

# On the Relationship between Self-Attention & Convolutional Layers

Jean-Baptiste Cordonnier, Andreas Loukas, Martin Jaggi

Presented by Yenson Lau — Layer 6 AI

August 11th, 2021

# Motivation

Lots of intuition available for convolutions

- Great for modeling functions on signals with *translational symmetries*

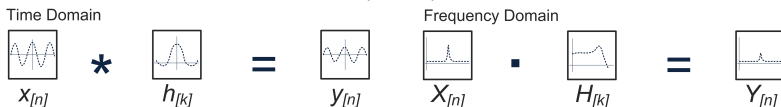
- Equivariance / invariance:

$$\text{Conv}(W, \mathcal{T}_\delta X) = \mathcal{T}_\delta \text{Conv}(W, X)$$

- Associativity / commutativity:

$$\text{Conv}(W', \text{Conv}(W, X)) = \text{Conv}(\text{Conv}(W', W), X) = \text{Conv}(\text{Conv}(W, W'), X)$$

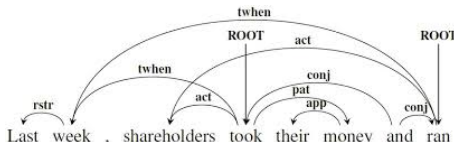
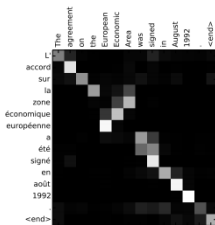
- Diagonalized by the Fourier Transform ([source](#))



- Great at detecting local patterns in some fixed proximity – as long as order is preserved. Naturally a core component of deep learning models for vision tasks.
  - e.g. [LeNet](#), [AlexNet](#), [ResNet](#), etc.
- For language tasks, however, a “sliding-window” approach alone is insufficient for identifying patterns.

# Motivation

- Text does have a loose sequential component but long-range context matters, ordering of words can change in different settings, etc. ([source](#))



- Unlike convolutions, attention can simultaneously attend to every word in an input sequence and produces the best language models, e.g. [Transformer](#), [BERT](#), etc.
- In fact, NN architectures using self-attention (SA) only (without convolution) can compete with SA + convolutional architectures on vision tasks.
  - [Stand-Alone Self-Attention in Vision Models, Ramachandran et. al. 2018.](#)
- Can attention be interpreted as a “generalized convolution”, and does this interpretation provide meaningful insight into how the mechanism works?

- A constructive theoretical proof that SA can express convolutional layers using relative positional encoding:

**Main theorem.** *A multi-head self-attention (MHSA) layer with  $N_h$  heads of dimension  $D_h$  each, final output dimension  $D_{out}$ , and a relative positional encoding of dimension  $D_p \geq 3$  can express any 2D convolutional layer of kernel size  $\sqrt{N_h} \times \sqrt{N_h}$  and  $\min(D_h, D_{out})$  output channels.*

- Experiments showing that SA behaves (somewhat) similarly to the theoretical construction for vision tasks.
- Attention probabilities possibly interpretable as some form of “generalized convolution”.

# Self-Attention

- Let  $\mathbf{X} \in \mathbb{R}^{T \times D_{in}}$  be an input matrix consisting of  $T$  tokens in  $D_{in}$  dimensions. For an output token  $t$ , the self-attention operation  $\text{SA} : T \times D_{in} \rightarrow D_{out}$  is expressed as

$$\text{SA}(\mathbf{X})_{t,:} \doteq \text{softmax}(\mathbf{A}_{t,:}) \mathbf{X} \mathbf{W}_{val}, \quad (1)$$

$$T \begin{bmatrix} \text{SA}(\mathbf{X}) \end{bmatrix}^{D_{out}} = T \begin{bmatrix} \text{softmax}(\mathbf{A}_{1,:}) \\ \vdots \\ \text{softmax}(\mathbf{A}_{T,:}) \end{bmatrix}^T T \begin{bmatrix} \mathbf{X} \end{bmatrix}^{D_{in}} \begin{bmatrix} \mathbf{W}_{val} \end{bmatrix}^{D_{out}}$$

- We refer to the the elements of the  $T \times T$  matrix

$$\mathbf{A} \doteq \mathbf{X} \mathbf{W}_{qry} \mathbf{W}_{key}^T \mathbf{X}^T \quad (2)$$

the *attention scores*, and the softmax outputs as *attention probabilities*.

