# Compositional generalization in a deep seq2seq model by separating syntax and semantics

**Jake Russin**
Department of Psychology and Neuroscience
University of Colorado Boulder
jacob.russin@colorado.edu

**Jason Jo**
MILA
Université de Montréal

**Randall C. O'Reilly**
Department of Psychology and Neuroscience
University of Colorado Boulder

**Yoshua Bengio**
MILA, Université de Montréal
CIFAR Senior Fellow

## Abstract

Standard methods in deep learning for natural language processing fail to capture the compositional structure of human language that allows for systematic generalization outside of the training distribution. However, human learners readily generalize in this way, e.g. by applying known grammatical rules to novel words. Inspired by work in neuroscience suggesting separate brain systems for syntactic and semantic processing, we implement a modification to standard approaches in neural machine translation, imposing an analogous separation. The novel model, which we call Syntactic Attention, substantially outperforms standard methods in deep learning on the SCAN dataset, a compositional generalization task, without any hand-engineered features or additional supervision. Our work suggests that separating syntactic from semantic learning may be a useful heuristic for capturing compositional structure.

## 1 Introduction

A crucial property underlying the expressive power of human language is its systematicity [16; 9]: syntactic or grammatical rules allow arbitrary elements to be combined in novel ways, making the number of sentences possible in a language to be exponential in the number of its basic elements. Recent work has shown that standard deep learning methods in natural language processing fail to capture this important property: when tested on unseen combinations of known elements, state-of-the-art models fail to generalize [15; 17; 4]. It has been suggested that this failure represents a major deficiency of current deep learning models, especially when they are compared to human learners [19; 16].

A recently published dataset called SCAN [15] (**S**implified version of the **C**ommAI **N**avigation tasks), tests compositional generalization in a sequence-to-sequence (seq2seq) setting by systematically holding out of the training set all inputs containing a basic primitive verb ("jump"), and testing on sequences containing that verb. Success on this difficult problem requires models to generalize knowledge gained about the other primitive verbs ("walk", "run" and "look") to the novel verb "jump," without having seen "jump" in any but the most basic context ("jump" → JUMP). It is trivial for human learners to generalize in this way (e.g. if I tell you that "dax" is a verb, you can generalize its usage to all kinds of constructions, like "dax twice and then dax again", without even knowing what the word means) [15]. However, standard recurrent seq2seq models fail miserably on this task, with the best-reported model (a gated recurrent unit augmented with an attention mechanism) achieving
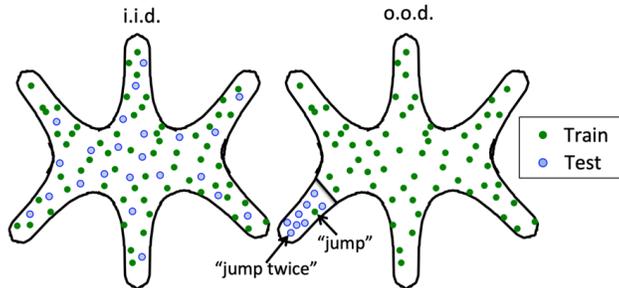
Figure 1: Simplified illustration of out-of-domain (o.o.d.) extrapolation required by SCAN compositional generalization task. Shapes represent the distribution of all possible command sequences. In a simple split, train and test data are independent and identically distributed (i.i.d.), but in the add-primitive splits, models are required to extrapolate out-of-domain from a single example.

only 12.5% accuracy on the test set [15; 4]. Recently, convolutional neural networks (CNN) were shown to perform better on this test, but still only achieved 69.2% accuracy on the test set.

From a statistical-learning perspective, this failure is quite natural. The neural networks trained on the SCAN task fail to generalize because they have memorized biases that do indeed exist in the training set. Because "jump" has never been seen with any adverb, it would not be irrational to assume that "jump twice" is an invalid sentence in this language. The SCAN task requires networks to make an inferential leap about the entire structure of part of the distribution that they have not seen - that is, it requires them to make an out-of-domain (o.o.d.) *extrapolation* [19], rather than merely *interpolate* according to the assumption that train and test data are independent and identically distributed (i.i.d.) (see Figure 1). Seen another way, the SCAN task and its analogues in human learning (e.g. "dax"), require models *not* to learn some of the correlations that are actually present in the training data [14].

Given that humans can perform well on certain kinds of o.o.d. extrapolation tasks, the human brain must be implementing principles that allow humans to generalize systematically, but which are lacking in current deep learning models. One prominent idea from neuroscience research on language processing that may offer such a principle is that the brain contains partially separate systems for processing syntax and semantics. In this paper, we motivate such a separation from a machine-learning perspective, and test a simple implementation on the SCAN dataset. Our novel model, which we call Syntactic Attention, encodes syntactic and semantic information in separate streams before producing output sequences. Our experiments show that our novel architecture achieves substantially improved compositional generalization performance over other recurrent networks on the SCAN dataset.

### 1.1 Syntax and prefrontal cortex

Syntax is the aspect of language underlying its systematicity [9]. When given a novel verb like "dax," humans can generalize its usage to many different constructions that they have never seen before, by applying known syntactic or grammatical rules about verbs (e.g. rules about how to conjugate to a different tense or about how adverbs modify verbs). It has long been thought that humans possess specialized cognitive machinery for learning the syntactic or grammatical structure of language [7]. A part of the prefrontal cortex called Broca's area, originally thought only to be involved in language production, was later found to be important for comprehending syntactically complex sentences, leading some to conclude that it is important for syntactic processing in general [6; 26]. For example, patients with lesions to this area showed poor comprehension on sentences such as "The girl that the boy is chasing is tall". Sentences such as this one require listeners to process syntactic information because semantics is not enough to understand their meanings - e.g. either the boy or the girl could be doing the chasing, and either could be tall.

A more nuanced view situates the functioning of Broca's area within the context of prefrontal cortex in general, noting that it may simply be a part of prefrontal cortex specialized for language [26]. The prefrontal cortex is known to be important for cognitive control, or the active maintenance of top-down attentional signals that bias processing in other areas of the brain [21] (see diagram on the
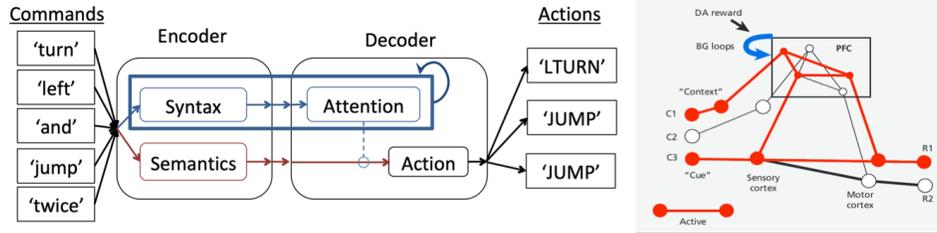
2

Figure 2: (left) Syntactic Attention architecture. Syntactic and semantic information are maintained in separate streams. The semantic stream processes words with a simple linear transformation, so that sequential information is not maintained. This information is used to directly produce actions. The syntactic stream processes inputs with a recurrent neural network, allowing it to capture temporal dependencies between words. This stream determines the attention over semantic representations at each time step during decoding. (right) Diagram of an influential computational model of prefrontal cortex (PFC) [21]. Prefrontal cortex dynamically modulates processes in other parts of the brain through top-down selective attention signals. A part of the prefrontal cortex, Broca's area, is thought to be important for syntactic processing [26]. Figure reproduced from [20].

right of Figure 2). In this framework, Broca's area can be thought of as a part of prefrontal cortex specialized for language, and responsible for selectively attending to linguistic representations housed in other areas of the brain [26].

The prefrontal cortex has received much attention from computational neuroscientists [21; 22], and one model even showed a capacity for compositional generalization [14]. However, these ideas have not been taken up in deep learning research. Here, we emphasize the idea that the brain contains two separate systems for processing syntax and semantics, where the semantic system learns and stores representations of the meanings of words, and the syntactic system, housed in Broca's area of the prefrontal cortex, learns how to selectively attend to these semantic representations according to grammatical rules.

