

On the Relationship between Self-Attention and Convolutional Layers

Jean-Baptiste Cordonnier, Andreas Loukas, Martin Jaggi
Ecole Polytechnique Federale de Lausanne (EPFL)

ICLR 2020

Presented by Eungyeup Kim

Vision Seminar
03 MAR 2020

Motivation

Attention mechanisms in vision is working well recently.

- *Bello et al.(2019)* replaced some convolutional layers with self-attention(SA) layers, leading to improvements on image classification and object detection tasks.
- *Ramachandran et al.(2019)* noticed that SOTA results are reached when attention and convolutional features are combined, or even when **only** SA only architecture is used.
- So, **do SA layers process images in a similar manner to convolutional layers?**

This work

- 1) Provides the constructive proof showing that multi-head SA layers with relative positional encoding can express any convolutional layers.
- 2) Provides experiments demonstrating the layers of attention-only architectures act like a convolutional layers, such as attending on grid-like pattern around each query pixel.

Backgrounds

Multi-head self-attention layer (MHSA)

- *Attention scores*

$$A := XW_{\text{query}}W_{\text{key}}^{\top}X^{\top}$$

- *Attention scores (with positional encoding)*

$$A := (X + P)W_{\text{query}}W_{\text{key}}^{\top}(X + P)^{\top}, \text{ where } P \text{ indicates positional encoding}$$

- Self – Attention $(X)_{q,:} = \sum_k \text{softmax}(A_{q,:})_k X_{k,:}W_{\text{value}}$, where $q = (i, j)$ indicating 2D coordinates.
- $\text{MHSA}(X) := \text{concat}_{h \in [N_h]} [\text{Self – Attention}_h(X)] W_{\text{out}} + b_{\text{out}}$, where $W_{\text{out}} \in R^{N_h D_h \times D_{\text{out}}}$

Convolutional layer

- $\text{Conv}(X)_{i,j,:} := \sum_{(\delta_1, \delta_2) \in \Delta_K} X_{i+\delta_1, j+\delta_2,:} W_{\delta_1, \delta_2,:} + b$, where $\Delta_K := \left[-\left[\frac{K}{2}\right], \dots, \left[\frac{K}{2}\right]\right] \times \left[-\left[\frac{K}{2}\right], \dots, \left[\frac{K}{2}\right]\right]$ containing all possible shifts with a $K \times K$ kernel.

Backgrounds

Positional Encoding

- Absolute encoding : a vector $P_{p,:}$ is assigned to each pixel p .

$$\begin{aligned} A_{q,k}^{abs} &:= (X_{q,:} + P_{q,:})^\top W_{\text{query}}^\top W_{\text{key}} (X_{k,:} + P_{k,:}) \\ &= X_{q,:}^\top W_{\text{query}}^\top W_{\text{key}} X_{k,:} + X_{q,:}^\top W_{\text{query}}^\top W_{\text{key}} P_{k,:} + P_{q,:}^\top W_{\text{query}}^\top W_{\text{key}} X_{k,:} + P_{q,:}^\top W_{\text{query}}^\top W_{\text{key}} P_{k,:} \end{aligned}$$

- Relative encoding : only consider the position difference between pixels. (Dai et al.)

$$A_{q,k}^{rel} = \underbrace{X_{q,:}^\top W_{\text{query}}^\top W_{\text{key}} X_{k,:}}_{\text{content-based key vectors}} + \underbrace{X_{q,:}^\top W_{\text{query}}^\top \hat{W}_{\text{key}} r_\delta}_{\text{location-based key vectors}} + \underbrace{u^\top W_{\text{key}} X_{k,:}}_{\text{global content bias}} + \underbrace{v^\top \hat{W}_{\text{key}} r_\delta}_{\text{global positional bias}}$$

where the attention scores only depend on the shift $\delta := k - q$.

