
[Re]: On the Relationship between Self-Attention and Convolutional Layers

Anonymous Author(s)

Affiliation

Address

email

1 Scope of Reproducibility

2 In this report, we perform a detailed study on the paper "On the Relationship between Self-Attention and Convolutional
3 Layers" [1] which provides theoretical and experimental evidence that self attention layers can behave like convolutional
4 layers. The proposed method does not obtain state-of-the-art performance but rather answers an interesting question - *do*
5 *self-attention layers process images in a similar manner to convolutional layers?*. This has inspired many recent works
6 like [2, 3] which propose a fully-attentional model for image recognition. We focus on experimentally validating the
7 claims of the original paper, highlight differences with other similar works, and propose a new variant of the attention
8 operation - *Hierarchical Attention* which shows improved performance with significantly lesser parameters. To facilitate
9 further study, all the code used in our experiments are publicly available here¹.

10 Methodology

11 We implement the original paper [1] from scratch in Pytorch and refer to the author's source code² for verification. In
12 our experiments involving SAN [2], we utilize the official implementation³ due to the available faster CUDA kernels
13 while we implement VIT [3] from scratch referring to the author's source code⁴. We then incorporate our proposed
14 hierarchical operation in all three methods for comparison. For all our experiments mentioned in this report, we use
15 the CIFAR10 dataset to benchmark the performance of the model in the image classification task. Each training run
16 required around 20 hours while the corresponding hierarchical versions took around 10-12 hours for convergence in an
17 Nvidia RTX 2060 GPU.

18 Results

19 We were able to reproduce all the results from the paper within 1% of the reported value, hence validating the claims
20 of the original paper. However, there seems to be some differences in the attention figures which lead to interesting
21 insights and the proposed Hierarchical Attention. In the case of VIT and SAN, we do not have a comparative baseline
22 as the corresponding papers do not evaluate performance on the CIFAR10 dataset (without pre-training).

23 What was easy

24 We did not face any major challenges in reproducing the results in the paper.

25 What was difficult

26 Most of the code in the official implementation seems to be borrowed from HuggingFace's repository⁵ which also
27 brought along a lot of unnecessary code making it difficult to read and understand quickly. Further the training time for
28 each run is quite significant making it difficult for us to experiment with multiple datasets and hyperparameter settings.

29 Communication with original authors

30 We have tried contacting the authors regarding the differences in the attention figures since the code for the same was
31 not available on the repository for verification. However, we have not received any response regarding the same.

¹Code available in the supplementary material during review period.

²<https://github.com/epfml/attention-cnn>

³<https://github.com/hszhao/SAN>

⁴https://github.com/google-research/vision_transformer

⁵<https://github.com/huggingface/transformers>

32 **1 Introduction**

33 In computer vision, convolutional architectures [4, 5, 6, 7] have dominated across various image recognition tasks
34 like classification, segmentation, etc. However, they have some limitations like lack of rotation invariance, inability to
35 aggregate information based on the image content, etc. This has inspired researchers to explore a different design space
36 and introduce models with interesting new capabilities. Self-attention based networks, in particular Transformers [8]
37 have become the model of choice for various Natural Language Processing (NLP) tasks. The major difference between
38 transformers and previous methods, such as recurrent neural networks and Convolutional Neural Networks (CNN), is
39 that the former can simultaneously learn to attend to various parts of the input sequence. To utilize this capacity to learn
40 meaningful inter-dependencies, many recent works have tried to incorporate self-attention [9, 10], some even replacing
41 convolutions entirely [11, 12], in networks for vision tasks. The main highlights of the paper "On the Relationship
42 between Self-Attention and Convolutional Layers" [1] are:

- 43 1. Theoretical proof that self-attention layers can behave similar to convolution layers.
44 2. Empirical validation of the performance of the self-attention model on the image classification task using the
45 CIFAR10 dataset. Visual evidence that self-attention applied to images learns convolutional filters around the
46 query pixels.
- 47 In this study, we perform a detailed analysis on the various experiments outlined in the paper. We observe certain
48 differences from the original paper which lead to interesting observations. We go beyond verifying the claims by trying
49 to solve the observed problems and propose a novel attention operation, which we refer to as *Hierarchical Attention*
50 (HA). Detailed experiments show that the proposed method significantly outperforms existing self-attention variants
51 with roughly $1/5^{th}$ the number of parameters and can be easily adapted to any existing attention-based architecture.