- So far the learnable parameters are the *query*, *key*, and *value* matrices

$$\mathbf{W}_{qry} \in \mathbb{R}^{D_{in} \times D_k}, \quad \mathbf{W}_{key} \in \mathbb{R}^{D_{in} \times D_k}, \quad \text{and} \quad \mathbf{W}_{val} \in \mathbb{R}^{D_{in} \times D_{out}}.$$

For simplicity, ignore residual connections, normalization or constant factors.

- Beneficial to replicate the SA mechanism into multiple heads: concatenate the output of  $N_h$  heads of size  $D_h$  and project it to dimension  $D_{out}$ .

$$\text{MHSA}(\mathbf{X}) \doteq \text{hcat}_{h \in [N_h]} [\text{SA}_h(\mathbf{X})] \mathbf{W}_{out} + \mathbf{b}_{out}. \quad (3)$$

Here  $\mathbf{W}_{out} \in \mathbb{R}^{N_h D_h \times D_{out}}$  is the projection matrix and  $\mathbf{b}_{out} \in \mathbb{R}^{D_{out}}$  is a bias term.

$$\overset{D_{out}}{\mathbf{T} \begin{bmatrix} \text{MHSA}(\mathbf{X}) \end{bmatrix}} = \overset{N_h \times D_h}{\mathbf{T} \begin{bmatrix} \text{SA}_1(\mathbf{X}) & \dots & \text{SA}_{N_h}(\mathbf{X}) \end{bmatrix}} \overset{\substack{N_h \\ \times \\ D_h}}{\mathbf{W}_{out}} + \mathbf{b}_{out}$$

- $\mathbf{A}$  is *permutation equivariant* – shuffling the order of the tokens (rows) in  $\mathbf{X}$  shuffles the scores in  $\mathbf{A}$ . *This is not desired behavior.*
  - e.g. "the cat ate the fish" is very different from "the fish ate the cat"
- Overcome this by adding a (fixed or learned) *positional encoding*  $\mathbf{P}_{t,:}$ , for each token in the sequence when computing attention scores

$$\mathbf{A} \doteq (\mathbf{X} + \mathbf{P}) \mathbf{W}_{qry} \mathbf{W}_{key}^T (\mathbf{X} + \mathbf{P})^T. \quad (4)$$

The encoding matrix  $\mathbf{P}$  has size  $T \times D_{in}$ .

- e.g. sinusoidal positional encoding.

- For a pair of query (the token to compute the representation of) and key tokens (the pixel we are attending to)  $q, k$ , we can write this out as

$$\begin{aligned} \mathbf{A}_{q,k}^{abs} &= (\mathbf{X}_{q,:} + \mathbf{P}_{q,:}) \mathbf{W}_{qry} \mathbf{W}_{key}^T (\mathbf{X}_{q,:} + \mathbf{P}_{q,:})^T \\ &= \mathbf{X}_{q,:} \mathbf{W}_{qry} \mathbf{W}_{key}^T \mathbf{X}_{k,:}^T + \mathbf{X}_{q,:} \mathbf{W}_{qry} \mathbf{W}_{key}^T \mathbf{P}_{k,:}^T \\ &\quad + \mathbf{P}_{q,:} \mathbf{W}_{qry} \mathbf{W}_{key}^T \mathbf{X}_{k,:}^T + \mathbf{P}_{q,:} \mathbf{W}_{qry} \mathbf{W}_{key}^T \mathbf{P}_{k,:}^T \end{aligned} \quad (5)$$

