in Neural Networks without Tree-Structured Architectures

**Samuel R. Bowman, Christopher D. Manning, and Christopher Potts**
Stanford University
Stanford, CA 94305-2150
{sbowman, manning, cgpotts}@stanford.edu

## Abstract

Tree-structured neural networks encode a particular tree geometry for a sentence in the network design. However, these models have at best only slightly outperformed simpler sequence-based models. We hypothesize that neural sequence models like LSTMs are in fact able to discover and implicitly use recursive compositional structure, at least for tasks with clear cues to that structure in the data. We demonstrate this possibility using an artificial data task for which recursive compositional structure is crucial, and find an LSTM-based sequence model can indeed learn to exploit the underlying tree structure. However, its performance consistently lags behind that of tree models, even on large training sets, suggesting that tree-structured models are more effective at exploiting recursive structure.

## 1 Introduction

Neural networks that encode sentences as real-valued vectors have been successfully used in a wide array of NLP tasks, including translation [1], parsing [2], and sentiment analysis [3]. These models are generally either sequence models based on recurrent neural networks, which build representations incrementally from left to right [4, 1], or tree-structured models based on *recursive* neural networks, which build representations incrementally according to the hierarchical structure of linguistic phrases [5, 6].

While both model classes perform well on many tasks, and both are under active development, tree models are often presented as the more principled choice, since they align with standard linguistic assumptions about constituent structure and the compositional derivation of complex meanings. Nevertheless, tree models have not shown the kinds of dramatic performance improvements over sequence models that their billing would lead one to expect: head-to-head comparisons with sequence models show either modest improvements [3] or none at all [7].

We propose a possible explanation for these results: standard sequence models can learn to exploit recursive syntactic structure in generating representations of sentence meaning, thereby learning to use the structure that tree models are explicitly designed around. This requires that sequence models be able to identify syntactic structure in natural language. We believe this is plausible on the basis of other recent research [8, 9]. In this paper, we evaluate whether LSTM sequence models are able to use such structure to guide interpretation, focusing on cases where syntactic structure is clearly indicated in the data.

We compare standard tree and sequence models on their handling of recursive structure by training the models on sentences whose length and recursion depth are limited, and then testing them on longer and more complex sentences, such that only models that exploit the recursive structure will be able to generalize in a way that yields correct interpretations for these test sentences. Our methods extend those of our earlier work in [10], which introduces an experiment and corresponding artificial dataset to test this ability in two tree models. We adapt that experiment to sequence models by

| | | |
|---:|:---:|:---|
| *not $p_3$* | $\wedge$ | *$p_3$* |
| *$p_3$* | $\sqsubset$ | *$p_3$ or $p_2$* |
| *(not $p_2$) and $p_6$* | $\mid$ | *not($p_6$ or ($p_5$ or $p_3$))* |
| *$p_4$ or (not(($p_1$ or $p_6$) or $p_4$))* | $\sqsubset$ | *not(((( not $p_6$) or (not $p_4$)) and (not $p_5$)) and ($p_6$ and $p_6$))* |

Table 1: Examples of short to moderate length pairs from the artificial data introduced in [10]. We only show the parentheses that are needed to disambiguate the sentences rather than the full binary bracketings that the models use.

decorating the statements with an explicit bracketing, and we use this design to compare an LSTM sequence model with three tree models, with a focus on what data each model needs in order to generalize well.

As in [10], we find that standard tree neural networks are able to make the necessary generalizations, with their performance decaying gradually as the structures in the test set grow in size. We additionally find that extending the training set to include larger structures mitigates this decay. Then considering sequence models, we find that a single-layer LSTM is also able to generalize to unseen large structures, but that it does this only when trained on a larger and more complex training set than is needed by the tree models to reach the same generalization performance.

Our results engage with those of [8] and [2], who find that sequence models can learn to recognize syntactic structure in natural language, at least when trained on explicitly syntactic tasks. The simplest model presented in [8] uses an LSTM sequence model to encode each sentence as a vector, and then generates a linearized parse (a sequence of brackets and constituent labels) with high accuracy using only the information present in the vector. This shows that the LSTM is able to identify the correct syntactic structures and also hints that it is able to develop a generalizable method for encoding these structures in vectors. However, the massive size of the dataset needed to train that model, 250M tokens, leaves open the possibility that it primarily learns to generate only tree structures that it has already seen, representing them as simple hashes—which would not capture unseen tree structures—rather than as structured objects. Our experiments, though, show that LSTMs can learn to understand tree structures when given enough data, suggesting that there is no fundamental obstacle to learning this kind of structured representation. We also find, though, that sequence models lag behind tree models across the board, even on training corpora that are quite large relative to the complexity of the underlying grammar, suggesting that tree models can play a valuable role in tasks that require recursive interpretation.