52 **1.1 Outline of this report**

53 We structure the report as follows:

- 54 1. Section 2 formally introduces the attention operation, followed by the fundamental working principle of the
55 transformer.
56 2. We validate the claims made by the paper in Section 3. We visualize the attention patterns of the suggested
57 models and compare with those mentioned in the paper, commenting on any similarities or differences.
58 3. We introduce a novel Hierarchical Attention operation described in Section 4. We compare the modified
59 operation on various methods and empirically show significant performance gains.
60 4. Section 5 concludes and suggests possible improvements for future research.

61 **2 Fundamentals**

62 In this section, we introduce the attention operation, first in terms of its origin in NLP and how it can be extended for
63 images. We also provide a short introduction to transformers since most methods described in this report heavily rely
64 on it.

65 **2.1 Attention**

66 Attention was first proposed for NLP, where the goal is to focus on a subset of important words. Consequently, relations
67 between inputs are highlighted that can be used to capture context and higher-order dependencies. The attention matrix
68 $A(\cdot)$ indicates a score between N queries Q and N_k keys, which indicates which part of the input sequence to focus on.
69 $\sigma(\cdot)$ is an activation function (generally $softmax(\cdot)$).

$$A(Q, K) = \sigma(QK^T) \quad (1)$$

70 To capture the relations among the input sequence, the values V are weighted by the scores from Equation 1. Therefore,
71 we have

$$\text{SelfAttention}(Q, K, V) = A(Q, K) \cdot V \quad (2)$$

72 While in NLP each element in the sequence corresponds to a word, the same idea is applicable for a sequence of
73 N discrete objects, like pixels of an image. A key property of the self-attention model described above is that it is
74 equivariant to the input order, i.e. it gives the same output independent of how the N input tokens are shuffled. This is

75 quite problematic in cases where the order actually matters like in images. Hence, a positional encoding is learnt for
76 each token in the sequence and added before the self-attention operation. Hence, the Q and K vectors are derived from
77 the summation of input X and positional encoding P .

78 **2.2 Transformer Attention**

79 The Transformer network is an extension of the attention mechanism from Eqn. 2 based on the Multi-Head Attention
80 (MHA) operation. Rather than computing the attention once, the MHA operation computes it multiple times (heads).
81 This helps the transformer jointly attend to different information derived from each head. The output from each of these
82 heads are concatenated before projecting onto a final output dimension. A transformer layer also contains a residual
83 connection followed by a layer normalization. The overall operation can be summarized as:

$$\begin{aligned} \text{MHA}(Q, K, V, \text{heads}) &= \text{concat}_{\text{heads}}[A(Q, K, V)] \\ \text{Transformer} &= \text{LayerNorm}(\text{MHA} + \text{MLP}(\text{MHA})) \end{aligned} \quad (3)$$

84 **2.3 Positional Encoding for Images**

85 There are two types of positional encodings used in transformer-based architectures: absolute and relative encoding.
86 Absolute encodings assign a (fixed or learned) vector P_p to every pixel p whereas the relative positional encoding
87 [13] considers only the position difference between the query pixel (pixel we compute the representation of) and the
88 key pixel (pixel we attend to).

89 The authors from the paper have elegantly proved how the attention operation mimics a convolution. The main result
90 is the following: *A multi-head self-attention layer with N_h heads of dimension D_h , output dimension D_{out} and a*
91 *relative positional encoding of dimension $D_p \geq 3$ can express any convolutional layer of kernel size $\sqrt{N_h} * \sqrt{N_h}$*
92 *and $\min(D_h, D_{out})$ output channels.* In this proposed construction, the attention scores of each head must attend to
93 different relative pixel shifts within a kernel (Lemma 1 from the paper). The above condition is satisfied for the relative
94 positional encoding referred to as *Quadratic Encoding*. However, experiments suggest that a relative position encoding
95 learned by a neural network (*Learned Relative Position Encoding*) can also satisfy the conditions of the lemma. We
96 strongly urge the readers to refer to the original paper to get a complete understanding.