- Relative positional encoding* is less restrictive and instead considers the positional difference between the query and key tokens:

$$\mathbf{A}_{q,k}^{rel} = \mathbf{X}_{q,:}^T \mathbf{W}_{qry} \mathbf{W}_{key}^T \mathbf{X}_{k,:} + \mathbf{X}_{q,:}^T \mathbf{W}_{qry} \hat{\mathbf{W}}_{key} \mathbf{r}_\delta + \mathbf{u}^T \mathbf{W}_{key} \mathbf{X}_{k,:}^T + \mathbf{v}^T \hat{\mathbf{W}}_{key} \mathbf{r}_\delta. \quad (6)$$

- $\mathbf{u}, \mathbf{v}$ : fixed or learned, unique for each head
  - $\mathbf{r}_\delta \in \mathbb{R}^{D_p}$ : depends only on the shift  $\delta \doteq k - q$ , shared by all layers and heads.
  - Key weights are split into  $\mathbf{W}_{key}$  for the input and  $\hat{\mathbf{W}}_{key}$  for positional encoding.
- Attention scores are now *translation equivariant* rather than permutation equivariant. This is the desired behavior for convolutional / vision tasks.



# Self-Attention for Images

From 1D sequence to 2D

- Replace query and key tokens with pixels  $\mathbf{q}, \mathbf{k} \in [W] \times [H]$ , and the input with  $\mathbf{X} \in \mathbb{R}^{W \times H \times D}$ . Each attention score now associates a query and key pixel.
- For a pixel  $\mathbf{p} \in (i, j)$ ,  $\mathbf{X}_{\mathbf{p},:} \doteq \mathbf{X}_{i,j,:}$  and  $\mathbf{A}_{\mathbf{p},:} \doteq \mathbf{A}_{i,j,:}$ .
- Then, analogously to the 1D case,

$$\text{SA}(\mathbf{A})_{\mathbf{q},:} = \sum_k \text{softmax}(\mathbf{A}_{\mathbf{q},:})_k \mathbf{X}_{\mathbf{k},:} \mathbf{W}_{val}, \quad (7)$$

and MHSA retains the same form as Equation (3). Similarly, relative positional encoding retains the same form as Equation (6).

# Self-Attention for Images

## Convolution

- Given an image tensor  $\mathbf{X} \in \mathbb{R}^{W \times H \times D_{in}}$  and kernel tensor  $\mathbf{W} \in \mathbb{R}^{K \times K \times D_{in} \times D_{out}}$ ,

$$\text{Conv}(\mathbf{X})_{i,j,:} \doteq \sum_{(\delta_1, \delta_2) \in \Delta_K} \mathbf{X}_{i-\delta_1, j-\delta_2, :} \mathbf{W}_{\delta_1, \delta_2, :, :} + \mathbf{b} \in \mathbb{R}^{D_{out}}, \quad (8)$$

where

$$\Delta_K \doteq \left[ -\left\lfloor \frac{K}{2} \right\rfloor, \dots, \left\lfloor \frac{K}{2} \right\rfloor \right] \times \left[ -\left\lfloor \frac{K}{2} \right\rfloor, \dots, \left\lfloor \frac{K}{2} \right\rfloor \right]$$

is the set of all shifts represented by a  $K \times K$  kernel.

- In the paper, the convolution is flipped so the summand is  $\mathbf{X}_{i+\delta_1, j+\delta_2, :} \mathbf{W}_{\delta_1, \delta_2, :, :}$ .  
Does not affect the outcome of the theorem.

# Self-Attention as a Convolutional Layer

## Main Theorem

- **Main theorem.** A multi-head self-attention (MHSA) layer with  $N_h$  heads of dimension  $D_h$  each, final output dimension  $D_{out}$ , and a relative positional encoding of dimension  $D_p \geq 3$  can express any 2D convolutional layer of kernel size  $\sqrt{N_h} \times \sqrt{N_h}$  and  $\min(D_h, D_{out})$  output channels.