## 2   Syntactic Attention

The Syntactic Attention model improves the compositional generalization capability of an existing attention mechanism [2] by implementing two separate streams of information processing for syntax and semantics (see Figure 2). Here, by "semantics" we mean the information in each word in the input that determines its *meaning* (in terms of target outputs), and by "syntax" we mean the information contained in the input sequence that should determine the *alignment* of input to target words. We describe the mechanisms of this separation and the other details of the model below, following the notation of [2], where possible.

### 2.1   Separation assumption

In the seq2seq problem, models must learn a mapping from arbitrary-length sequences of inputs $\mathbf{x} = \{x_1, x_2, ..., x_{T_x}\}$ to arbitrary-length sequences of outputs $\mathbf{y} = \{y_1, y_2, ..., y_{T_y}\}$: $p(\mathbf{y}|\mathbf{x})$. The attention mehcanism of [2] models the conditional probability of each target word given the input sequence and previous targets: $p(y_i|y_1, y_2, ..., y_{i-1}, \mathbf{x})$. This is accomplished by processing the input sequence with a recurrent neural network (RNN) in the encoder. The outputs of this RNN are used both for encoding individual words in the input for later translation, and for determining their alignment to targets during decoding.

The underlying assumption made by the Syntactic Attention architecture is that the dependence of target words on the input sequence can be separated into two independent factors. One factor, $p(y_i|x_j)$, which we refer to as "semantics," models the conditional distribution from individual words in the input to individual words in the target. Note that, unlike in the model of Bahdanau et al. [2], these $x_j$ do not contain any information about the other words in the input sequence because they are not processed with an RNN. They are "semantic" in the sense that they contain the information relevant to translating into the target language. The other factor, $p(j \rightarrow i|\mathbf{x})$, which we refer to as

3

"syntax," models the conditional probability that word $j$ in the input is relevant to word $i$ in the target sequence, given the entire input sequence. This alignment is accomplished from encodings of the inputs produced by an RNN. This factor is "syntactic" in the sense that it must capture all of the temporal information in the input that is relevant to determining the serial order of outputs. The crucial architectural assumption, then, is that any temporal dependency between individual words in the input that can be captured by an RNN should only be relevant to their alignment to words in the target sequence, and not to the translation of individual words. This assumption will be made clearer in the model description below.

## 2.2 Encoder

The encoder produces two separate vector representations for each word in the input sequence. Unlike the previous attention model [2]), we separately extract the semantic information from each word with a linear transformation:

$$m_j = W_m x_j, \tag{1}$$

where $W_m$ is a learned weight matrix that multiplies the one-hot encodings $\{x_1, ..., x_{T_x}\}$. Note that the semantic representation of each word does not contain any information about the other words in the sentence. As in the previous attention mechanism [2], we use a bidirectional RNN (biRNN) to extract what we now interpret as the syntactic information from each word in the input sequence. The biRNN produces a vector for each word on the forward pass, $(\overrightarrow{h_1}, ..., \overrightarrow{h_{T_x}})$, and a vector for each word on the backward pass, $(\overleftarrow{h_1}, ..., \overleftarrow{h_{T_x}})$. The syntactic information (or "annotations" [2]) of each word $x_j$ is determined by the two vectors $\overrightarrow{h_{j-1}}, \overleftarrow{h_{j+1}}$ corresponding to the words surrounding it:

$$h_j = [\overrightarrow{h_{j-1}}; \overleftarrow{h_{j+1}}] \tag{2}$$

In all experiments, we used a bidirectional Long Short-Term Memory (LSTM) for this purpose. Note that because there is no sequence information in the semantic representations, all of the information required to parse (i.e. align) the input sequence correctly (e.g. phrase structure, modifying relationships, etc.) must be encoded by the biRNN.

## 2.3 Decoder

The decoder models the conditional probability of each target word given the input and the previous targets: $p(y_i|y_1, y_2, ..., y_{i-1}, \mathbf{x})$, where $y_i$ is the target translation and $\mathbf{x}$ is the whole input sequence. As in the previous model, we use an RNN to determine an attention distribution over the inputs at each time step (i.e. to align words in the input to the current target). However, our decoder diverges from this model in that the mapping from inputs to outputs is performed from a weighted average of the *semantic* representations of the input words:

$$d_i = \sum_{j=1}^{T_x} \alpha_{ij} m_j \qquad p(y_i|y_1, y_2, ..., y_{i-1}, \mathbf{x}) = f(d_i) \tag{3}$$

where $f$ is parameterized by a linear function with a softmax nonlinearity, and the $\alpha_{ij}$ are the weights determined by the attention model. We note again that the $m_j$ are produced directly from corresponding $x_j$, and do not depend on the other inputs. The attention weights are computed by a function measuring how well the syntactic information of a given word in the input sequence aligns with the current hidden state of the decoder RNN, $s_i$:

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})} \qquad e_{ij} = a(s_i, h_j) \tag{4}$$

where $e_{ij}$ can be thought of as measuring the importance of a given input word $x_j$ to the current target word $y_i$, and $s_i$ is the current hidden state of the decoder RNN. Bahdanau et al. [2] model the function $a$ with a feedforward network, but following [11], we choose to use a simple dot product:

$$a(s_i, h_j) = s_i \cdot h_j, \tag{5}$$

4

relying on the end-to-end backpropagation during training to allow the model to learn to make appropriate use of this function. Finally, the hidden state of the RNN is updated with the same weighted combination of the *syntactic* representations of the inputs:

$$s_i = g(s_{i-1}, c_i) \qquad c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j \tag{6}$$

where $g$ is the decoder RNN, $s_i$ is the current hidden state, and $c_i$ can be thought of as the information in the attended words that can be used to determine what to attend to on the next time step. Again, in all experiments an LSTM was used.

## 3 Experiments

### 3.1 SCAN dataset

| | | |
|---|---|---|
| jump | ⇒ | JUMP |
| jump left | ⇒ | LTURN JUMP |
| jump around right | ⇒ | RTURN JUMP RTURN JUMP RTURN JUMP RTURN JUMP |
| turn left twice | ⇒ | LTURN LTURN |
| jump thrice | ⇒ | JUMP JUMP JUMP |
| jump opposite left and walk thrice | ⇒ | LTURN LTURN JUMP WALK WALK WALK |
| jump opposite left after walk around left | ⇒ | LTURN WALK LTURN WALK LTURN WALK LTURN WALK LTURN LTURN JUMP |

Figure 3: Examples from SCAN dataset. Figure reproduced from [15].

The SCAN[1] dataset is composed of sequences of commands that must be mapped to sequences of actions [15] (see Figure 3 and supplementary materials for further details). The dataset is generated from a simple finite phrase-structure grammar that includes things like adverbs and conjunctions. There are 20,910 total examples in the dataset that can be split systematically into training and testing sets in different ways. These splits include the following:

- Simple split: training and testing data are split randomly
- Length split: training includes only shorter sequences
- Add primitive split: a primitive command (e.g. "turn left" or "jump") is held out of the training set, except in its most basic form (e.g. "jump" → JUMP)

Here we focus on the most difficult problem in the SCAN dataset, the add-jump split, where "jump" is held out of the training set. The best test accuracy reported in the original paper [15], using standard seq2seq models, was 1.2%. More recent work has tested other kinds of seq2seq models, including Gated Recurrent Units (GRU) augmented with attention [4] and convolutional neural networks (CNNs) [8]. Here, we compare the Syntactic Attention model to the best previously reported results.

### 3.2 Implementation details

Experimental procedure is described in detail in the supplementary materials. Train and test sets were kept as they were in the original dataset, but following [4], we used early stopping by validating on a 20% held out sample of the training set. All reported results are from runs of 200,000 iterations with a batch size of 1. Unless stated otherwise, each architecture was trained 5 times with different random seeds for initialization, to measure variability in results. All experiments were implemented in PyTorch. Details of the hyperparameter search are given in supplementary materials. Our best model used LSTMs, with 2 layers and 200 hidden units in the encoder, and 1 layer and 400 hidden units in the decoder, and 120-dimensional semantic vectors. The model included a dropout rate of 0.5, and was optimized using an Adam optimizer [13] with a learning rate of 0.001.