97 **3 Reproducibility**

98 The aim of this section is to validate the results claimed by the paper - to examine whether self-attention layers in
99 practice do actually learn to operate like convolutional layers when trained on the standard image classification task. For
100 all our experiments mentioned in this report, we use the CIFAR10 dataset to benchmark the performance of the model.

101 **3.1 Dataset - CIFAR10**

102 The CIFAR-10 dataset [14] consists of 60,000 colour images of size 32x32 split in 10 classes. There are 6,000 images
103 per class split into 5,000 training and 1,000 validation samples.

104 **3.2 Experiments and Results**

105 The results mentioned in the paper uses a fully-attentional model consisting of 6 multi-head self-attention layers each
106 with 9 heads. In all the experiments, the input image undergoes a 2x2 down-sampling operation to reduce its size. The
107 final image vector is derived by average-pooling the representations derived from the last layer and then passed to a
108 linear layer for classification. Please refer to Table. 4 (Appendix) for a detailed list of hyperparameters used in each
109 experiment. We closely refer to the official implementation and were able to reproduce all the results within 1% of the
110 reported value. Table. 1 compares the results mentioned in the paper and the ones obtained using our implementation.
111 Fig. 1a visualizes the test accuracy on CIFAR10 at every 10 epochs for each model and it is quite evident that fully
112 convolutional networks like ResNet18 tend to converge faster. The following subsections describe these results in detail.

Model	Accuracy (paper)	Accuracy (ours)	# of params
ResNet18	0.938	0.946	11.2M
Quadratic embedding (isotropic)	0.938	0.943	12.1M
Quadratic embedding (non-isotropic)	0.934	0.939	12.1M
Learned embedding w/o content	0.918	0.904	12.3M
Learned embedding w/ content	0.871	0.864	29.5M