## 2 Recursive structure in artificial data

**Reasoning about entailment**  The data that we use define a version of the *recognizing textual entailment* task, in which the goal is to determine what kind of logical consequence relation holds between two sentences, drawing on a small fixed vocabulary of relations such as entailment, contradiction, and synonymy. This task is well suited to evaluating neural network models for sentence interpretation: models must develop comprehensive representations of the meanings of each sentence to do well at the task, but the data do not force these representations to take a specific form, allowing the model to learn whatever kind of representations it can use most effectively.

The data we use are labeled with the seven mutually exclusive logical relations of [11], which distinguish entailment in two directions ($\sqsubset$, $\sqsupset$), equivalence ($\equiv$), exhaustive and non-exhaustive contradiction ($\wedge$, $\mid$), and two types of semantic independence ($\#$, $\smile$).

**The artificial language**  The language described in [10] (§4) is designed to highlight the use of recursive structure with minimal additional complexity. Its vocabulary consists only of six unanalyzed word types ($p_1, p_2, p_3, p_4, p_5, p_6$), *and*, *or*, and *not*. Sentences of the language can be straightforwardly interpreted as statements of propositional logic (where the six unanalyzed words types are variable names), and labeled sentence pairs can be interpreted as theorems of that logic. Some example pairs are provided in Table 1.

Crucially, the language is defined such that any sentence can be embedded under negation or conjunction to create a new sentence, allowing for arbitrary-depth recursion, and such that the scope of

negation and conjunction are determined only by bracketing with parentheses (rather than bare word order). The compositional structure of each sentence can thus be an arbitrary tree, and interpreting a sentence correctly requires using that structure.

The data come with parentheses representing a complete binary bracketing. Our models use this information in two ways. For the tree models, the parentheses are not word tokens, but rather are used in the expected way to build the tree. For the sequence model, the parentheses are word tokens with associated learned embeddings. This approach provides the models with equivalent data, so their ability to handle unseen structures can be reasonably compared.

**The data** Our sentence pairs are divided into thirteen bins according to the number of logical connectives (*and, or, not*) in the longer of the two sentences in each pair. We test each model on each bin separately (58k total examples, using an 80/20% train/test split) in order to evaluate how each model's performance depends on the complexity of the sentences. In three experiments, we train our models on the training portions of bins 0–3 (62k examples), 0–4 (90k), and 0–6 (160k), and test on every bin but the trivial bin 0. Capping the size of the training sentences allows us to evaluate how the models interpret the sentences: if a model's performance falls off abruptly above the cutoff, it is reasonable to conclude that it relies heavily on specific sentence structures and cannot generalize to new structures. If a model's performance decays gradually[1] with no such abrupt change, then it must have learned a more generally valid interpretation function for the language which respects its recursive structure.

## 3   Testing sentence models on entailment

We use the architecture depicted in Figure 1a, which builds on the one used in [10]. The model architecture uses two copies of a single sentence model (a tree or sequence model) to encode the premise and hypothesis (left and right side) expressions, and then uses those encodings as the features for a multilayer classifier which predicts one of the seven relations. Since the encodings are computed separately, the sentence models must encode complete representations of the meanings of the two sentences for the downstream model to be able to succeed.

**Classifier** The classifier component of the model consists of a combining layer which takes the two sentence representations as inputs, followed by two neural network layers, then a softmax classifier. For the combining layer, we use a neural tensor network (NTN, [12]) layer, which sums the output of a plain recursive/recurrent neural network layer with a vector computed using two multiplications with a learned (full rank) third-order tensor parameter:
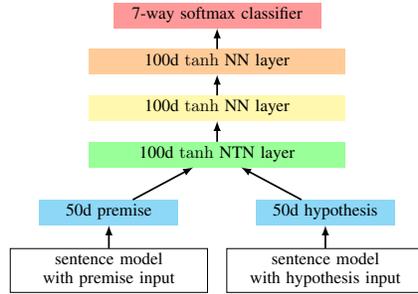
$$
(1) \qquad \vec{y}_{NN} = \tanh(\mathbf{M} \begin{bmatrix} \vec{x}^{(l)} \\ \vec{x}^{(r)} \end{bmatrix} + \vec{b})
$$

$$
(2) \qquad \vec{y}_{NTN} = \vec{y}_{NN} + \tanh(\vec{x}^{(l)T} \mathbf{T}^{[1...n]} \vec{x}^{(r)})
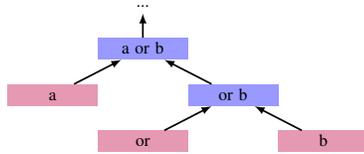$$

Our model is largely identical to the model from [10], but adds the two additional $\tanh$ NN layers, which we found help performance across the board, and also uses the NTN combination layer when evaluating all four models, rather than just the TreeRNTN model, so as to ensure that the sentence models are compared in as similar a setting as possible.

