

Kris Korrel, Dieuwke Hupkes, Verna Dankers, Elia Bruni



1. With specifically designed sequence-to-sequence tasks, the human-like compositional understanding is tested for standard seq2seq deep learning models.
2. A new design of the seq2seq model is created that should improve the compositional skills on the selected tasks by putting more emphasize on the attention module.
3. (i) is that all and (ii) are these objectives or contributions?

Seq2seq models have become ubiquitous in the field of machine translation, language modeling, speech recognition and other tasks which can be casted into temporally-dependent sequences with possibly varying input and output lengths. Although their generalization capabilities in these fields might indicate an understanding of the rules and hierarchies that underlie the tasks, specialized tasks/synonym that test specifically for the **compositional understanding** of these systems have shown evidence that this is not the case. We take these specialized tasks(same synonym) as a basis for understanding where seq2seq models lack human-like generalization skills, and (i) extend the testing methodologies for compositional understanding by examining **overgeneralization** capabilities, and (ii) propose a new architecture with a specialized focus on its attention module, which we dub **seq2attn**, with the aim to improve compositional understanding.

The diagram illustrates a sequence-to-sequence model architecture. On the left, the **Encoder** consists of three pink boxes labeled  $E_0$ ,  $E_1$ , and  $E_2$ , which process **Input embeddings** (Lorum, ipsum, dolor). On the right, the **Transcoder** consists of two green boxes labeled  $T_0$  and  $T_1$ . The **Decoder** consists of two purple boxes labeled  $D_0$  and  $D_1$ , which produce **Output symbols** (<eos>, Nonsense, <eos>). The **Context vector** layer consists of two yellow circles labeled  $C_0$  and  $C_1$ . Arrows show the flow of information: solid blue arrows from encoder states to transcoder states, solid green arrows from transcoder states to decoder states, and dotted blue and green arrows representing various connections between the layers.

In summary, the proposed seq2attn model could be explained as a combination of four adaptations of the standard seq2seq architecture with attention.

- |          | held-out<br>inputs               | held-out<br>compositions         | held-out<br>tables               |
|----------|----------------------------------|----------------------------------|----------------------------------|
| Baseline | 38.25 $\pm$ 0.04                 | 43.28 $\pm$ 0.09                 | 7.86 $\pm$ 0.02                  |
| Seq2attn | <b>100 <math>\pm</math> 0.00</b> | <b>100 <math>\pm</math> 0.00</b> | <b>100 <math>\pm</math> 0.00</b> |

A bar chart comparing the sequence accuracy of two models, 'look around left' (solid green bars) and 'look around right' (red hatched bars), across different numbers of training examples. The x-axis represents the 'Number of training examples containing look around right' with values 1, 2, 4, 8, 16, 32, 64, 128, 256, 512, and 1024. The y-axis represents 'sequence accuracy' from 0 to 100. The 'look around left' model consistently shows higher accuracy than the 'look around right' model, especially at lower example counts. Both models show an upward trend in accuracy as the number of training examples increases, with the 'look around right' model's accuracy converging towards the 'look around left' model's accuracy at higher example counts.

Number of training examples containing look around right	look around left (sequence accuracy)	look around right (sequence accuracy)
1	~43	~3
2	~34	~4
4	~55	~11
8	~66	~18
16	~82	~38
32	~79	~50
64	~78	~64
128	~84	~69
256	~83	~73
512	~88	~88
1024	~74	~93