Table 1: Test accuracy (paper vs ours) on CIFAR-10 and model sizes; 9 heads

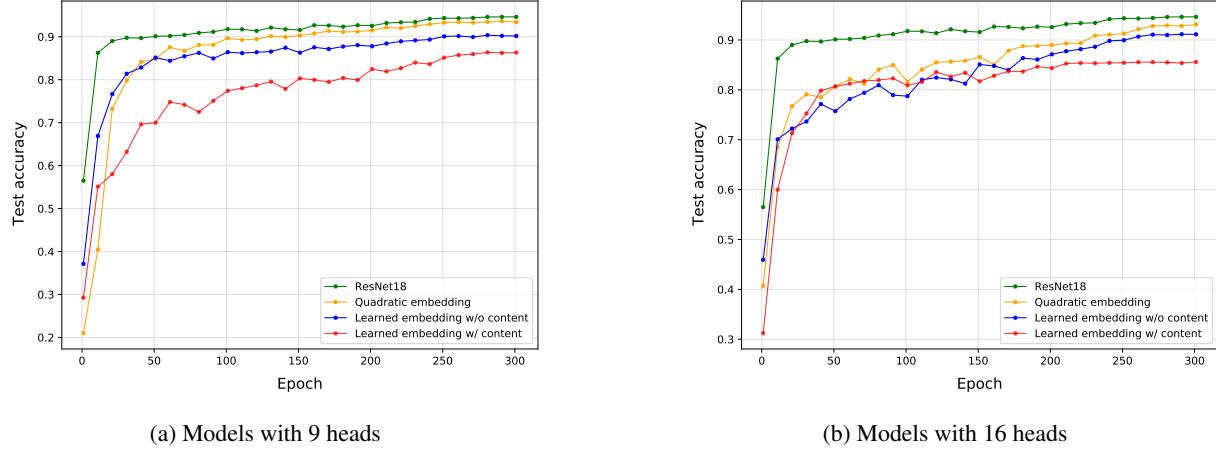


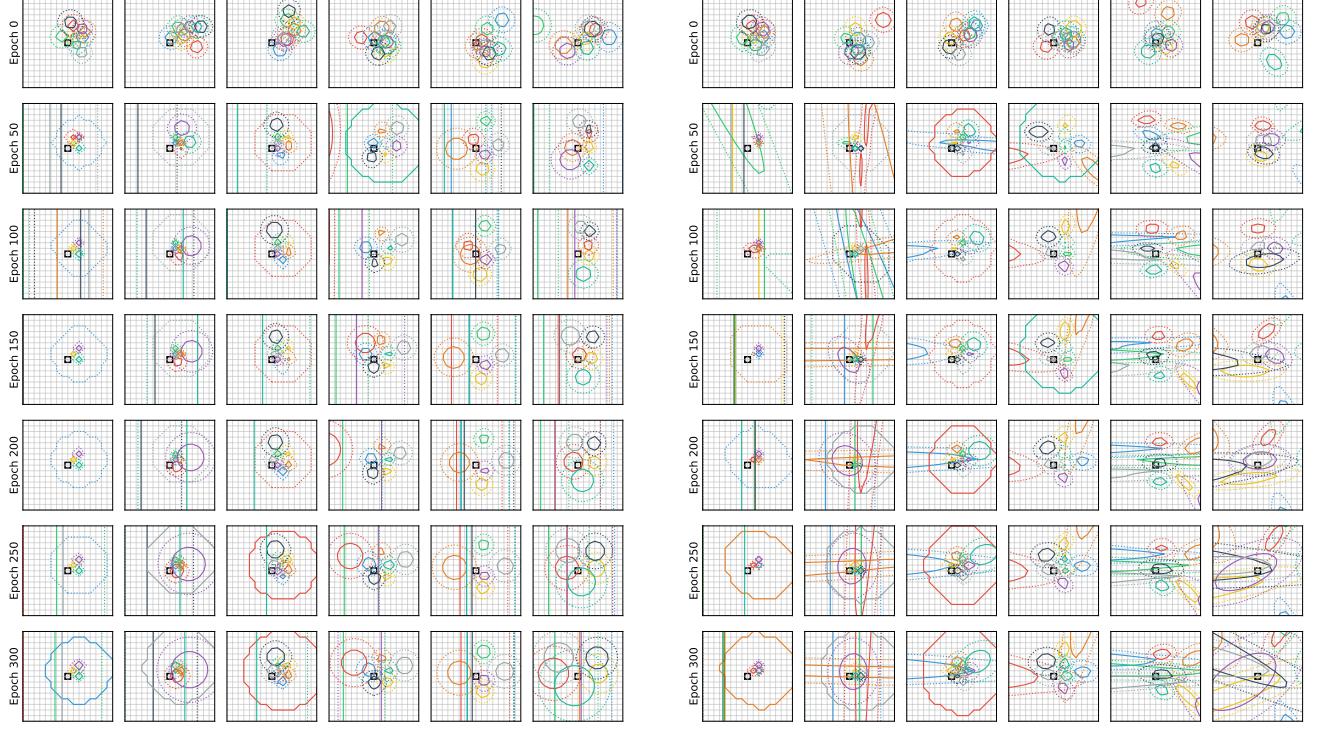
Figure 1: Test performance on CIFAR10 at every 10 epochs

113 3.2.1 Quadratic Encoding

114 The authors show that the attention probabilities in the quadratic positional encoding is similar to an isotropic bivariate
 115 Gaussian distribution with bounded support. Hence to validate their claims, all the attention matrices in the model are
 116 replaced with these Gaussian priors, with learnable parameters to determine the center and width of each attention
 117 head. Further, this is extended to a non-isotropic distribution over pixel positions as it might be interesting to see if the
 118 model would learn to attend to such groups of pixels - thus forming unseen representations in CNNs. Fig. 2 visualizes
 119 the attention centers for each head for all the layers and at different epochs. After optimization, we can see that the
 120 heads attend to a specific pixel of the image forming a grid around the query pixel. This confirms the intuition that
 121 self-attention applied to images learns convolution-like filters around the query pixel. Also, it can be seen that the
 122 initial layers (1-2) focus at local patterns while the deeper layers (3-6) attend to larger patterns by positioning the center
 123 of attention further from the queried pixel position. Fig. 2b shows that the network did learn non-isotropic attention
 124 patterns especially in the last layers. However, there is no performance improvement suggesting that it is not particularly
 125 helpful in practice.