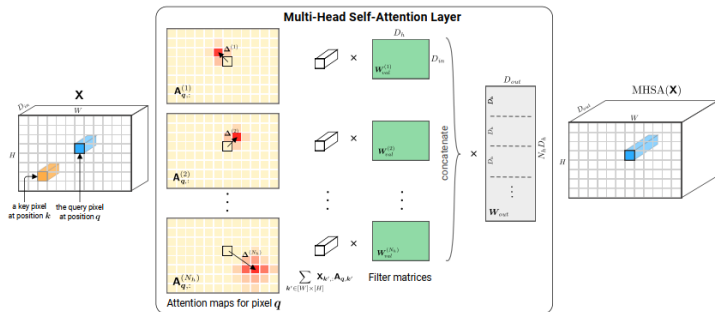


Figure 1: Illustration of a Multi-Head Self-Attention layer applied to a tensor image  $\mathbf{X}$ . Each head  $h$  attends pixel values around shift  $\Delta^{(h)}$  and learn a filter matrix  $\mathbf{W}_{val}^{(h)}$ . We show attention maps computed for a query pixel at position  $q$ .

- The theorem is a consequence of two lemmas.

# Self-Attention as a Convolutional Layer

## Lemma 1

- **Lemma 1.** Consider a MHSA layer consisting of  $N_h = K^2$  heads,  $D_h \geq D_{out}$  and let  $f : [N_h] \rightarrow \Delta_K$  be a bijective mapping of heads onto shifts. Further, suppose that for every head

$$\text{softmax}(\mathbf{A}_{q,:}^{(h)})_k = \begin{cases} 1 & \text{if } f(h) = q - k, \\ 0 & \text{otherwise.} \end{cases} \quad (9)$$

Then for any convolutional layer with a  $K \times K$  kernel and  $D_{out}$  output channels, there exists  $\{\mathbf{W}_{val}^{(h)}\}_{h \in [N_h]}$  such that  $\text{MHSA}(\mathbf{X}) = \text{Conv}(\mathbf{X})$ , for any kernel tensor  $\mathbf{W} \in \mathbb{R}^{K \times K \times D_{in} \times D_{out}}$ , and for all  $\mathbf{X} \in \mathbb{R}^{W \times H \times D_{in}}$ .

- *Proof of Lemma 1.* Rewrite Equation (3) as

$$\text{MHSA}(\mathbf{X}) = \mathbf{b}_{out} + \sum_{h \in [N_h]} \text{softmax}(\mathbf{A}^{(h)}) \mathbf{X} \underbrace{\mathbf{W}_{val}^{(h)} \mathbf{W}_{out}^{(h)}}_{\mathbf{W}^{(h)}}. \quad (10)$$

Here  $\mathbf{W}_{out}^{(h)} \in \mathbb{R}^{D_h \times D_{out}}$  is the output matrix for head  $h$ , or in MATLAB notation,

$$\mathbf{W}_{out}^{(h)} = \mathbf{W}[(h-1) * D_h + 1 : hD_h, :] \in \mathbb{R}^{D_h \times D_{out}}.$$

# Self-Attention as a Convolutional Layer

## Lemma 1

- *Proof of Lemma 1 (cont).* Consider a single "query" pixel from MHSA( $\mathbf{X}$ ):

$$\text{MHSA}(\mathbf{X})_{q,:} = \mathbf{b}_{out} + \sum_{h \in [N_h]} \left( \sum_{k \in [T]} \text{softmax}(\mathbf{A}_{q,:}^{(h)})_k \mathbf{X}_{k,:} \right) \mathbf{W}^{(h)}. \quad (11)$$

The summand  $\sum_k \text{softmax}(\mathbf{A}_{q,:}^{(h)})_k \mathbf{X}_{k,:}$  is a weighted average on the rows of  $\mathbf{X}$ .