- Trainable u and v are unique for each head, while non-trainable r_δ is shared by all layers and heads.
- We deliberately separate the W_{key} and \hat{W}_{key} for producing content-based key vectors and location-based key vectors respectively.
- Content-based addressing, content-dependent positional bias, global content bias, global positional bias

Self-Attention vs Convolutional Layers

Theorem 1. A multi-head self-attention layer with N_h heads of dimension D_h , output dimension D_{out} and a relative positional encoding of dimension $D_p \geq 3$ can express any convolutional layer of kernel size $\sqrt{N_h} \times \sqrt{N_h}$ and $\min(D_h, D_{out})$ output channels.

Lemma 1. Consider a multi-head self-attention 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 the following holds:

$$\text{softmax}\left(A_{q,:}^{(h)}\right)_k = \begin{cases} 1 & \text{if } f(h) = q - k \\ 0 & \text{o/w} \end{cases}$$

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

Self-Attention vs Convolutional Layers

proof

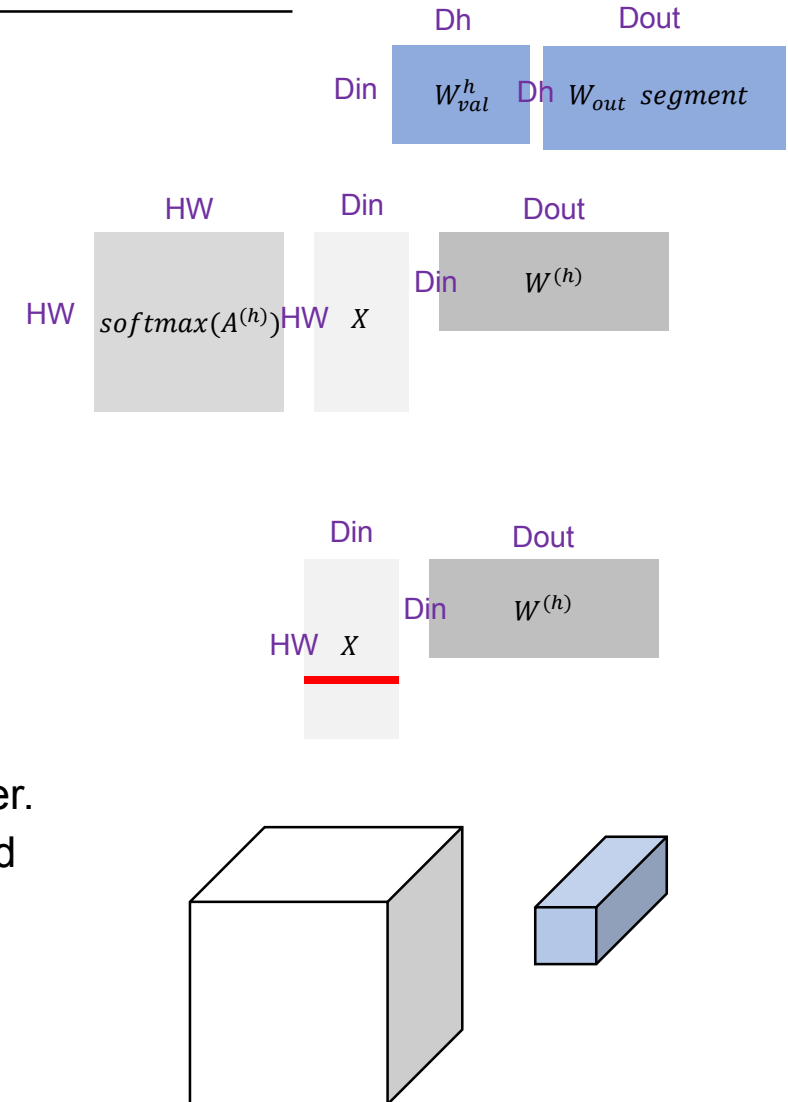
$$MHSA(X) = b_{out} + \sum_{h \in [N_h]} softmax(A^{(h)}) X \underline{W_{val}^h W_{out}[(h-1)D_h + 1 : hD_h + 1]}$$

$$MHSA(X)_{q,:} = \sum_{h \in [N_h]} (\sum_K softmax(A_{q,:}^{(h)}) X_{k,:}) \underline{W^{(h)}} + b_{out}$$