We only study models that encode entire sentences in fixed length vectors, and we set aside models with attention [13], a technique which gives the downstream model (here, the classifier) the potential to access each input token individually through a soft content addressing system. While attention simplifies the problem of learning complex correspondences between input and output, there is no apparent reason to believe that it should improve or harm a model's ability to track structural information like a given token's position in a tree. As such, we expect our results to reflect the same basic behaviors that would be seen in attention-based models.
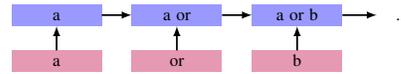
---

[1]Since sentences are fixed-dimensional vectors of fixed-precision floating point numbers, all models will make errors on sentences above some length, and L2 regularization (which helps overall performance) exacerbates this by discouraging the model from using the kind of numerically precise, nonlinearity-saturating functions that generalize best.

(a) The general architecture shared across models.



(b) The architecture for the tree-structured sentence models. Terminal nodes are learned embeddings and nonterminal nodes are NN, NTN, or TreeLSTM layers.

(c) The architecture for the sequence sentence model. Nodes in the lower row are learned embeddings and nodes in the upper row are LSTM layers.

Figure 1: In our model, two copies of a sentence model—based on either tree (b) or sequence (c) models—encode the two input sentences. A multilayer classifier component (a) then uses the resulting vectors to predict a label that reflects the logical relationship between the two sentences.

**Sentence models**  The sentence encoding component of the model transforms the (learned) embeddings of the input words for each sentence into a single vector representing that sentence. We experiment with tree-structured models (Figure 1b) with TreeRNN (eqn. 1), TreeRNTN (eqn. 2), and TreeLSTM [3] activation functions. In addition, we use a sequence model (Figure 1c) with an LSTM activation function [14] implemented as in [15]. In experiments with a simpler non-LSTM RNN sequence model, the model tended to badly underfit the training data, and those results are not included here.

**Training**  We randomly initialize all embeddings and layer parameters, and train them using mini-batch stochastic gradient descent with AdaDelta [16] learning rates. Our objective is the standard negative log likelihood classification objective with L2 regularization (tuned on a separate train/test split). All models were trained for 100 epochs, after which all had largely converged without significantly declining from their peak performances.

## 4   Results and discussion

The results are shown in Figure 2. The tree models fit the training data well, reaching 98.9, 98.8, and 98.4% overall accuracy respectively in the ≤6 setting, with the LSTM underfitting slightly at 94.8%. In that setting, all models generalized well to structures of familiar length, with the tree models all surpassing 97% on examples in bin 4, and the LSTM reaching 94.8%. On the longer test sentences, the tree models decay smoothly in performance across the board, while the LSTM decays more quickly and more abruptly, with a striking difference in the ≤4 setting, where LSTM performance falls 10% from bin 4 to bin 5, compared to 4.4% for the next worse model. However, the LSTM improves considerably with more ample training data in the ≤6 condition, showing only a 3% drop and generalization results better than the best model's in the ≤3 setting.

All four models robustly beat the simple baselines reported in [10]: the most frequent class occurs just over 50% of the time and a neural bag of words model does reasonably on the shortest examples but falls below 60% by bin 4.

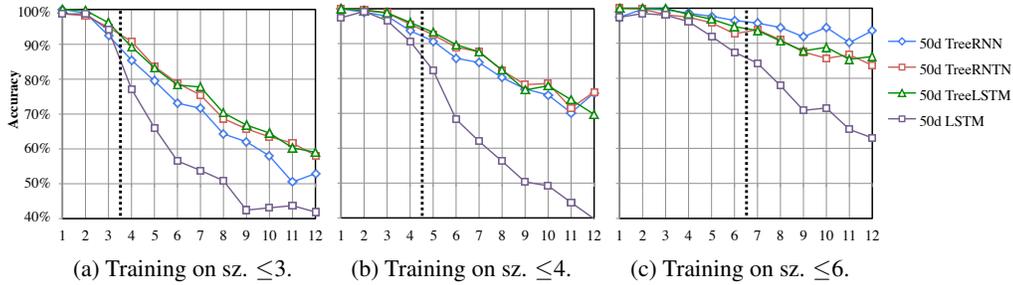|  | (a) Training on sz. $\leq 3$. | (b) Training on sz. $\leq 4$. | (c) Training on sz. $\leq 6$. |

Figure 2: Test accuracy on three experiments with increasingly rich training sets. The horizontal axis on each graph divides the test set expression pairs into bins by the number of logical operators in the more complex of the two expressions in the pair. The dotted line shows the size of the largest examples in the training set in each experiment.
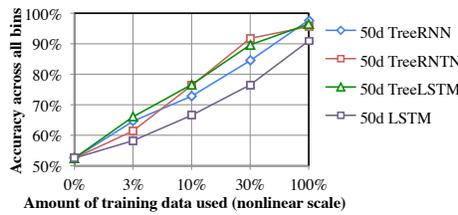


Figure 3: Learning curve for the $\leq 6$ experiment.