- Applying the assumption from Equation (9), each of the softmax weights pick out a single pixel  $k = q - f(h)$ :

$$\text{MHSA}(\mathbf{X})_{q,:} = \mathbf{b}_{out} + \sum_{h \in [N_h]} \mathbf{X}_{q-f(h),:} \mathbf{W}^{(h)}. \quad (12)$$

Set  $K = \sqrt{N_h}$  and index all each shift using one  $h \in [N_h]$ , i.e.  $f([N_h]) = \Delta_K$ . Then setting  $\mathbf{W}_{f(h),:,} = \mathbf{W}^{(h)}$  we get

$$\text{MHSA}(\mathbf{X})_{q,:} = \mathbf{b}_{out} + \sum_{h \in [N_h]} \mathbf{X}_{q-f(h),:} \mathbf{W}_{f(h),:,} = \text{Conv}(\mathbf{X})_{q,:}. \quad \square \quad (13)$$

# Self-Attention as a Convolutional Layer

## Lemma 1: Summary

- Lemma 1 says that, with the right choice of relative positional encoding (and a very restricted set of attention parameters), *a MHSA layer is exactly equivalent to a convolutional layer*.
- The number of linearly independent filters expressible by  $\mathbf{W}$  is clearly limited by  $\min(D_h, D_{out})$ . Therefore, to express  $D_h$  convolutional filters, it is best to simply let  $D_{out} = N_h D_h$  and have  $\mathbf{W}^{(h)} \in \mathbb{R}^{D_h \times D_h}$ .
- There is a generalized version of this lemma in the Appendix of the paper that delves into the space of filters spannable for a given set of dimensions  $D_h, N_h, D_{out}$ , etc., for interested readers.

# Self-Attention as a Convolutional Layer

## Lemma 2

- **Lemma 2.** *There exists a relative encoding scheme  $\{\mathbf{r}_\delta \in \mathbb{R}^{D_p}\}_{\delta \in \mathbb{Z}^2}$  with  $D_p \geq 3$  and parameters  $\mathbf{W}_{qry}$ ,  $\mathbf{W}_{key}$ ,  $\hat{\mathbf{X}}_{key}$ ,  $\mathbf{u}$  with  $D_p \leq D_k$  such that, for every  $\tau \in \Delta_K$  there exists some vector  $\mathbf{v}(\tau)$  that yields  $\text{softmax}(\mathbf{A}_{q,:})_k = 1$  if  $k - q = \tau$ , and zero otherwise.*
- *Proof of Lemma 2.* Start with the relative positional encoding of Equation (6):

$$\mathbf{A}_{q,k}^{rel} = \mathbf{X}_{q,:}^T \mathbf{W}_{qry} \mathbf{W}_{key}^T \mathbf{X}_{k,:} + \mathbf{X}_{q,:}^T \mathbf{W}_{qry} \hat{\mathbf{W}}_{key}^T \mathbf{r}_\delta + \mathbf{u}^T \mathbf{W}_{key} \mathbf{X}_{k,:} + \mathbf{v}^T \hat{\mathbf{W}}_{key} \mathbf{r}_\delta.$$

Since the desired attention probabilities (9) from Lemma 1 are independent of the input  $\mathbf{X}$ , set  $\mathbf{W}_{key} = \mathbf{W}_{qry} = 0$ .

This leaves only the final term. Setting  $\hat{\mathbf{W}}_{key} = [\mathbf{I}_{D_p} ; 0] \in \mathbb{R}^{D_k \times D_p}$  (recall  $D_p \leq D_k$ ) yields

$$\mathbf{A}_{q,k} = \mathbf{v}_{1:D_p}^T \mathbf{r}_\delta.$$

# Self-Attention as a Convolutional Layer

## Lemma 2

- *Proof of Lemma 2 (cont).* Recall that  $\delta \doteq \mathbf{k} - \mathbf{q}$ . Now suppose some  $\mathbf{v}, \mathbf{r}_\delta$  could be found such that