Due to the conditions of Lemma 1, for the h^{th} attention head the attention probability is one when $k = q - f(h)$ and zero otherwise. The layer's output at pixel q is thus equal to

$$MHSA(X)_q = \sum_{h \in [N_h]} X_{q-f(h),:} W^{(h)} + b_{out}$$

- For $K = \sqrt{N_h}$, the above can be seen to be equivalent to a convolutional layer.
- There is a one to one mapping (implied by f) between the matrices $W^{(h)}$ and the matrices $W_{k1,k2,:}$ for all $(k1, k2) \in [K]^2$



Self-Attention vs Convolutional Layers

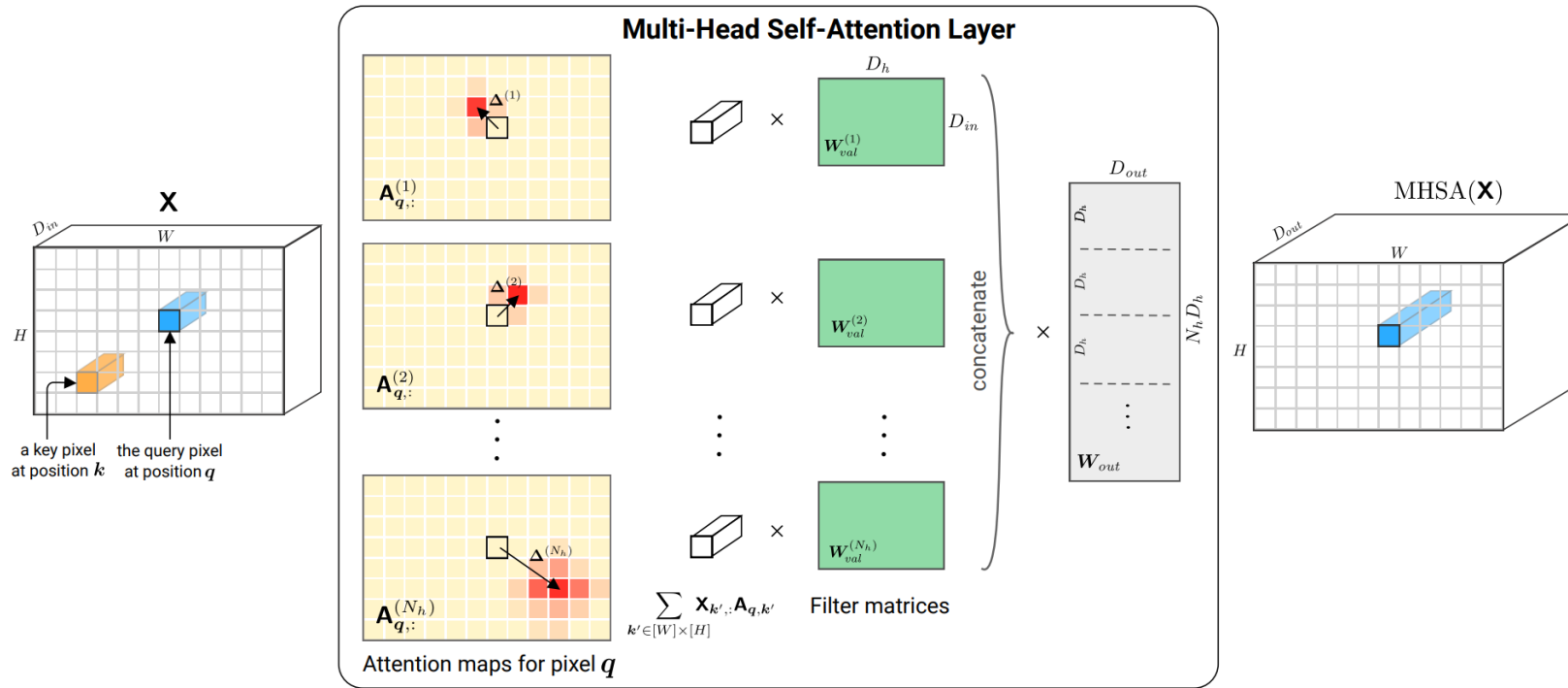


Figure 1: Illustration of a Multi-Head Self-Attention layer applied to a tensor image \mathbf{X} . Each head h attends pixel values around shift $\Delta^{(h)}$ and learn a filter matrix $\mathbf{W}_{val}^{(h)}$. We show attention maps computed for a query pixel at position q .