126 3.2.2 Learned Relative Position Encoding

127 In this experiment, the authors try to study the positional encoding generally used in fully-attentional models [11].
 128 The positional encoding vector for each row and column pixel shift is learnt. The final relative position encoding
 129 of a key pixel with a query pixel is derived as the concatenation of row and column shift embeddings. First, the
 130 authors completely discard the input data and compute the attention weights solely with the derived encoding (Learned
 131 embedding w/o content). Fig. 3a visualizes the attention probabilities for a given query pixel, confirming the hypothesis
 132 that: even when left to learn positional encoding from randomly initialized vectors, certain self-attention heads learn to
 133 attend to individual pixels while the others learn non-localized patterns and long range dependencies. In another setting
 134 (Learned embedding w/ content), both the positional and content based attention information is used which corresponds
 135 to a full-blown stand alone self-attention model. Fig. 3b visualizes the attention probabilities for a given query pixel in
 136 this setting and it is interesting to note that even when left to learn the encoding from the data, some attention heads
 137 exploit positional information like CNNs while the others focus on the content. In Fig. 4, we visualize the attention



(a) isotropic Gaussian parameterization

(b) non-isotropic Gaussian parameterization

Figure 2: Centers of attention of each attention head (different colors) for all 6 layers (*columns*) at various training epochs (*rows*). The central black square is the query pixel, whereas solid and dotted circles represent the 50% and 90% percentiles of each Gaussian, respectively.

138 probabilities averaged across an entire batch of images to understand the focus of each head and remove dependency on
139 the input image for both the experiments.

140 **Average Attention Visualization:** The authors from the original paper visualize the attention probabilities for a single
141 image or across a batch images for a specific query pixel. A single query pixel does not convey information regarding
142 where the model focuses on in the entire image and it is not practical to plot individual figures for every query pixel.
143 Hence, we also visualize the attention probabilities using what we refer to as *Average Attention*, to identify which
144 portions of the entire image the model attends to. Given a softmax normalized attention matrix of size $N \times N$, every row
145 represents the relationships between a query pixel and the others. First, every element $(\alpha_{i,j})$ is divided by the row-wise
146 sum to ensure that all the values are in scale. Then the row-wise mean is computed to determine the importance value
147 for each pixel. If a pixel is strongly correlated to multiple pixels, the importance value will be higher determining
148 that the model has a stronger focus on that given pixel. We mathematically describe the operation in Eqn. 4. Fig. 5
149 visualizes the average attention for the learned embedding with and without content. In Fig. 5a, since the content data is
150 discarded, the model is clearly focusing on positional patterns while in Fig. 5b, the model attends to both positional and
151 content information. We visualize additional figures in Sections. B.1-B.3 (Appendix).

$$\begin{aligned} \alpha_{i,j} &= \text{softmax}(\alpha_{i,j}) = \frac{\exp \alpha_{i,j}}{\sum_k (\alpha_{i,k})} \\ \alpha_{i,j} &= \frac{\alpha_{i,j}}{\sum_k (\alpha_{k,i})} \\ \text{Avg Attn}_i &= \frac{\sum_k (\alpha_{i,k})}{N} \end{aligned} \tag{4}$$