$$\mathbf{A}_{\mathbf{q}, \mathbf{k}} = -\alpha(\|\delta - \tau\|^2 + c) \quad (14)$$

so that the maximum score  $-\alpha c$  is achieved when  $\delta = \tau$ , then

$$\begin{aligned} \lim_{\alpha \rightarrow \infty} \text{softmax}(\mathbf{A}_{\mathbf{q}, :})_{\mathbf{k}} &= \lim_{\alpha \rightarrow \infty} \frac{e^{-\alpha(\|\delta - \tau\|^2 + c)}}{\sum_{\delta' \in \Delta_K} e^{-\alpha(\|\delta' - \tau\|^2 + c)}} \\ &= \frac{1_{\delta = \tau}}{1 + \lim_{\alpha \rightarrow \infty} \sum_{\delta' \neq \delta} e^{-\alpha\|\delta' - \tau\|^2}} = 1_{\delta = \tau}. \end{aligned}$$

- Choosing  $\mathbf{v}_{1:D_p} = -\alpha(1, -2\tau_1, -2\tau_2)$  and  $\mathbf{r}_\delta = (\|\delta\|^2, \delta_1, \delta_2)$ , yields

$$\begin{aligned} \mathbf{A}_{\mathbf{q}, \mathbf{k}} &= \mathbf{v}_{1:D_p}^T \mathbf{r}_\delta \\ &= -\alpha(\|\delta\|^2 - 2\tau_1\delta_1 - 2\tau_2\delta_2) \\ &= -\alpha(\|\delta - \tau\|^2 - \|\tau\|^2). \end{aligned}$$

Setting  $c = -\|\tau\|^2$  recovers Equation (14) and completes the proof. □



# Self-Attention as a Convolutional Layer

## Lemma 2: Summary

- The encoding scheme satisfying Lemma 2 is a *quadratic encoding scheme*, and is achieved by setting

$$\begin{aligned}\mathbf{v}_{1:D_p}^{(h)} &\doteq -\alpha^{(h)}(1, -2\tau_1^{(h)}, -2\tau_2^{(h)}), \\ \mathbf{r}_\delta &\doteq (\|\delta\|^2, \delta_1, \delta_2), \\ \mathbf{W}_{qry} &= \mathbf{W}_{key} \doteq 0, \\ \hat{\mathbf{W}}_{key} &\doteq [\mathbf{I}_{D_p} ; 0].\end{aligned}\tag{15}$$

The learned parameters  $\boldsymbol{\tau}^{(h)} = (\tau_1^{(h)}, \tau_2^{(h)})$  and  $\alpha^{(h)}$  denote the center and width of attention of each head.  $\delta = (\delta_1, \delta_2)$  and therefore  $\mathbf{r}_\delta$  is fixed and expresses the relative shift between query and key pixels.

- Although the proof requires  $\alpha \rightarrow \infty$  to satisfy the assumption (9) of Lemma 1, finite precision arithmetic performs hard attention with sufficiently large enough  $\alpha$ . E.g. for Float32, set  $\alpha \geq 46$ .
- The lemma, and thus the theorem, can be extended in a straightforward manner to cover  $K$ -dimensional convolutions, with  $D_p = K + 1$ .

- **Baseline.** Standard ResNet18 on CIFAR-10
- **Challenger models.**
  - 6 MHSA layers. Each layer consists of a MHSA step, followed by dense, dropout and LayerNorm sublayers.
  - Use 2x2 invertible downsampling to reduce image size (attention coefficients scale quadratically to image size)
  - The last layer representation is average pooled and fed to a linear classifier
- **MHSA variations.** Different types of relative positional encoding

$$\mathbf{A}_{q,k}^{rel} = \mathbf{X}_{q,:}^T \mathbf{W}_{qry} \mathbf{W}_{key}^T \mathbf{X}_{k,:} + \mathbf{X}_{q,:}^T \mathbf{W}_{qry} \hat{\mathbf{W}}_{key} \mathbf{r}_{\delta} + \mathbf{u}^T \mathbf{W}_{key} \mathbf{X}_{k,:}^T + \mathbf{v}^T \hat{\mathbf{W}}_{key} \mathbf{r}_{\delta}.$$