---

[1]The SCAN dataset can be downloaded at `https://github.com/brendenlake/SCAN`

## 3.3 Results

The Syntactic Attention model achieves state-of-the-art performance on the key compositional generalization task of the SCAN dataset (see table 1). The table shows results (mean test accuracy (%) $\pm$ standard deviation) on the test splits of the dataset. Syntactic Attention is compared to the previous best models, which were a CNN [8], and GRUs augmented with an attention mechanism ("+ attn"), which either included or did not include a dependency ("- dep") in the decoder on the previous action [4]. The best model from the hyperparameter search showed strong compositional generalization performance, attaining a mean accuracy of 91.1% (median = 98.5%) on the test set of the add-jump split. However, as in Dessì and Baroni [8], we found that our model showed variance across initialization seeds. We suggest that this may be due to the nature of the add-jump split: since "jump" has only been encountered in the simplest context, it may be that slight changes to the way that this verb is encoded can make big differences when models are tested on more complicated constructions. For this reason, we ran the best model 25 times on the add-jump split to get a more accurate assessment of performance. These results were highly skewed, with a mean accuracy of 78.4 % but a median of 91.0 % (see supplementary materials for detailed results). Overall, this represents an improvement over the best previously reported results on this task [4; 8], and does so without any hand-engineered features or additional supervision.

Table 1: Compositional generalization results. The Syntactic Attention model achieves an improvement on the compositional generalization tasks of the SCAN dataset, compared to the best previously reported models [4; 8]. Star[*] indicates median of 25 runs.

| Model | Simple | Length | Add turn left | Add jump |
|---|---|---|---|---|
| GRU + attn [4] | $100.0 \pm 0.0$ | $18.1 \pm 1.1$ | $59.1 \pm 16.8$ | $12.5 \pm 6.6$ |
| GRU + attn - dep [4] | $100.0 \pm 0.0$ | $17.8 \pm 1.7$ | $90.8 \pm 3.6$ | $0.7 \pm 0.4$ |
| CNN [8] | $100.0 \pm 0.0$ | - | - | $69.2 \pm 8.2$ |
| Syntactic Attention | $100.0 \pm 0.0$ | $15.2 \pm 0.7$ | $\mathbf{99.9 \pm 0.16}$ | $\mathbf{91.0^{*} \pm 27.4}$ |

## 3.4 Additional experiments

To test our hypothesis that compositional generalization requires a separation between syntax (i.e. sequential information used for alignment), and semantics (i.e. the mapping from individual source words to individual targets), we conducted two more experiments:

- *Sequential semantics*. An additional biLSTM was used to process the semantics of the sentence: $m_j = [\overrightarrow{m_j}; \overleftarrow{m_j}]$, where $\overrightarrow{m_j}$ and $\overleftarrow{m_j}$ are the vectors produced for the source word $x_j$ by a biLSTM on the forward and backward passes, respectively. These $m_j$ replace those generated by the simple linear layer in the Syntactic Attention model (in equation (1)).

- *Syntax-action*. Syntactic information was allowed to directly influence the output at each time step in the decoder: $p(y_i|y_1, y_2, ..., y_{i-1}, \mathbf{x}) = f([d_i; c_i])$, where again $f$ is parameterized with a linear function and a softmax output nonlinearity.

The results of the additional experiments (mean test accuracy (%) $\pm$ standard deviations) are shown in table 2. These results partially confirmed our hypothesis: performance on the jump-split test set was worse when the strict separation between syntax and semantics was violated by allowing sequential information to be processed in the semantic stream. However, "syntax-action," which included sequential information produced by a biLSTM (in the syntactic stream) in the final production of actions, maintained good compositional generalization performance. We hypothesize that this was because in this setup, it was easier for the model to learn to use the semantic information to directly translate actions, so it largely ignored the syntactic information. This experiment suggests that the separation between syntax and semantics does not have to be perfectly strict, as long as non-sequential semantic representations are available for direct translation.

## 4 Dicussion

The Syntactic Attention model was designed to incorporate a key principle that has been hypothesized to describe the organization of the linguistic brain: mechanisms for learning rule-like or syntactic

Table 2: Results of additional experiments. Star[*] indicates median of 25 runs.

| Model | Simple | Length | Add turn left | Add jump |
|---|---|---|---|---|
| *Sequential semantics* | $99.3 \pm 0.7$ | $13.1 \pm 2.5$ | $99.4 \pm 1.1$ | $42.3 \pm 32.7$ |
| *Syntax-action* | $99.3 \pm 0.85$ | $15.2 \pm 1.9$ | $98.2 \pm 2.2$ | $88.7 \pm 14.2$ |
| Syntactic Attention | $100.0 \pm 0.0$ | $15.2 \pm 0.7$ | $99.9 \pm 0.16$ | $91.0^{*} \pm 27.4$ |

information are separated from mechanisms for learning semantic information. Our experiments confirm that this simple organizational principle encourages systematicity in recurrent neural networks in the seq2seq setting, as shown by the substantial improvement in the model's performance on the compositional generalization tasks in the SCAN dataset.

The model makes the assumption that the translation of individual words in the input should be independent of their alignment to words in the target sequence. To this end, two separate encodings are produced for the words in the input: semantic representations in which each word is not influenced by other words in the sentence, and syntactic representations which are produced by an RNN that can capture temporal dependencies in the input sequence (e.g. modifying relationships, binding to grammatical roles). Just as Broca's area of the prefrontal cortex is thought to play a role in syntactic processing through a dynamic selective-attention mechanism that biases processing in other areas of the brain, the syntactic system in our model encodes serial information and is constrained to influence outputs through an attention mechanism alone.

Patients with lesions to Broca's area are able to comprehend sentences like "The girl is kicking a green ball", where semantics can be used to infer the grammatical roles of the words (e.g. that the girl, not the ball, is doing the kicking) [6]. However, these patients struggle with sentences such as "The girl that the boy is chasing is tall", where the sequential order of the words, rather than semantics, must be used to infer grammatical roles (e.g. either the boy or the girl could be doing the chasing). In our model, the syntactic stream can be seen as analogous to Broca's area, because without it the model would not be able to learn about the temporal dependencies that determine the grammatical roles of words in the input.

The separation of semantics and syntax, which is in the end a *constraint*, forces the model to learn, in a relatively independent fashion, 1) the individual meanings of words and 2) how the words are being used in a sentence (e.g. how they can modify one another, what grammatical role each is playing, etc.). This encourages systematic generalization because, even if a word has only been encountered in a single context (e.g. "jump" in the add-jump split), as long as its syntactic role is known (e.g. that it is a verb that can be modified by adverbs such as "twice"), it can be used in many other constructions that follow the rules for that syntactic role (see supplementary materials for visualizations). Additional experiments confirmed this intuition, showing that when sequential information is allowed to be processed by the semantic system ("sequential semantics"), systematic generalization performance is substantially reduced.

The Syntactic Attention model bears some resemblance to a symbolic system - the paradigm example of systematicity - in the following sense: in symbolic systems, representational content (e.g. the value of a variable stored in memory) is maintained separately from the computations that are performed on that content. This separation ensures that the *manipulation* of the content stored in variables is fairly independent of the content itself, and will therefore generalize to arbitrary elements. Our model implements an analogous separation, but in a purely neural architecture that does not rely on hand-coded rules or additional supervision. In this way, it can be seen as transforming a difficult out-of-domain (o.o.d.) generalization problem into two separate i.i.d. generalization problems - one where the individual meanings of words are learned, and one where *how words are used* (e.g. how adverbs modify verbs) is learned (see Figure 4).

