

FORMATTING INSTRUCTIONS FOR ICLR 2020 CONFERENCE SUBMISSIONS

Anonymous authors

Paper under double-blind review

ABSTRACT

The abstract paragraph should be indented 1/2 inch (3 picas) on both left and right-hand margins. Use 10 point type, with a vertical spacing of 11 points. The word ABSTRACT must be centered, in small caps, and in point size 12. Two line spaces precede the abstract. The abstract must be limited to one paragraph.

1 INTRODUCTION

Attention mechanism were introduced in Neural Machine Translation to better handle long range dependencies and align pairs of sentences (Bahdanau et al., 2015).

2 BACKGROUND ON ATTENTION MECHANISM AND RELATED WORK

We briefly recall the formulation of standard CNN layers and Self-Attention layers using unified notation.

2.1 CONVOLUTION LAYER

A convolution layer is defined by its kernel size K (assuming square kernel for simplicity), the number of input channel C_{in} and the number of output channel C_{out} . It is parametrized by a kernel tensor $\mathbf{W} \in \mathbb{R}^{K \times K \times C_{out} \times C_{in}}$ and a bias vector $\mathbf{b} \in \mathbb{R}^{C_{out}}$. Given an input image $\mathbf{X} \in \mathbb{R}^{W \times H \times C}$ of width W , height H and C channels, the output of the convolution layer is given by,

$$\mathbf{Y}_{i,j,:} = \mathbf{b} + \sum_{(m,n) \in \{[(K-1)/2], \dots, \lfloor K/2 \rfloor\}^2} \mathbf{W}_{m,n,:} \mathbf{X}_{i+m,j+n,:} \quad (1)$$

2.2 SELF-ATTENTION LAYER

A Self-Attention layer is defined by its hidden dimension d_h and its number of heads K . It is parametrized by $\mathbf{W}^Q \in \mathbb{R}^{K \times d \times d_h}$ query layer, $\mathbf{W}^V \in \mathbb{R}^{K \times d \times d_h}$ key layer, $\mathbf{W}^V \in \mathbb{R}^{K \times d \times d_h}$ value layer, a projection layer $\mathbf{W}^Q \in \mathbb{R}^{K \times d_h \times d}$ and output layer $\mathbf{W}^O \in \mathbb{R}^{(K d_h) \times d}$. Applied to an input $\mathbf{X} \in \mathbb{R}^{T \times d}$, the output of the Self-Attention layer is computed as,

$$\mathbf{A}_k = \text{softmax} \left((\mathbf{X} \mathbf{W}_k^K) (\mathbf{X} \mathbf{W}_k^Q)^\top \right), \quad (2)$$

$$\mathbf{H}_{t,:} = \sum_{k \in 1, \dots, K} \mathbf{A}_k \mathbf{X} \mathbf{W}_k^P, \quad (3)$$

$$\mathbf{Y}_{t,:} = \mathbf{H}_{t,:} \mathbf{W}^O + \mathbf{X}_{t,:}, \quad (4)$$

where the $\text{softmax}(\cdot)$ is taken over all the dimensions of the tensor but the first one. For simplicity, we exclude the batch normalization layers.

JB: same as size of the CNN kernel, we point that they are related (K^2) but need clarify notation

JB: check dimensions and readability

2.3 ATTENTION MECHANISM IN VISION

Successful work to use attention on images (Ramachandran et al., 2019; Bello et al., 2019).

Model	relative	sinusoids	learned
Vaswani et al. (2017)		✓	
Radford et al. (2018)			✓
Devlin et al. (2018)			✓
Dai et al. (2019)	✓	✓	
Yang et al. (2019)	✓	✓	
Bello et al. (2019)	✓		✓
Ramachandran et al. (2019)	✓		✓

Table 1: Type of position encoding used by transformers models applied to text (*top*) and images (*bottom*). When multiple encoding types have been tried, we report the one advised by the authors.

3 ATTENTION LAYER CAN IMPLEMENT CONVOLUTION LAYER

Our goal is to show that attention layers have the expressive power to learn convolution filter and that we can learn such filters in practice. The ability of the Attention mechanism to encode spacial filters (3×3 kernels) depends heavily on how we encode position in the image.