- **SA with quadratic embedding:** Retain final term only and fix the variables using Equation (15). The attention widths  $\alpha^{(h)}$  and centers  $\tau^{(h)}$  are still learnt.
- **SA with learned embedding:** Retain final term only but learn  $\mathbf{v}$ ,  $\hat{\mathbf{W}}_{key}$ ,  $\mathbf{r}_{\delta}$ , with  $D_p = D_{out} = 400$ . Set  $D_h = D_{out}$ .
- **SA with content-based attention:** All terms retained (might actually be just first two terms) and all variables learnable. Same dimensions as above.

- **ResNet converges faster:**  
Probably because SA's inductive bias is not as strong as ResNet, but may also be due to different optimization setup.
- SA models with more learnable parameters converge slower and ends up with lower testing accuracy.

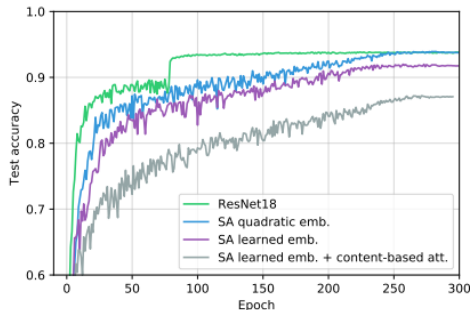


Figure 2: Test accuracy on CIFAR-10.

# Do Self-Attention Layers Actually Learn Convolutions? Maybe.

- Recall that from Lemma 1, the critical assumption is

$$\text{softmax}(\mathbf{A}_{q,:}^{(h)})_k = 1_{f(h)=q-k}, \quad \forall h.$$

If  $N_h = |\Delta_K|$ , mapping each head to a single shift  $f(h) = \tau \in \Delta_K$  yields

$$\text{softmax}(\mathbf{A}_{q,:}^{\tau})_k = 1_{\tau=q-k}, \quad \forall \tau.$$

- MHSA expresses a single convolution operation when attention probabilities are all elements of a *translation group*. For instance, this would also suffice:

$$\text{softmax}(\mathbf{A}_{q,:}^{\tau})_k \propto 1_{|(q-k)-\tau| \leq R}, \quad \forall \tau.$$

- Expresses “sparse convolution” if we remove some pixels from being attended to in  $\Delta_K$ ; the heads are still elements of a translation group.
- Expresses “mixture of convolutions” if there were multiple distribution profiles in the attention probabilities.
- Takeaway:** *It's difficult to judge whether a MHSA layer has a preference for learning convolutional operators or not.*

# Do Self-Attention Layers Actually Learn Convolutions? Maybe.

## Quadratic encoding

- Here the attention probabilities are Gaussians shifted from origin  $\mathbf{q}$ , so a single convolution is expressed if they have the same width and different centers.
- The attention probabilities for each of the  $N_h = 9$  heads after training:

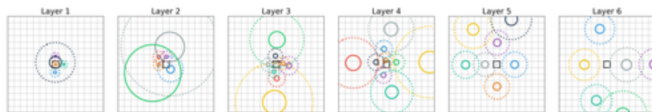


Figure 4: Centers of attention of each attention head (different colors) for the 6 self-attention layers using quadratic positional encoding. The central black square is the query pixel, whereas solid and dotted circles represent the 50% and 90% percentiles of each Gaussian, respectively.

- A significant portion of the heads have similar widths, esp. at lower layers. Higher layers have more varied widths. Notice that the centers tend to separate.
- Again, should we interpret this as a mixture of convolutions? Sparse convolutions? Convolutions with strides? Hard to say.