Self-Attention vs Convolutional Layers

On which condition can $\text{softmax}(A_{q,:})_k = \begin{cases} 1 & \text{if } f(h) = q - k \\ 0 & \text{o/w} \end{cases}$ be derived?

Lemma 2. There exists a relative encoding scheme $\{r_\delta \in R^{D_p}\}_{\delta \in Z^2}$ with $D_p \geq 3$ and parameter $W_{\text{query}}, W_{\text{key}}, u$ with $D_p \leq D_k$ such that, for every $\Delta \in \Delta_K$ there exists some vector v yielding $\text{softmax}(A_{q,:})_k = 1$ if $k - q = \Delta$ and zero, otherwise.

proof

$$A_{q,k} = X_{q,:}^\top W_{\text{query}}^\top W_{\text{key}} X_{k,:} + X_{q,:}^\top W_{\text{query}}^\top \hat{W}_{\text{key}} r_\delta + u^\top W_{\text{key}} X_{k,:} + v^\top \hat{W}_{\text{key}} r_\delta$$

As the attention probabilities are independent of the input tensor X , we set $W_{\text{query}} = W_{\text{key}} = 0$.
Thus

$$A_{q,k} = v^\top \hat{W}_{\text{key}} r_\delta = v^\top r_\delta \text{ } (\hat{W}_{\text{key}} \text{ is set to be identity})$$

where $r_\delta := k - q$

Suppose we could write

$$A_{q,k} = -\alpha(\|\delta - \Delta\|^2 + c)$$

for some constant c .

Self-Attention vs Convolutional Layers

In this way, we have

$$\lim_{\alpha \rightarrow \infty} \text{softmax}(A_{q,:})_k = \lim_{\alpha \rightarrow \infty} \frac{e^{-\alpha(\|\delta - \Delta\|^2 + c)}}{\sum_{k'} e^{-\alpha(\|(k-q) - \Delta\|^2 + c)}} = \lim_{\alpha \rightarrow \infty} \frac{e^{-\alpha(\|\delta - \Delta\|^2)}}{\sum_{k'} e^{-\alpha(\|(k-q) - \Delta\|^2)}}$$

For $\delta = \Delta$,

$$\lim_{\alpha \rightarrow \infty} \text{softmax}(A_{q,:})_k = \lim_{\alpha \rightarrow \infty} \frac{e^{-\alpha(\|\delta - \Delta\|^2)}}{\sum_{k'} e^{-\alpha(\|(k-q) - \Delta\|^2)}} = \frac{1}{1 + \lim_{\alpha \rightarrow \infty} \sum_{k' \neq k} e^{-\alpha(\|(k-q) - \Delta\|^2)}} = 1$$

For $\delta \neq \Delta$,

$$\lim_{\alpha \rightarrow \infty} \text{softmax}(A_{q,:})_k = \lim_{\alpha \rightarrow \infty} \frac{e^{-\alpha(\|\delta - \Delta\|^2)}}{\sum_{k'} e^{-\alpha(\|(k-q) - \Delta\|^2)}} = 0$$

So, which v and $\{r_\delta\}_{\delta \in \mathbb{Z}^2}$ satisfy $A_{q,k} = -\alpha(\|\delta - \Delta\|^2 + c)$?

If we set

$$v = -\alpha(1, -\Delta_1, -\Delta_2), \quad r_\delta = (\|\delta\|, \delta_1, \delta_2)$$

Then, $A_{q,k} = v^\top r_\delta = -\alpha(\|\delta - \Delta\|^2 + c)$

Self-Attention vs Convolutional Layers

Quadratic Encoding

$$v = -\alpha(1, -\Delta_1, -\Delta_2), \quad r_\delta = (\|\delta\|, \delta_1, \delta_2)$$