The learning curve (Figure 3) suggests that additional data is unlikely to change these basic results. The LSTM lags behind the tree models across the curve, but appears to gain accuracy at a similar rate.

## 5   Conclusion

We find that all four models are able to effectively exploit a recursively defined language to interpret sentences with complex unseen structures. We find that tree models' biases allow them to do this with greater efficiency, outperforming sequence-based models substantially in every experiment. However, our sequence model is nonetheless able to generalize smoothly from seen sentence structures to unseen ones, showing that its lack of explicit recursive structure does not prevent it from recognizing recursive structure in our artificial language.

We interpret these results as evidence that both tree and sequence architectures can play valuable roles in the construction of sentence models over data with recursive syntactic structure. Tree architectures provide an explicit bias that makes it possible to efficiently learn to compositional interpretation, which is difficult for sequence models. Sequence models, on the other hand, lack this bias, but have other advantages. Since they use a consistent graph structure across examples, it is easy to accelerate minibatch training in ways that yield substantially faster training times than are possible with tree models, especially with GPUs. In addition, when sequence models integrate each word into a partial sentence representation, they have access to the entire sentence representation up to that point, which may provide valuable cues for the resolution of lexical ambiguity, which is not present in our artificial language, but is a serious concern in natural language text.

Finally, we suggest that, because of the well-supported linguistic claim that the kind of recursive structure that we study here is key to the understanding of real natural languages, there is likely to be value in developing sequence models that can more efficiently exploit this structure without fully sacrificing the flexibility that makes them succeed.

### Acknowledgments

## References

[1] Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. Sequence to sequence learning with neural networks. In *Proc. NIPS*, 2014.

[2] Chris Dyer, Miguel Ballesteros, Wang Ling, Austin Matthews, and Noah A. Smith. Transition-based dependency parsing with stack long short-term memory. In *Proc. ACL*, 2015.

[3] Kai Sheng Tai, Richard Socher, and Christopher D. Manning. Improved semantic representations from tree-structured long short-term memory networks. In *Proc. ACL*, 2015.

[4] Jeffrey L. Elman. Finding structure in time. *Cognitive science*, 14(2), 1990.

[5] Christoph Goller and Andreas Kuchler. Learning task-dependent distributed representations by backpropagation through structure. In *Proc. IEEE International Conference on Neural Networks*, 1996.

[6] Richard Socher, Jeffrey Pennington, Eric H. Huang, Andrew Y. Ng, and Christopher D. Manning. Semi-supervised recursive autoencoders for predicting sentiment distributions. In *Proc. EMNLP*, 2011.

[7] Minh-Thang Luong Li, Jiwei, Dan Jurafsky, and Eudard Hovy. When are tree structures necessary for deep learning of representations? *Proc. EMNLP*, 2015.

[8] Oriol Vinyals, Lukasz Kaiser, Terry Koo, Slav Petrov, Ilya Sutskever, and Geoffrey Hinton. Grammar as a foreign language. In *Proc. NIPS*, 2015.

[9] Andrej Karpathy, Justin Johnson, and Fei-Fei Li. Visualizing and understanding recurrent networks. *arXiv:1506.02078*, 2015.

[10] Samuel R. Bowman, Christopher Potts, and Christopher D. Manning. Recursive neural networks can learn logical semantics. In *Proc. of the 3rd Workshop on Continuous Vector Space Models and their Compositionality*, 2015.

[11] Bill MacCartney and Christopher D. Manning. An extended model of natural logic. In *Proc. of the Eighth International Conference on Computational Semantics*, 2009.

[12] Danqi Chen, Richard Socher, Christopher D. Manning, and Andrew Y. Ng. Learning new facts from knowledge bases with neural tensor networks and semantic word vectors. In *Proc. ICLR*, 2013.

[13] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. In *Proc. ICLR*, 2015.

[14] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9 (8), 1997.

[15] Wojciech Zaremba, Ilya Sutskever, and Oriol Vinyals. Recurrent neural network regularization. In *Proc. ICLR*, 2015.

[16] Matthew D. Zeiler. ADADELTA: an adaptive learning rate method. *arXiv:1212.5701*, 2012.