Understanding how attention can apply convolutions can help to:

- find better position encoding scheme,
- explain the need of multiple heads.

We take inspiration from (Dai et al., 2019, TransformerXL) to define relative positional encoding. They decompose the computation of the attention coefficients as follow:

$$\mathbf{A}_{i,j}^{\text{abs}} = (\mathbf{E}_{x_i} + \mathbf{U}_i)^\top \mathbf{W}_q^\top \mathbf{W}_k (\mathbf{E}_{x_j} + \mathbf{U}_j) \quad (5)$$

$$= \underbrace{\mathbf{E}_{x_i}^\top \mathbf{W}_q^\top \mathbf{W}_k \mathbf{E}_{x_j}}_{(a)} + \underbrace{\mathbf{E}_{x_i}^\top \mathbf{W}_q^\top \mathbf{W}_k \mathbf{U}_j}_{(b)} + \underbrace{\mathbf{U}_i^\top \mathbf{W}_q^\top \mathbf{W}_k \mathbf{E}_{x_j}}_{(c)} + \underbrace{\mathbf{U}_i^\top \mathbf{W}_q^\top \mathbf{W}_k \mathbf{U}_j}_{(d)} \quad (6)$$

with \mathbf{E}_{x_i} the embedding of token x_i , \mathbf{U}_i the absolute position encoding, \mathbf{R}_{i-j} relative position encoding. They replace all absolute position encoding with relative ones:

$$\mathbf{A}_{i,j}^{\text{rel}} = \underbrace{\mathbf{E}_{x_i}^\top \mathbf{W}_q^\top \mathbf{W}_{k,E} \mathbf{E}_{x_j}}_{(a)} + \underbrace{\mathbf{E}_{x_i}^\top \mathbf{W}_q^\top \mathbf{W}_{k,R} \mathbf{R}_{i-j}}_{(b)} + \underbrace{\mathbf{U}_i^\top \mathbf{W}_q^\top \mathbf{W}_{k,E} \mathbf{E}_{x_j}}_{(c)} + \underbrace{\mathbf{U}_i^\top \mathbf{W}_q^\top \mathbf{W}_{k,R} \mathbf{R}_{i-j}}_{(d)} \quad (7)$$

To mimic the CNN computation, each attention head focuses on a pixel at a given relative position. Hence reproducing a 3×3 kernel requires 9 attention heads and the ability for each of them to perfectly attend on one pixel at a relative position. Given that spatial convolution filters are not conditioned on the input data, we set $\mathbf{W}_{k,E}$ and $\mathbf{W}_{k,R}$ to 0, leaving only the (d) term in equation 7. We further write the attention computation in the 2D setting,

$$\mathbf{A}_{i,j}^{\text{conv}} = \mathbf{v}^\top \mathbf{W}_{k,R} \mathbf{R}_{\mathbf{p}_i - \mathbf{p}_j} \quad (8)$$

with \mathbf{p}_i the row and column position of pixel i . As both \mathbf{v} and $\mathbf{W}_{k,R}$ are learnable parameters, we simplify the expression to learn a single vector \mathbf{t}_h for each head h .

3.1 ATTENTION ON A PIXEL AT ANY RELATIVE POSITION

How to discern relative position? The goal is to reproduce the 3×3 pattern of CNN kernel with attention instead. Each head should be able to chose which pixel to attend based on its relative position to the query token x_i . The most basic attention mechanism is the dot product $(\mathbf{p}_i - \mathbf{p}_j)^\top \mathbf{t}_h$ between the relative positions of the input pixels and the relative target positions $\mathbf{t}_h = (x_h, y_h)^\top$. The dot product is large for pixels in the same direction as \mathbf{t}_h , however it also grows for farther pixels. To be able to focus on a relative pixel position (direction and distance), we propose the following relative position encoding

$$\mathbf{R}_\Delta = (\Delta_x \quad \Delta_y \quad \Delta_x^2 \quad \Delta_y^2 \quad 1 \quad 1)^\top \quad (9)$$