It is unlikely that the human brain has such a strict separation between semantic and syntactic processing, and in the end, there must be more of an interaction between the two streams. We expect that the separation between syntax and semantics in the brain is only a relative one, but we have shown here that this kind of separation can be useful for encouraging systematicity and allowing for compositional generalization.
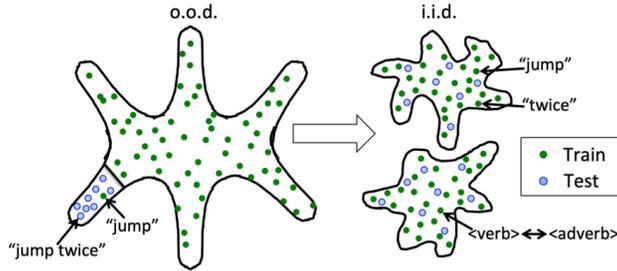
Figure 4: Illustration of the transformation of an out-of-domain (o.o.d.) generalization problem into two independent, identically distributed (i.i.d.) generalization problems. This transformation is accomplished by the Syntactic Attention model without hand-coding grammatical rules or supervising with additional information such as parts-of-speech tags.

## 5 Other related work

Our model integrates ideas from computational and cognitive neuroscience [26; 22; 14; 21], into the neural machine translation framework. Much of the work in neural machine translation uses an encoder-decoder framework, where one RNN is used to encode the source sentence, and then a decoder neural network decodes the representations given by the encoder to produce the words in the target sentence [25]. Earlier work attempted to encode the source sentence into a single fixed-length vector (the final hidden state of the encoder RNN), but it was subsequently shown that better performance could be achieved by encoding each word in the source, and using an attention mechanism to *align* these encodings with each target word during the decoding process [2]. The current work builds directly on this attention model, while incorporating a separation between syntactic and semantic information streams.

The principle of compositionality has recently regained the attention of deep learning researchers [1; 3; 16; 15; 5; 12] . In particular, the issue has been explored in the visual-question answering (VQA) setting [1; 11; 12; 23; 10; 24; 27]. Many of the successful models in this setting learn hand-coded operations [1; 10], use highly specialized components [11; 24], or use additional supervision [10; 27]. In contrast, our model uses standard recurrent networks and simply imposes the additional constraint that syntactic and semantic information are processed in separate streams.

Some of the recent research on compositionality in machine learning has had a special focus on the use of attention. For example, in the Compositional Attention Network, built for VQA, a strict separation is maintained between the representations used to encode images and the representations used to encode questions [11]. This separation is enforced by restricting them to interact only through attention distributions. Our model utilizes a similar restriction, reinforcing the idea that compositionality is enhanced when information from different modalities (in our case syntax and semantics) are only allowed to interact through discrete probability distributions.

Previous research on compositionality in machine learning has also focused on the incorporation of symbol-like processing into deep learning models [1; 10; 27]. These methods generally rely on hand-coding or additional supervision for the symbolic representations or algorithmic processes to emerge. For example, in neural module networks [1], a neural network is constructed out of composable neural modules that each learn a specific operation. These networks have shown an impressive capacity for systematic generalization on VQA tasks [3]. These models can be seen as accomplishing a similar transformation as depicted in Figure 4, because the learning in each module is somewhat independent of the mechanism that composes them. However, Bahdanau et al. [3] find that when these networks are trained end-to-end (i.e. without hand-coded parameterizations and layouts) their systematicity is significantly degraded.

In contrast, our model learns in an end-to-end way to generalize systematically without any explicit symbolic processes built in. This offers an alternative way in which symbol-like processing can be achieved with neural networks - by enforcing a separation between mechanisms for learning representational content (semantics) and mechanisms for learning how to dynamically attend to or manipulate that content (syntax) in the context of a cognitive operation or reasoning problem.

## 6  Conclusion

The Syntactic Attention model incorporates ideas from cognitive and computational neuroscience into the neural machine translation framework, and produces the kind of systematic generalization thought to be a key component of human language-learning and intelligence. The key feature of the architecture is the separation of sequential information used for alignment (syntax) from information used for mapping individual inputs to outputs (semantics). This separation allows the model to generalize the usage of a word with known syntax to many of its valid grammatical constructions. This principle may be a useful heuristic in other natural language processing tasks, and in other systematic or compositional generalization tasks. The success of our approach suggests a conceptual link between dynamic selective-attention mechanisms in the prefrontal cortex and the systematicity of human cognition, and points to the untapped potential of incorporating ideas from cognitive science and neuroscience into modern approaches in deep learning and artificial intelligence [18].

## References

[1] J. Andreas, M. Rohrbach, T. Darrell, and D. Klein. Neural Module Networks. *arXiv:1511.02799 [cs]*, Nov. 2015.

[2] D. Bahdanau, K. Cho, and Y. Bengio. Neural Machine Translation by Jointly Learning to Align and Translate. *arXiv:1409.0473 [cs, stat]*, Sept. 2014.

[3] D. Bahdanau, S. Murty, M. Noukhovitch, T. H. Nguyen, H. de Vries, and A. Courville. Systematic Generalization: What Is Required and Can It Be Learned? *arXiv:1811.12889 [cs]*, Nov. 2018.

[4] J. Bastings, M. Baroni, J. Weston, K. Cho, and D. Kiela. Jump to better conclusions: SCAN both left and right. *arXiv:1809.04640 [cs]*, Sept. 2018.

[5] P. W. Battaglia, J. B. Hamrick, V. Bapst, A. Sanchez-Gonzalez, V. Zambaldi, M. Malinowski, A. Tacchetti, D. Raposo, A. Santoro, R. Faulkner, C. Gulcehre, F. Song, A. Ballard, J. Gilmer, G. Dahl, A. Vaswani, K. Allen, C. Nash, V. Langston, C. Dyer, N. Heess, D. Wierstra, P. Kohli, M. Botvinick, O. Vinyals, Y. Li, and R. Pascanu. Relational inductive biases, deep learning, and graph networks. *arXiv:1806.01261 [cs, stat]*, June 2018.

[6] A. Caramazza and E. B. Zurif. Dissociation of algorithmic and heuristic processes in language comprehension: Evidence from aphasia. *Brain and Language*, 3(4):572–582, Oct. 1976. ISSN 0093-934X. doi: 10.1016/0093-934X(76)90048-1.

[7] N. Chomsky, editor. *Syntactic Structures*. Mouton & Co., The Hague, Jan. 1957.

[8] R. Dessì and M. Baroni. CNNs found to jump around more skillfully than RNNs: Compositional generalization in seq2seq convolutional networks. *arXiv:1905.08527 [cs]*, May 2019.

[9] J. A. Fodor and Z. W. Pylyshyn. Connectionism and cognitive architecture: A critical analysis. *Cognition*, 28(1-2):3–71, Apr. 1988.

[10] R. Hu, J. Andreas, M. Rohrbach, T. Darrell, and K. Saenko. Learning to Reason: End-to-End Module Networks for Visual Question Answering. *arXiv:1704.05526 [cs]*, Apr. 2017.

[11] D. A. Hudson and C. D. Manning. Compositional attention networks for machine reasoning. *arXiv:1803.03067 [cs]*, Mar. 2018.

[12] J. Johnson, B. Hariharan, L. van der Maaten, L. Fei-Fei, C. L. Zitnick, and R. Girshick. CLEVR: A Diagnostic Dataset for Compositional Language and Elementary Visual Reasoning. *arXiv:1612.06890 [cs]*, Dec. 2016.

[13] D. P. Kingma and J. Ba. Adam: A Method for Stochastic Optimization. *arXiv:1412.6980 [cs]*, Dec. 2014.

[14] T. Kriete, D. C. Noelle, J. D. Cohen, and R. C. O'Reilly. Indirection and symbol-like processing in the prefrontal cortex and basal ganglia. *Proceedings of the National Academy of Sciences*, 110 (41):16390–16395, Oct. 2013. ISSN 0027-8424, 1091-6490. doi: 10.1073/pnas.1303547110.

[15] B. M. Lake and M. Baroni. Generalization without systematicity: On the compositional skills of sequence-to-sequence recurrent networks. *arXiv:1711.00350 [cs]*, Oct. 2017.

[16] B. M. Lake, T. D. Ullman, J. B. Tenenbaum, and S. J. Gershman. Building machines that learn and think like people. *The Behavioral and Brain Sciences*, 40:e253, Jan. 2017. ISSN 1469-1825. doi: 10.1017/S0140525X16001837.