# Do Self-Attention Layers Actually Learn Convolutions? Maybe.

Quadratic / non-isotropic Gaussian encoding

- $N_h = 16$  heads. Similar to behavior with 9 heads, but now lower layers have groups of heads with small widths + a few heads with widely varying widths.

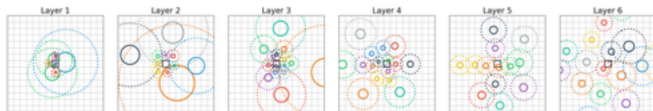


Figure 14: Centers of attention for 16 attention heads (different colors) for the 6 self-attention layers using quadratic positional encoding. The central black square is the query pixel, whereas solid and dotted circles represent the 50% and 90% percentiles of each Gaussian, respectively.

- Non-isotropic Gaussians ( $N_h = 9$ ). Large variations in shapes and sizes, but perhaps useful for representing rotations?

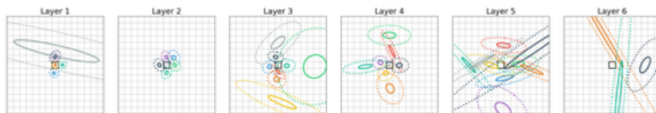


Figure 12: Centers of attention of each attention head (different colors) for the 6 self-attention layers using non-isotropic Gaussian parametrization. The central black square is the query pixel, whereas solid and dotted circles represent the 50% and 90% percentiles of each Gaussian, respectively.

# Do Self-Attention Layers Actually Learn Convolutions? Maybe.

Learned relative positional encoding (w/o content)

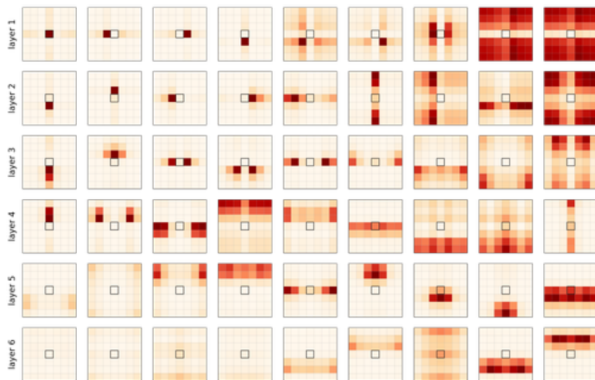


Figure 5: Attention probabilities of each head (*column*) at each layer (*row*) using learned relative positional encoding without content-based attention. The central black square is the query pixel. We reordered the heads for visualization and zoomed on the 7x7 pixels around the query pixel.

- Each layer seems to mostly learn one or two attention profiles.

# Do Self-Attention Layers Actually Learn Convolutions? Maybe.

Learned relative positional encoding with content

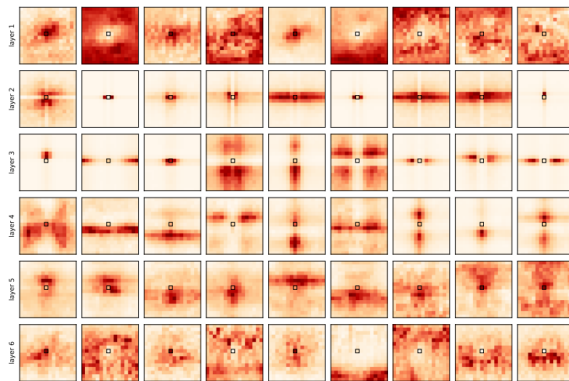


Figure 6: Attention probabilities for a model with 6 layers (rows) and 9 heads (columns) using learned relative positional encoding and content-content based attention. Attention maps are averaged over 100 test images to display head behavior and remove the dependence on the input content. The black square is the query pixel. More examples are presented in Appendix [A](#)

- Attention probabilities are averaged over images. Still demonstrates some of the behavior shown for the learned embeddings without content.
- Attention profiles do seem to shift with the query pixel  $q$  (see Appendix).



Thanks for listening!