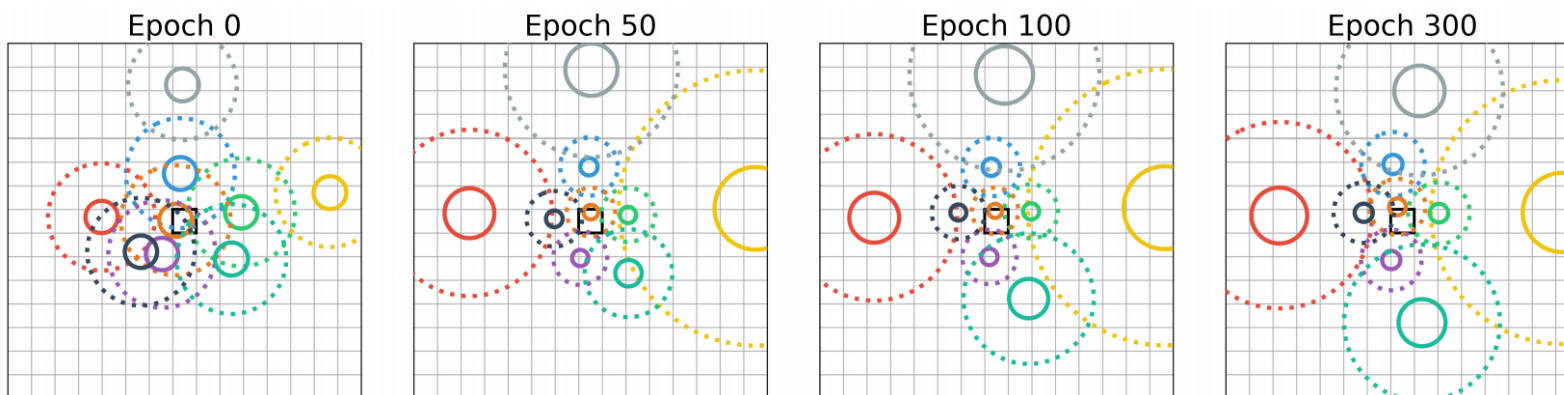
It is important that above encoding is not the only one for which the conditions of Lemma 1 are satisfied. The relative encoding learned by the neural network also matched the conditions of the lemma.

Nevertheless, the encoding defined above is very efficient in terms of size, as only $D_p = 3$ dimensions suffice to encode the relative position of pixels, while also reaching similar or better empirical performance than the learned one.

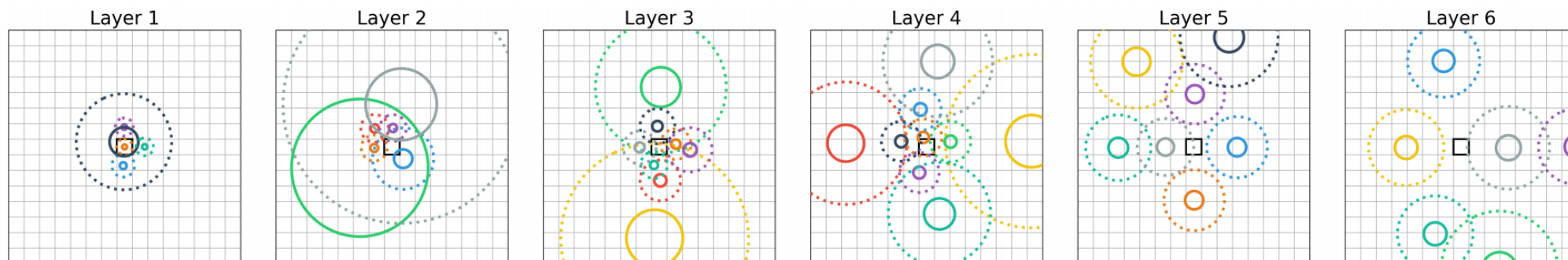
Experiments

A fully attentional model consisting of six multi-head self-attention layers is trained on Cifar-100 and ImageNet for image classification.

Using quadratic encoding : $\Delta^h \sim N(\mu, \sigma)$



[This figure shows how the initial positions of the heads at **layer 4** changed during training.]



[This figure shows that the first few layers tend to focus on local patterns, while deeper layers also attend to larger patterns by positioning the center of attention further from the query pixel.]