[17] J. Loula, M. Baroni, and B. M. Lake. Rearranging the familiar: Testing compositional generalization in recurrent networks. *arXiv:1807.07545 [cs]*, July 2018.

[18] A. H. Marblestone, G. Wayne, and K. P. Kording. Toward an Integration of Deep Learning and Neuroscience. *Frontiers in Computational Neuroscience*, 10, 2016. ISSN 1662-5188. doi: 10.3389/fncom.2016.00094.

[19] G. Marcus. Deep learning: A critical appraisal. Jan. 2018.

[20] E. K. Miller. The "working" of working memory. *Dialogues in Clinical Neuroscience*, 15(4): 411–418, Dec. 2013. ISSN 1294-8322.

[21] E. K. Miller and J. D. Cohen. An integrative theory of prefrontal cortex function. *Annual Review of Neuroscience*, 24:167–202, 2001.

[22] R. C. O'Reilly and M. J. Frank. Making working memory work: A computational model of learning in the prefrontal cortex and basal ganglia. *Neural Computation*, 18(2):283–328, 2006.

[23] E. Perez, F. Strub, H. de Vries, V. Dumoulin, and A. Courville. FiLM: Visual Reasoning with a General Conditioning Layer. *arXiv:1709.07871 [cs, stat]*, Sept. 2017.

[24] A. Santoro, D. Raposo, D. G. T. Barrett, M. Malinowski, R. Pascanu, P. Battaglia, and T. Lillicrap. A simple neural network module for relational reasoning. *arXiv:1706.01427 [cs]*, June 2017.

[25] I. Sutskever, O. Vinyals, and Q. V. Le. Sequence to Sequence Learning with Neural Networks. *arXiv:1409.3215 [cs]*, Sept. 2014.

[26] S. L. Thompson-Schill. Dissecting the language organ : A new look at the role of Broca ' s area in language processing. 2004.

[27] K. Yi, J. Wu, C. Gan, A. Torralba, P. Kohli, and J. B. Tenenbaum. Neural-Symbolic VQA: Disentangling Reasoning from Vision and Language Understanding. *arXiv:1810.02338 [cs]*, Oct. 2018.

# 7 Supplementary materials

## 7.1 SCAN dataset details

The SCAN dataset [15] generates sequences of commands using the pharase-structure grammar described in Figure 5. This simple grammar is not recursive, and so can generate a finite number of command sequences (20,910 total).

$$
\begin{array}{lll}
C \rightarrow S \text{ and } S & V \rightarrow D[1] \text{ opposite } D[2] & D \rightarrow \text{turn left} \\
C \rightarrow S \text{ after } S & V \rightarrow D[1] \text{ around } D[2] & D \rightarrow \text{turn right} \\
C \rightarrow S & V \rightarrow D & U \rightarrow \text{walk} \\
S \rightarrow V \text{ twice} & V \rightarrow U & U \rightarrow \text{look} \\
S \rightarrow V \text{ thrice} & D \rightarrow U \text{ left} & U \rightarrow \text{run} \\
S \rightarrow V & D \rightarrow U \text{ right} & U \rightarrow \text{jump}
\end{array}
$$

Figure 5: Phrase-structure grammar used to generate SCAN dataset. Figure reproduced from [15].

These commands are interpreted according to the rules shown in Figure 6. Although the grammar used to generate and interpret the commands is simple compared to any natural language, it captures the basic properties that are important for testing compositionality (e.g. modifying relationships, discrete grammatical roles, etc.). The add-primitive splits (described in main text) are meant to be analogous to the capacity of humans to generalize the usage of a novel verb (e.g. "dax") to many constructions [15].

$$
\begin{array}{ll}
[\![\text{walk}]\!] = \text{WALK} & [\![u \text{ opposite left}]\!] = [\![\text{turn opposite left}]\!] \ [\![u]\!] \\
[\![\text{look}]\!] = \text{LOOK} & [\![u \text{ opposite right}]\!] = [\![\text{turn opposite right}]\!] \ [\![u]\!] \\
[\![\text{run}]\!] = \text{RUN} & [\![\text{turn around left}]\!] = \text{LTURN LTURN LTURN LTURN} \\
[\![\text{jump}]\!] = \text{JUMP} & [\![\text{turn around right}]\!] = \text{RTURN RTURN RTURN RTURN} \\
[\![\text{turn left}]\!] = \text{LTURN} & [\![u \text{ around left}]\!] = \text{LTURN} \ [\![u]\!] \ \text{LTURN} \ [\![u]\!] \ \text{LTURN} \ [\![u]\!] \ \text{LTURN} \ [\![u]\!] \\
[\![\text{turn right}]\!] = \text{RTURN} & [\![u \text{ around right}]\!] = \text{RTURN} \ [\![u]\!] \ \text{RTURN} \ [\![u]\!] \ \text{RTURN} \ [\![u]\!] \ \text{RTURN} \ [\![u]\!] \\
[\![u \text{ left}]\!] = \text{LTURN} \ [\![u]\!] & [\![x \text{ twice}]\!] = [\![x]\!] \ [\![x]\!] \\
[\![u \text{ right}]\!] = \text{RTURN} \ [\![u]\!] & [\![x \text{ thrice}]\!] = [\![x]\!] \ [\![x]\!] \ [\![x]\!] \\
[\![\text{turn opposite left}]\!] = \text{LTURN LTURN} & [\![x_1 \text{ and } x_2]\!] = [\![x_1]\!] \ [\![x_2]\!] \\
[\![\text{turn opposite right}]\!] = \text{RTURN RTURN} & [\![x_1 \text{ after } x_2]\!] = [\![x_2]\!] \ [\![x_1]\!]
\end{array}
$$

Figure 6: Rules for interpreting command sequences to generate actions in SCAN dataset. Figure reproduced from [15].

## 7.2 Experimental procedure details

The cluster used for all experiments consists of 3 nodes, with 68 cores in total (48 times Intel(R) Xeon(R) CPU E5-2650 v4 at 2.20GHz, 20 times Intel(R) Xeon(R) CPU E5-2650 v3 at 2.30GHz), with 128GB of ram each, connected through a 56Gbit infiniband network. It has 8 pascal Titan X GPUs and runs Ubuntu 16.04.

All experiments were conducted with the SCAN dataset as it was originally published [15]. No data were excluded, and no preprocessing was done except to encode words in the input and action sequences into one-hot vectors, and to add special tokens for start-of-sequence and end-of-sequence tokens. Train and test sets were kept as they were in the original dataset, but following [4], we used early stopping by validating on a 20% held out sample of the training set. All reported results are from runs of 200,000 iterations with a batch size of 1. Except for the additional batch of 25 runs for the add-jump split, each architecture was trained 5 times with different random seeds for initialization, to measure variability in results. All experiments were implemented in PyTorch.

Initial experimentation included different implementations of the assumption that syntactic information be separated from semantic information. After the architecture described in the main text showed promising results, a hyperparameter search was conducted to determine optimization (stochastic gradient descent vs. Adam), RNN-type (GRU vs. LSTM), regularizers (dropout, weight decay), and

number of layers (1 vs. 2 layers for encoder and decoder RNNs). We found that the Adam optimizer [13] with a learning rate of 0.001, two layers in the encoder RNN and 1 layer in the decoder RNN, and dropout worked the best, so all further experiments used these specifications. Then, a grid-search was conducted to find the number of hidden units (in both semantic and syntactic streams) and dropout rate. We tried hidden dimensions ranging from 50 to 400, and dropout rates ranging from 0.0 to 0.5.

The best model used an LSTM with 2 layers and 200 hidden units in the encoder, and an LSTM with 1 layer and 400 hidden units in the decoder, and used 120-dimensional semantic vectors, and a dropout rate of 0.5. The results for this model are reported in the main text. All additional experiments were done with models derived from this one, with the same hyperparameter settings.

All evaluation runs are reported in the main text: for each evaluation except for the add-jump split, models were trained 5 times with different random seeds, and performance was measured with means and standard deviations of accuracy. For the add-jump split, we included 25 runs to get a more accurate assessment of performance. This revealed a strong skew in the distribution of results, so we included the median as the main measure of performance. Occasionally, the model did not train at all due to an unknown error (possibly very poor random initialization, high learning rate or numerical error). For this reason, we excluded runs in which training accuracy did not get above 10%. No other runs were excluded.