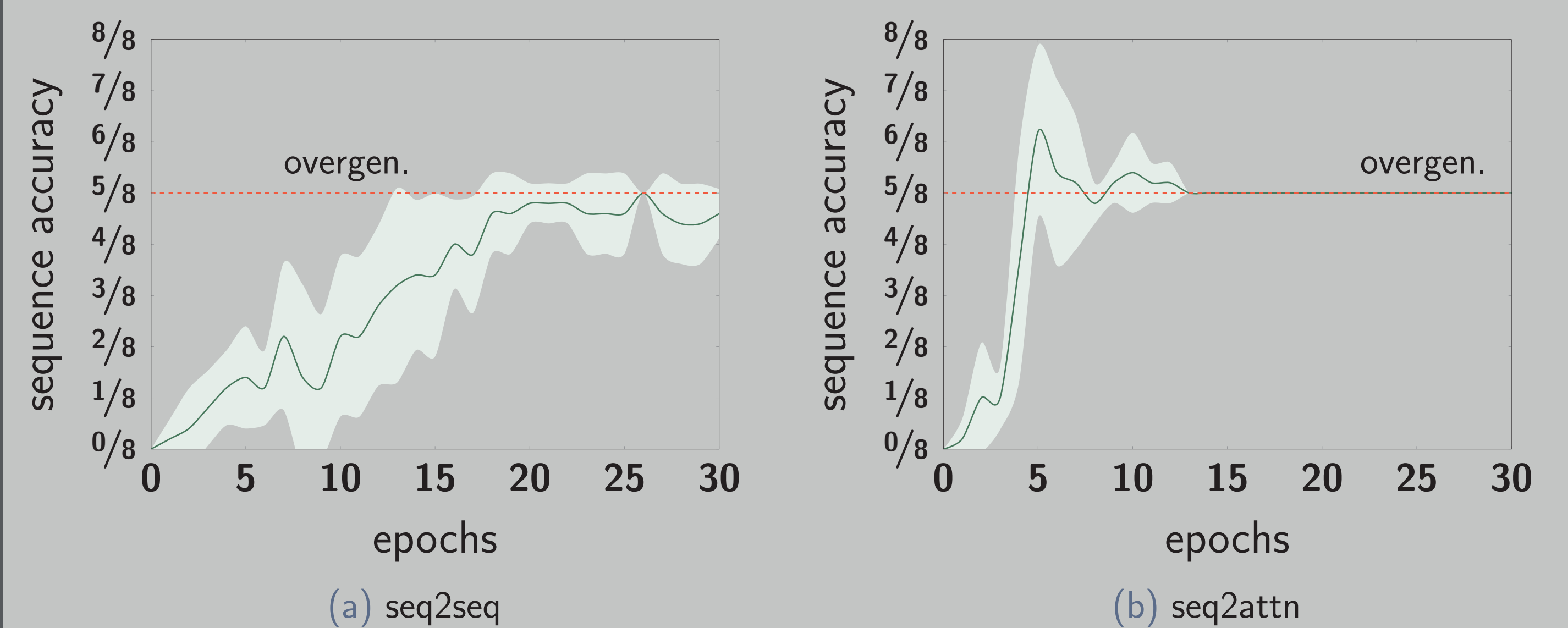
Figure 1 displays attention weights for three models: (a) seq2seq, (b) seq2seq, and (c) seq2attn. Each model's attention weights are visualized as a 4x4 grid. The columns represent time steps  $t_0, t_1, t_2, t_3$  and the rows represent input tokens 000, 101, 000, and  $\langle \text{eos} \rangle$ . The attention weights are color-coded: white for high attention, black for low attention, and red for the highest attention. In (a) and (b), the highest attention is concentrated on the diagonal elements (e.g., 000 to  $t_0$ , 101 to  $t_1$ , 000 to  $t_2$ , and  $\langle \text{eos} \rangle$  to  $t_3$ ). In (c), the attention is more distributed, with high attention also appearing on off-diagonal elements (e.g., 000 to  $t_1$ , 101 to  $t_2$ , 000 to  $t_3$ , and  $\langle \text{eos} \rangle$  to  $t_0$ ).

(a) seq2seq

(b) seq2attn

(c) seq2attn

We suspect that compositional generalization leads to more *overgeneralization*. This is a phenomenon where a learner generalizes rules even to exceptions where these rules do not apply. For the lookup tables task, we test this by making the output of  $x \text{ } t_1 \text{ } t_2$  random for 3 of the 8 bit-strings.



- ▶ Lookup Tables
  - ▷ The input sequences consist in a three-bit string followed by a series of function identifiers for lookup tables. The objective is to apply these bijective functions over the (intermediate) outputs successively
  - ▷ Example: 001 t1 t2 → 001 010 111
- ▶ SCAN
  - ▷ The input consists in one or two sub-sequences which can be joined by the conjunctives and or after. Each sub-sequence contains a primitive command and possibly modifiers that act on this command. The learner must interpret the commands and mentally apply them in a 2D game grid.
  - ▷ Example: jump after walk left twice → I\_TURN LEFT I\_WALK I\_TURN LEFT I\_WALK I\_JUMP

- ▶ increased general performance on specialized tasks indicate improved compositional generalization by seq2attn.
- ▶ The general ideas implemented could be extended to several variants of encoder-decoder architectures.
- ▶ Sparser attentional patterns improve interpretability of the found solutions,
- ▶ Preliminary results show neither increased or decreased performance on NMT (without the Straight-Through estimator of Gumbel-Softmax)
- ▶ Strictly sparse attention vectors might limit the expressiveness of the model. Possible improvements may be made by research on alternatives methods.