

Pervasive Attention 2D Convolutional Neural Networks for Sequence-to-Sequence Prediction Maha Elbayad[†] Laurent Besacier[‡] Jakob Verbeek[†] firstname.lastname@inria.fr

firstname.lastname@univ-grenoble-alpes.fr

UNIVERSITÉ Grenoble

Code available at: github.com/elbayadm/attn2d

Overview

In state-of-the-art encoder-decoder models, the source and target sequences are processed **separately**. The decoder, equipped with an **attention mechanism**, focuses on different parts of the source at each decoding step. However, the attention is limited to assigning weights to the **once and for all** computed encoder states.

Contributions: build an architecture from the get go around attention by **jointly** encoding the source and target sequences and allowing for different source representations for every target position.

Pervasive attention: the aggregation



To aggregate activations across source positions e.g. with $H_3 = [h_{31}^L, ..., h_{3|x|}^L] \in \mathbb{R}^{d \times |x|}$, we can use:

- ► Max/average pooling.
- ► Self-attention:
 - $\rho = softmax(H_3^TW + b).$ $h_{3} = H_{3}\rho$.
- A combination of the above.

Encoder-Decoders

Encoder

Inputs: source sequence $\mathbf{x} = (x_1, x_2, \dots, x_{|\mathbf{x}|})$. Depending on the chosen architecture, the encoder computes the source representations.



Decoder

Inputs: source codes $(s_1, \ldots, s_{|\mathbf{x}|})$ and target sequence $\mathbf{y} = (y_1, y_2, \ldots, y_{|\mathbf{y}|})$. At every step t:

- Under the architecture, compute the hidden state h_t causally.
- Given the new state, the attention mechanism yields a context c_t .
- ▶ $h_t := \operatorname{combine}(h_t, c_t).$

Pervasive attention: the input



Experimental results

Pre-processing:

Lower-casing.

Benchmark: IWSLT'14 German \leftrightarrow English translation.



		#-pinis		#-pinis
RNNsearch (Bahdanau et al., 2015)	29.98	13M	25.04	15M
Varational attention (Deng et al., 2018)	33.10	-	_	_
ConvS2S (MLE) (Gehring et al., 2017)	31.59	21M	27.18	22M
ConvS2S (MLE+SLE) (Edunov et al., 2018)	32.84	-	_	-
Transformer (Vaswani et al., 2017), V1	34.42	46M	28.23	48M
Transformer (Vaswani et al., 2017), V2	34.44	52M	28.07	52M
Pervasive attention (ours), V1	33.86	11M	27.21	11M
Pervasive attention (ours), V2	34.05	22M	27.97	22M

models we trained using either our implementation or fairseq. 10 averaged checkpoints.

The initial 2D grid: $\forall i, j : \{1, ..., |y|\} \times \{1, ..., |x|\}$

 $= U_{v_i}$ (target embedding) U_i $= V_{x_i}$ (source embedding) V_i $= concat(u_i, v_i)$ h_{ii}

Pervasive attention: the convolutional network

- **Causality:** with masked filters in the target direction.
- **Context:** grown with stacked convolutions.
- **Padding:** throughout the network to maintain source/target resolution.



Alignment visualization



BLEU per sequence length





Using a single-block DenseNet (Huang et al., 2017) with L layers. Each layer grows its input channels by **g**.

References

Bahdanau, D., Cho, K., and Bengio, Y. (2015). Neural machine translation by jointly learning to align and translate. In ICLR. Deng, Y., Kim, Y., Chiu, J., Guo, D., and Rush, A. (2018). Latent alignment and variational attention. arXiv preprint arXiv:1807.03756. Edunov, S., Ott, M., Auli, M., Grangier, D., and Ranzato, M. (2018). Classical structured prediction losses for sequence to sequence learning. In NAACL. Gehring, J., Auli, M., Grangier, D., Yarats, D., and Dauphin, Y. (2017). Convolutional sequence to sequence learning. In ICML Huang, G., Liu, Z., van der Maaten, L., and Weinberger, K. (2017). Densely connected convolutional networks. In CVPR. Sennrich, R., Haddow, B., and Birch, A. (2016). Neural machine translation of rare words with subword units. In ACL. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A., Kaiser, L., and Polosukhin, I. (2017). Attention is all you need. In NIPS.



Due to memory/compute limitations, $\mathcal{O}(|x|,|y|)$ instead of $\mathcal{O}(|y| + |x|)$, we truncate sequences longer than 80 tokens when training which affects the performance on long sequences.

Takeaways

We have competitive results with:

- A sequence-to-sequence model outside of the encoder-decoder paradigm. A convolutional architecture, proving they work well for NLP problems. An implicit attention via re-encoding the source sequence then simply maxpooling the representations.
- ► Less parameters (at least 1/2 compared to transformer).