## 7.3   Skew of add-jump results

As mentioned in the results section of the main text, we found that test accuracy on the add-jump split was variable and highly skewed. Figure 7 shows a histogram of these results (proportion correct). The model performs near-perfectly most of the time, but is also prone to catastrophic failures. This may be because, at least for our model, the add-jump split represents a highly nonlinear problem in the sense that slight differences in the way the primitive verb "jump" is encoded during training can have huge differences for how the model performs on more complicated constructions. We recommend that future experiments with this kind of compositional generalization problem take note of this phenomenon, and conduct especially comprehensive analyses of variability in results. Future research will also be needed to better understand the factors that determine this variability, and whether it can be overcome with other priors or regularization techniques.
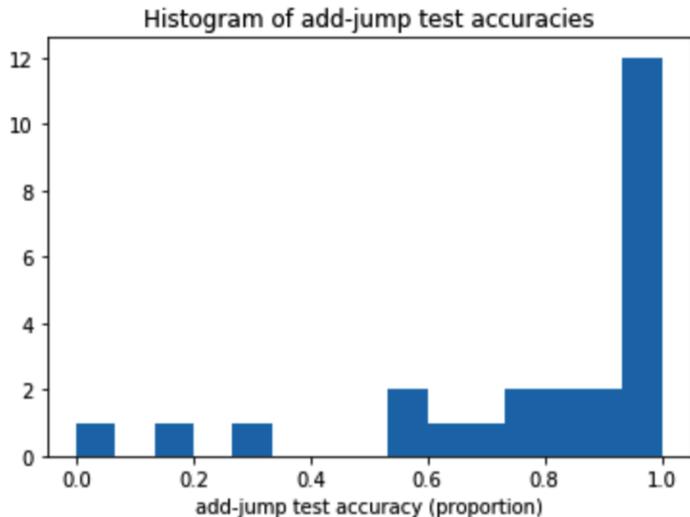


Figure 7: Histogram of test accuracies across all 25 runs of add-jump split.

### 7.4 Supplementary experiments

#### 7.4.1 Testing nonlinear semantics

Our main hypothesis is that the separation between sequential information used for alignment (syntax) and information about the meanings of individual words (semantics) encourages systematicity. The results reported in the main text are largely consistent with this hypothesis, as shown by the performance of the Syntactic Attention model on the composotional generalization tests of the SCAN dataset. However, it is also possible that the simplicity of the semantic stream in the model is also important for improving compositional generalization. To test this, we replaced the linear layer in the semantic stream with a nonlinear neural network. From the model description in the main text:

$$p(y_i|y_1, y_2, ..., y_{i-1}, \mathbf{x}) = f(d_i), \tag{7}$$

In the original model, $f$ was parameterized with a simple linear layer, but here we use a two-layer feedforward network with a ReLU nonlinearity, before a softmax is applied to generate a distribution over the possible actions. We tested this model on the add-primitive splits of the SCAN dataset. The results (mean (%) with standard deviations) are shown in Table 3, with comparison to the baseline Syntactic Attention model.

Table 3: Results of nonlinear semantics experiment. Star[*] indicates median of 25 runs.

| Model | Add turn left | Add jump |
|---|---|---|
| *Nonlinear semantics* | $99.0 \pm 1.7$ | $84.4 \pm 14.1$ |
| Syntactic Attention | $99.9 \pm 0.16$ | $91.0^* \pm 27.4$ |

The results show that this modification did not substantially degrade compositional generalization performance, suggesting that the success of the Syntactic Attention model does not depend on the parameterization of the semantic stream with a simple linear function.

#### 7.4.2 Add-jump split with additional examples

The original SCAN dataset was published with compositional generalization splits that have more than one example of the held-out primitive verb [15]. The training sets in these splits of the dataset include 1, 2, 4, 8, 16, or 32 random samples of command sequences with the "jump" command, allowing for a more fine-grained measurement of the ability to generalize the usage of a primitive verb from few examples. For each number of "jump" commands included in the training set, five different random samples were taken to capture any variance in results due to the selection of particular commands to train on.

Lake and Baroni [15] found that their best model (an LSTM without an attention mechansim) did not generalize well (below 39%), even when it was trained on 8 random examples that included the "jump" command, but that the addition of further examples to the training set improved performance. Subsequent work showed better performance at lower numbers of "jump" examples, with GRU's augmented with an attention mechanism ("+ attn"), and either with or without a dependence in the decoder on the previous target ("- dep") [4]. Here, we compare the Syntactic Attention model to these results.

The Syntactic Attention model shows a substantial improvement over previously reported results at the lowest numbers of "jump" examples used for training (see Figure 8 and Table 4). Compositional generalization performance is already quite high at 1 example, and at 2 examples is almost perfect (99.997% correct).

#### 7.4.3 Template splits

The compositional generalization splits of the SCAN dataset were originally designed to test for the ability to generalize known primitive verbs to valid unseen constructions [15]. Further work with SCAN augmented this set of tests to include compositional generalization based not on known verbs but on known *templates* [17]. These template splits included the following (see Figure 9 for examples):

- *Jump around right*: All command sequences with the phrase "jump around right" are held out of the training set and subsequently tested.
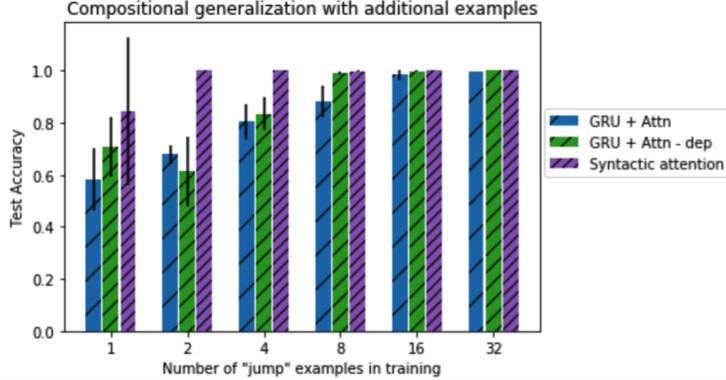
Figure 8: Compositional generalization performance on add-jump split with additional examples. Syntactic Attention model is compared to previously reported models [4] on test accuracy as command sequences with "jump" are added to the training set. Mean accuracy (proportion correct) was computed with 5 different random samples of "jump" commands. Error bars represent standard deviations.

Table 4: Results of Syntactic Attention compared to models of Bastings et al. [4] on jump-split with additional examples. Mean accuracy (% - rounded to tenths) is shown with standard deviations. Same data as depicted in Figure 8.

| Model | Number of jump commands in training set | | | | | |
|---|---|---|---|---|---|---|
| | 1 | 2 | 4 | 8 | 16 | 32 |
| GRU + attn | $58.2_{\pm12.0}$ | $67.8_{\pm3.4}$ | $80.3_{\pm7.0}$ | $88.0_{\pm6.0}$ | $98.3_{\pm1.8}$ | $99.6_{\pm0.2}$ |
| GRU + attn - dep | $70.9_{\pm11.5}$ | $61.3_{\pm13.5}$ | $83.5_{\pm6.1}$ | $99.0_{\pm0.4}$ | $99.7_{\pm0.2}$ | $100.0_{\pm0.0}$ |
| Syntactic Attention | $84.4_{\pm28.5}$ | $100.0_{\pm0.01}$ | $100.0_{\pm0.02}$ | $99.9_{\pm0.2}$ | $100.0_{\pm0.01}$ | $99.9_{\pm0.2}$ |

- Primitive *right*: All command sequences containing primitive verbs modified by "right" are held out of the training set and subsequently tested.
- Primitive *opposite right*: All command sequences containing primitive verbs modified by "opposite right" are held out of the training set and subsequently tested.
- Primitive *around right*: All command sequences containing primitive verbs modified by "around right" are held out of the training set and subsequently tested.

| Condition | Example train commands | Example test commands |
|---|---|---|
| *jump around right* | *"jump left"*, *"jump around left"*, *"walk around right"* | *"jump around right"*, *"jump around right and walk"* |
| Primitive *right* | *"jump left"*, *"walk around right"* | *"jump right"*, *"walk right"* |
| Primitive *opposite right* | *"jump left"*, *"jump opposite left"*, *"walk right"* | *"jump opposite right"*, *"walk opposite right"* |
| Primitive *around right* | *"jump left"*, *"jump around left"*, *"walk right"* | *"jump around right"*, *"walk around right"* |

Figure 9: Table of example command sequences for each template split. Reproduced from [17]
.