Experiments

Using learned encoding(without content-based attention) :

Satisfying lemma 1,
thus theorem as well.
(local pattern)



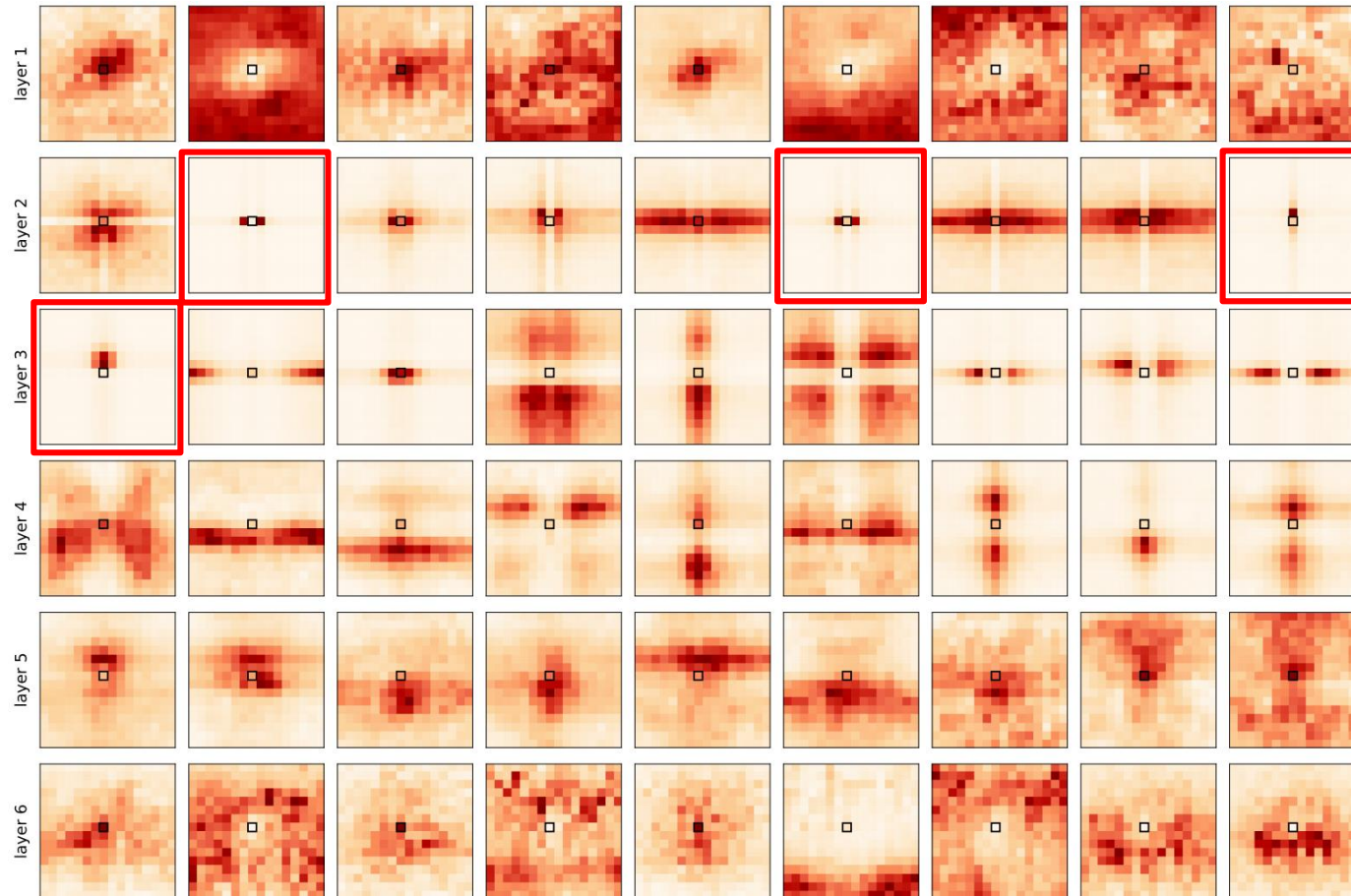
Capturing the horizontally-symmetric but not localized patterns. (long-range pixel inter-dependency)

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.