Instead of learning the target vector \mathbf{t}_h , we parametrize it for each head to attend to the relative pixel position (x_h, y_h) ,

$$\mathbf{t}_h = -\alpha_h \begin{pmatrix} -2x_h & -2y_h & 1 & 1 & x_h^2 & y_h^2 \end{pmatrix}^\top \quad (10)$$

The attention coefficient for head h between pixel i and j is given by the dot product,

$$\mathbf{R}_\Delta^\top \mathbf{t}_h = -\alpha_h ((\Delta_x - x_h)^2 + (\Delta_y - y_h)^2), \quad (11)$$

with $\Delta = \mathbf{p}_i - \mathbf{p}_j$. The maximum attention coefficient is 0 when $\Delta = (x_h, y_h)$, i.e. the center of attention of a given head. The α coefficient controls how spiky the attention is (analogous to temperature in the softmax). Figure 1 gives a representation of the attention weights for a fixed query pixel i and different α and (x_h, y_h) .

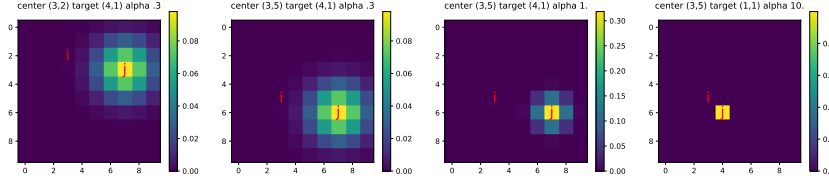


Figure 1: Attention coefficient for different center i and target j with varying α parameters. Target designate the relative position that the head attends to.

Discussion. Some relations can be drawn to RBF kernels centered at the target pixels. We are expressing a more general form of convolutions where (i) the input pixels do not follow a grid shape but has any pattern and (ii) the inputs are not single pixels but weighted averages of local patches in the image.

Simplification. One further simplification is to remove the constant terms of \mathbf{R}_Δ because the softmax (renormalization) is shift invariant and constant terms for a head are useless.

3.2 NON ISOTROPIC GAUSSIAN

In this section, we present an extension of the relative positional encoding introduced above to model non-isotropic Gaussians.

$$f_{\mu, \Sigma}(\mathbf{x}) = \exp \left(-\frac{1}{2} (\mathbf{x} - \mu)^\top \Sigma^{-1} (\mathbf{x} - \mu) \right) \quad (12)$$

We drop the normalization factor of the Gaussian distribution as it is replaced by the softmax of the attention mechanism. The matrix Σ^{-1} must be positive semi-definite. We parametrize it as

$$\Sigma^{-1} = (\Sigma^{-1/2})^\top \Sigma^{-1/2} = \begin{pmatrix} a & b \\ b & c \end{pmatrix} \quad (13)$$

Our goal is to write the exponent as a dot product of two vectors:

- one vector $R_{\mathbf{x}}$ dependent only on \mathbf{x} , the relative positional encoding,
- one vector \mathbf{t}_h dependent only on μ and Σ , the attention head.

$$R_{\mathbf{x}} = (x_1 \quad x_2 \quad x_1 \quad x_2 \quad x_1^2 \quad x_2^2 \quad x_1 x_2 \quad 1 \quad 1 \quad 1) \quad (14)$$

$$\mathbf{t}_h = -\frac{1}{2} (-2a\mu_1 \quad -2c\mu_2 \quad -2b\mu_2 \quad -2b\mu_1 \quad a \quad c \quad 2b \quad a\mu_1^2 \quad c\mu_2^2 \quad 2b\mu_1\mu_2) \quad (15)$$

$$R_{\mathbf{x}}^\top \mathbf{t}_h = -\frac{1}{2} [a(x_1 - \mu_1)^2 + 2b(x_1 - \mu_1)(x_2 - \mu_2) + c(x_2 - \mu_2)^2] = -\frac{1}{2} (\mathbf{x} - \mu)^\top \Sigma^{-1} (\mathbf{x} - \mu) \quad (16)$$

We can check that setting $b = 0$ and $a = c = \alpha$ recovers the expression of the isotropic Gaussian derived above.