Results of the Syntactic Attention model on these template splits are compared to those originally published [17] in Table 5. The model, like the one reported in [17], performs well on the *jump around right* split, consistent with the idea that this task does not present a problem for neural networks. The rest of the results are mixed: Syntactic Attention shows good compositional generalization performance on the Primitive *right* split, but fails on the Primitive *opposite right* and Primitive *around right* splits. All of the template tasks require models to generalize based on the symmetry between

"left" and "right" in the dataset. However, in the *opposite right* and *around right* splits, this symmetry is substantially violated, as one of the two prepositional phrases in which they can occur is never seen with "right." Further research is required to determine whether a model implementing similar principles to Syntactic Attention can perform well on this task.

Table 5: Results of Syntactic Attention compared to models of Loula et al. [17] on template splits of SCAN dataset. Mean accuracy (%) is shown with standard deviations. **P** = Primitive

| Model | Template split | | | |
|---|---|---|---|---|
| | *jump around right* | **P** *right* | **P** *opposite right* | **P** *around right* |
| LSTM (Loula et al. [17]) | 98.43±0.54 | 23.49±8.09 | 47.62±17.72 | 2.46±2.68 |
| Syntactic Attention | 98.9±2.3 | 99.1±1.8 | 10.5±8.8 | 28.9±34.8 |

## 7.5 Visualizing attention

The way that the attention mechanism of Bahdanau et al. [2] is set up allows for easy visualization of the model's attention. Here, we visualize the attention distributions over the words in the command sequence at each step during the decoding process. In the following figures (Figures 10 to 15), the attention weights on each command (in the columns of the image) is shown for each of the model's outputs (in the rows of the image) for some illustrative examples. Darker blue indicates a higher weight. The examples are shown in pairs for a model trained and tested on the add-jump split, with one example drawn from the training set and a corresponding example drawn from the test set. Examples are shown in increasing complexity, with a failure mode depicted in Figure 15.

In general, it can be seen that although the attention distributions on the test examples are not exactly the same as those from the corresponding training examples, they are usually good enough for the model to produce the correct action sequence. This shows the model's ability to apply the same syntactic rules it learned on the other verbs to the novel verb "jump." In the example shown in Figure 15, the model fails to attend to the correct sequence of commands, resulting in an error.
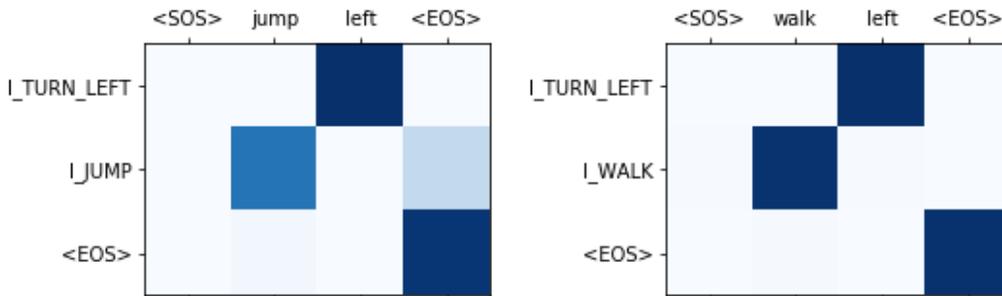


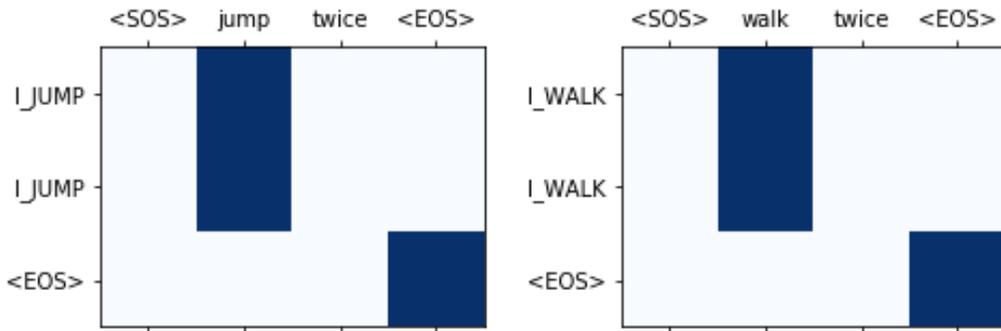Figure 10: Attention distributions: correct example

Figure 11: Attention distributions: correct example
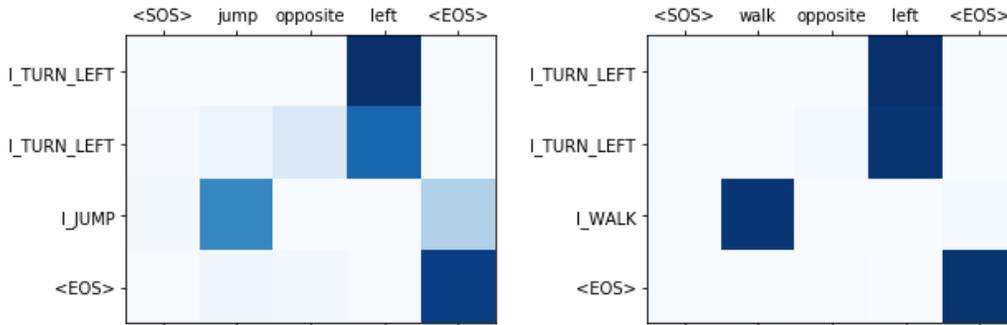


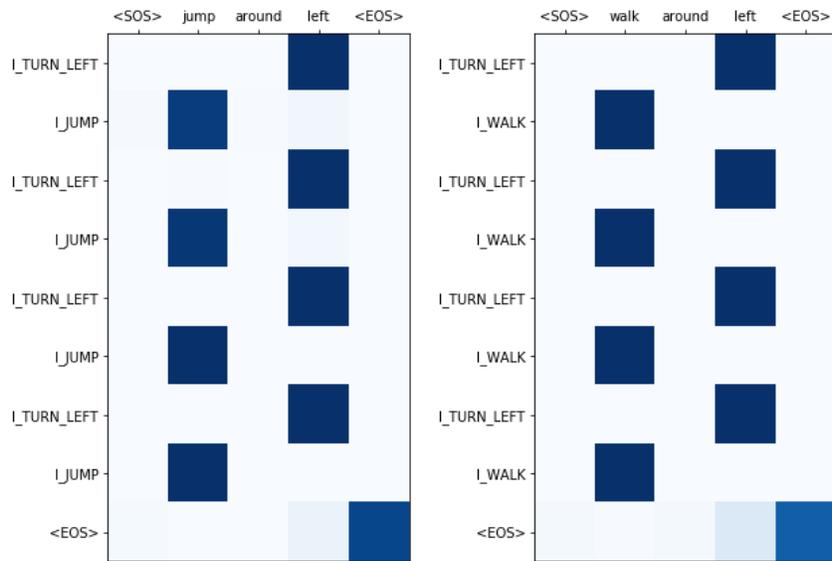Figure 12: Attention distributions: correct example