Experiments

Using learned encoding(with content-based attention) :

Satisfying lemma 1,
thus theorem as well.
(local pattern)



- Some heads use more **content-based attention** heads, leveraging the advantage of self-attention over CNN. This effectiveness was shown by *Bello et al.(2019)* that combining CNN and self-attention features outperforms each taken separately.

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.

Experiments

Using learned encoding(with content-based attention) :

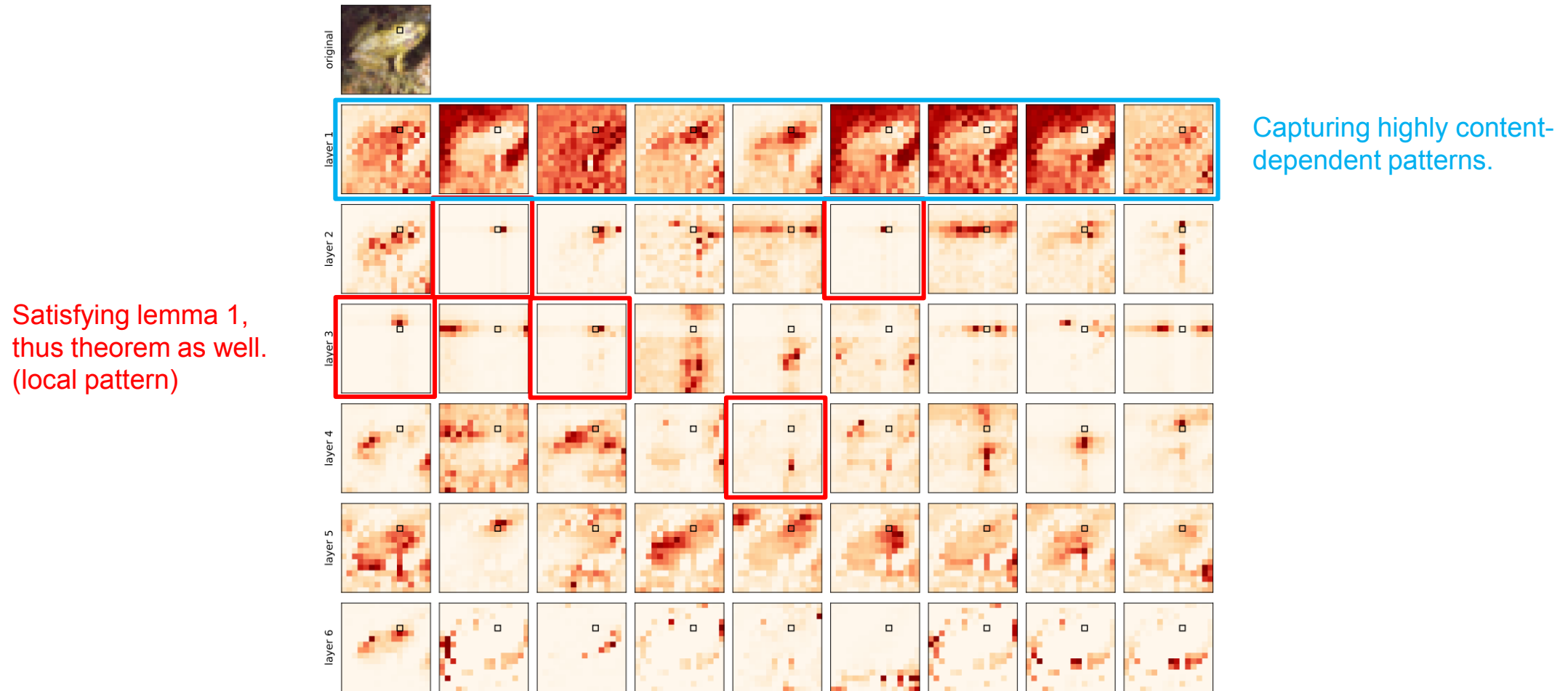


Figure 8: Attention probabilities for a model with 6 layers (*rows*) and 9 heads (*columns*) using learned relative positional encoding and content-content based attention. The query pixel (black square) is on the frog head.

Thank you