Model	Number of parameters
ResNet10	4.9M
BERT 4 layers	1.85M
BERT 6 layers	9.56M

Table 2: Number of parameters per model.

3.3 IMPLICATIONS TO SINUSOID OR LEARNED POSITIONAL ENCODING

3.4 EXPLOIT ATTENTION ON FEATURES?

- Does the terms (a) (b) (c) allow to condition the CNN filters on the input data?

4 EXPERIMENTS

We want to validate that not only Transformer architecture can express CNN filters but that it can also learn such filters with SGD from the data. It has already been shown by (Bello et al., 2019) that adding attention features to CNNs can improve performance on Cifar-100 and ImageNet. We want to stick to a fully attentional model, and show that it matches its siblings ResNet of equivalent depth.

4.1 PERFORMANCE AGAINST RESNET

We design a smaller ResNet (i.e. ResNet10) which have similar number of layers as our transformer. The sizes of the models are displayed in Table 2.

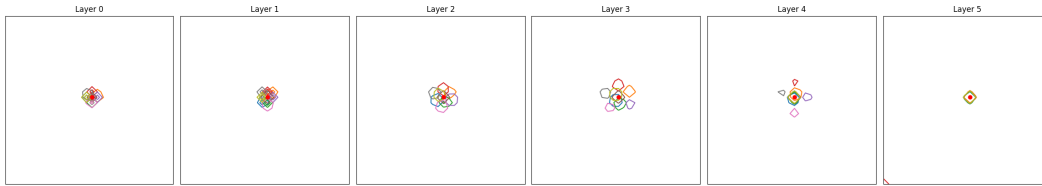


Figure 2: Contours of attention weights per layer for the 6 layers fully attentional model.

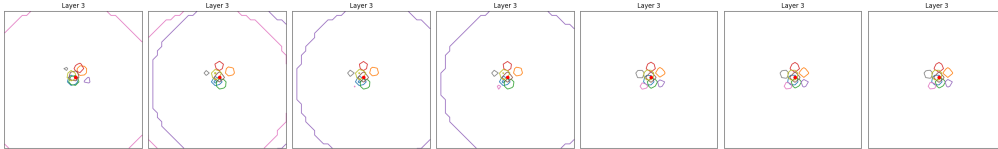


Figure 3: Contours of attention weights during training (300 epochs) at layer 3.

5 DISCUSSION

ACKNOWLEDGMENTS

Use unnumbered third level headings for the acknowledgments. All acknowledgments, including those to funding agencies, go at the end of the paper.

REFERENCES

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. In *3rd International Conference on Learning Representations*,



Figure 4: Evaluation accuracy on CIFAR-10 of a small ResNet and two Gaussian Attention models.

ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings, 2015. URL <http://arxiv.org/abs/1409.0473>.

Irwan Bello, Barret Zoph, Ashish Vaswani, Jonathon Shlens, and Quoc V. Le. Attention Augmented Convolutional Networks. *arXiv:1904.09925 [cs]*, April 2019.

Zihang Dai, Zhilin Yang, Yiming Yang, Jaime G. Carbonell, Quoc V. Le, and Ruslan Salakhutdinov. Transformer-XL: Attentive Language Models Beyond a Fixed-Length Context. *CoRR*, abs/1901.02860, 2019.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. *CoRR*, abs/1810.04805, 2018. URL <http://arxiv.org/abs/1810.04805>.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. 2018. URL <https://d4mucfpksywv.cloudfront.net/better-language-models/language-models.pdf>.

Prajit Ramachandran, Niki Parmar, Ashish Vaswani, Irwan Bello, Anselm Levskaya, and Jonathon Shlens. Stand-alone self-attention in vision models. *CoRR*, abs/1906.05909, 2019. URL <http://arxiv.org/abs/1906.05909>.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *CoRR*, abs/1706.03762, 2017. URL <http://arxiv.org/abs/1706.03762>.

Zhilin Yang, Zihang Dai, Yiming Yang, Jaime G. Carbonell, Ruslan Salakhutdinov, and Quoc V. Le. Xlnet: Generalized autoregressive pretraining for language understanding. *CoRR*, abs/1906.08237, 2019. URL <http://arxiv.org/abs/1906.08237>.