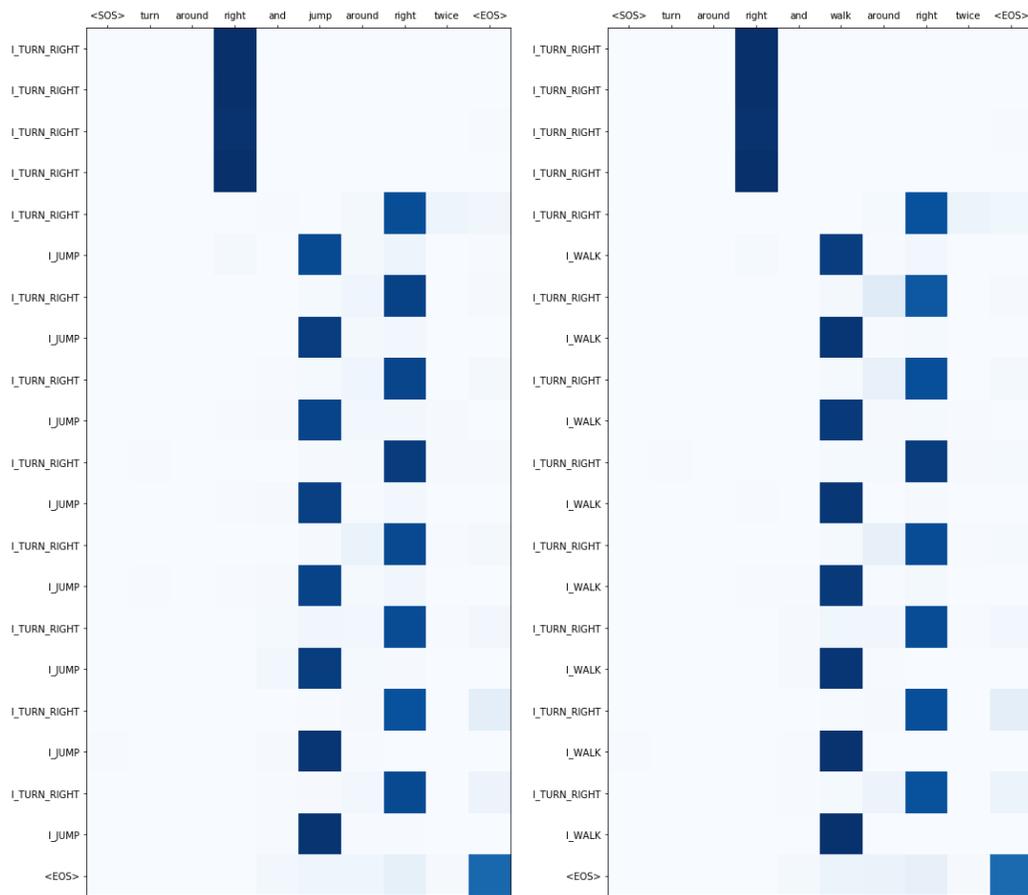Figure 13: Attention distributions: correct example
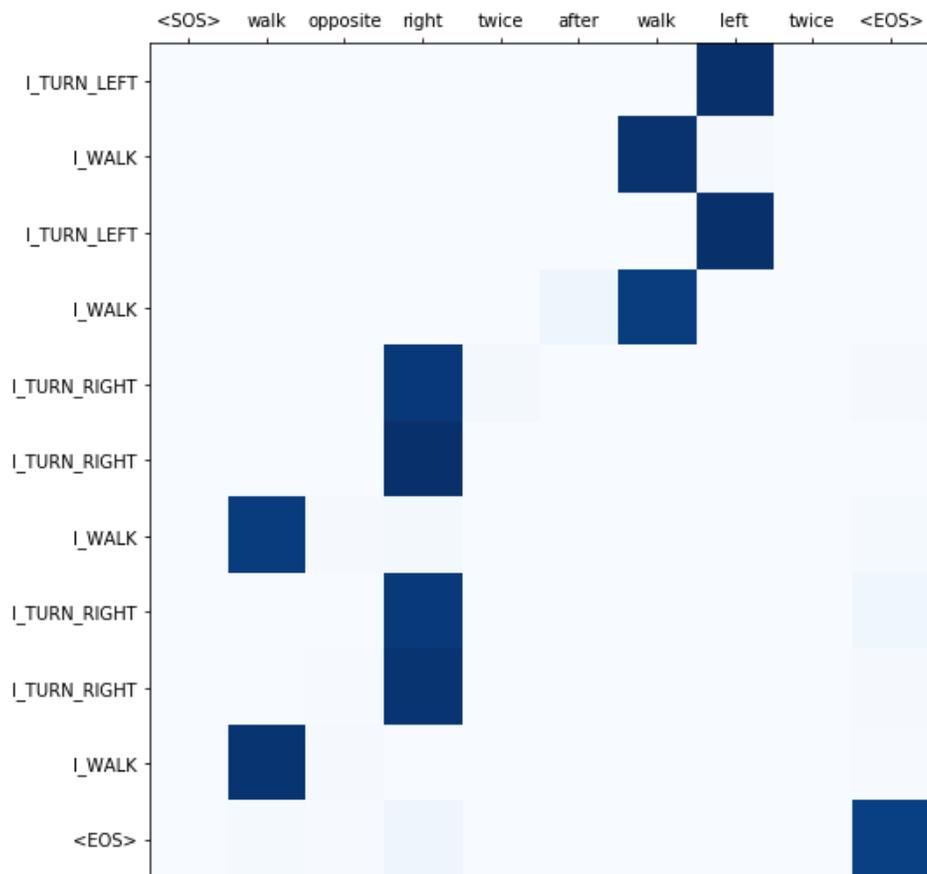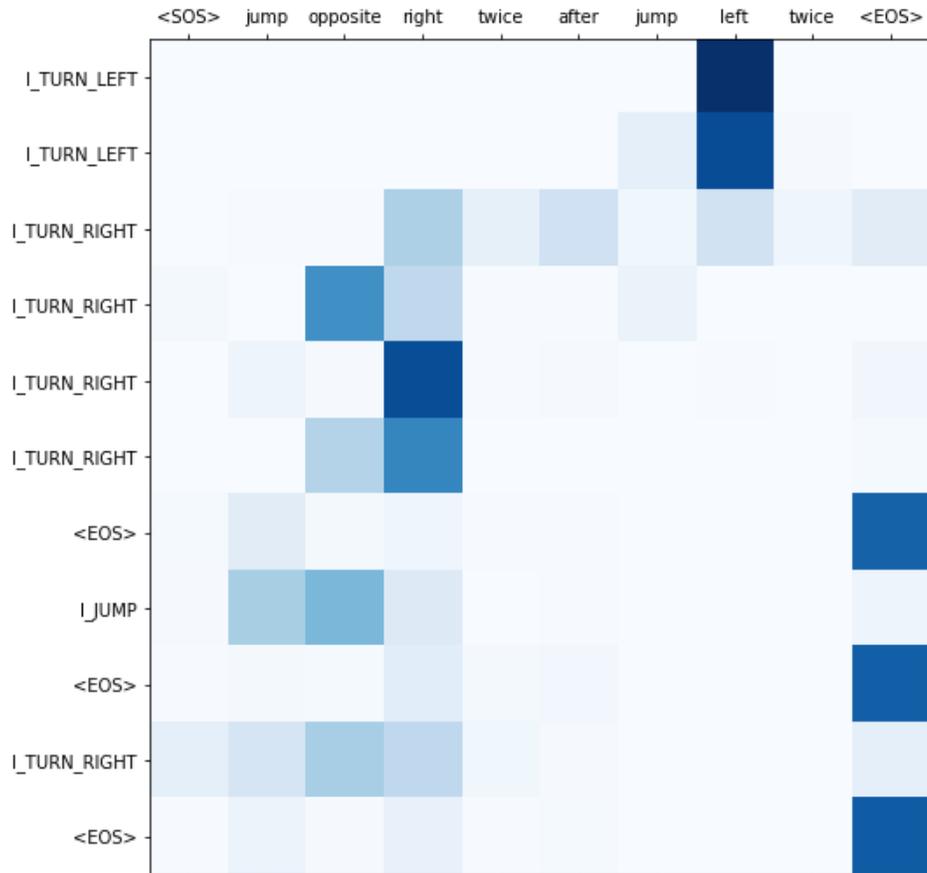
Figure 14: Attention distributions: correct example

Figure 15: Attention distributions: incorrect example