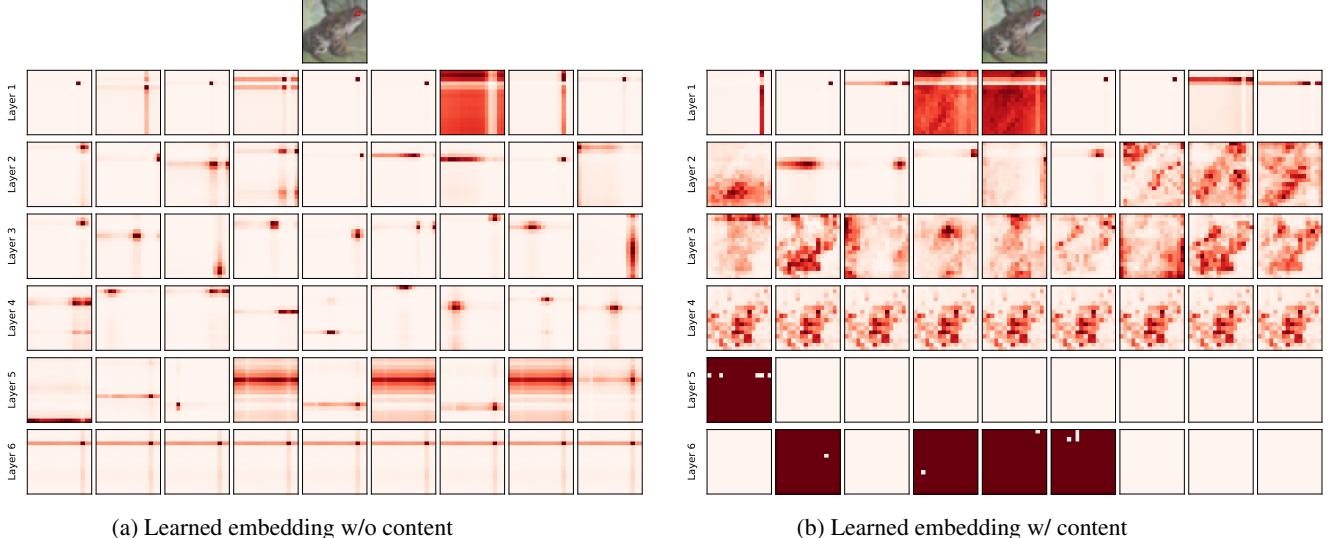


Figure 3: Attention probabilities for a model with 6 layers (rows) and 9 heads (columns) using learned relative positional encoding (with and without content). The query pixel (red square) is on the frog head.

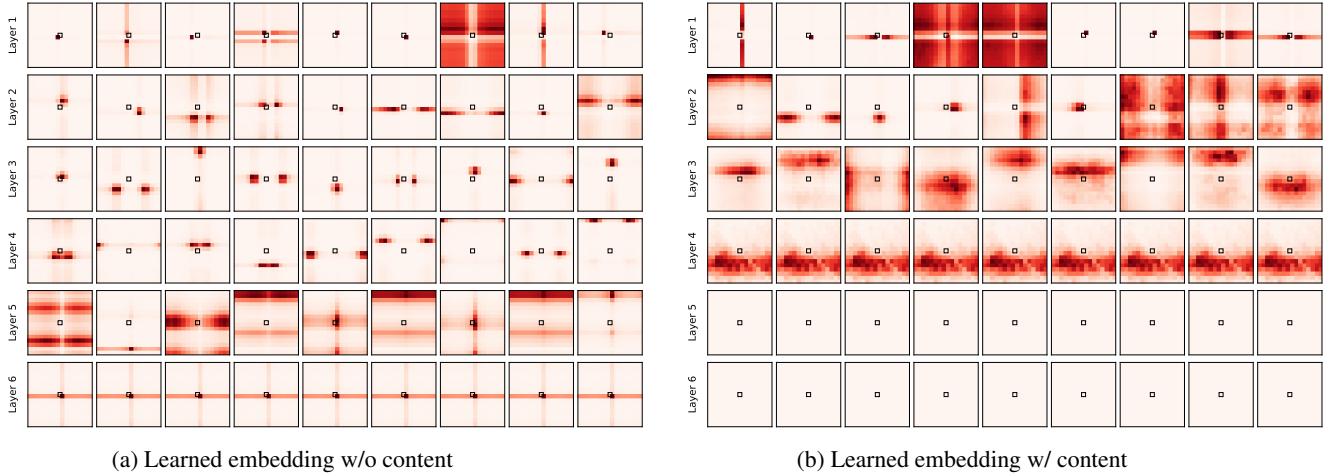
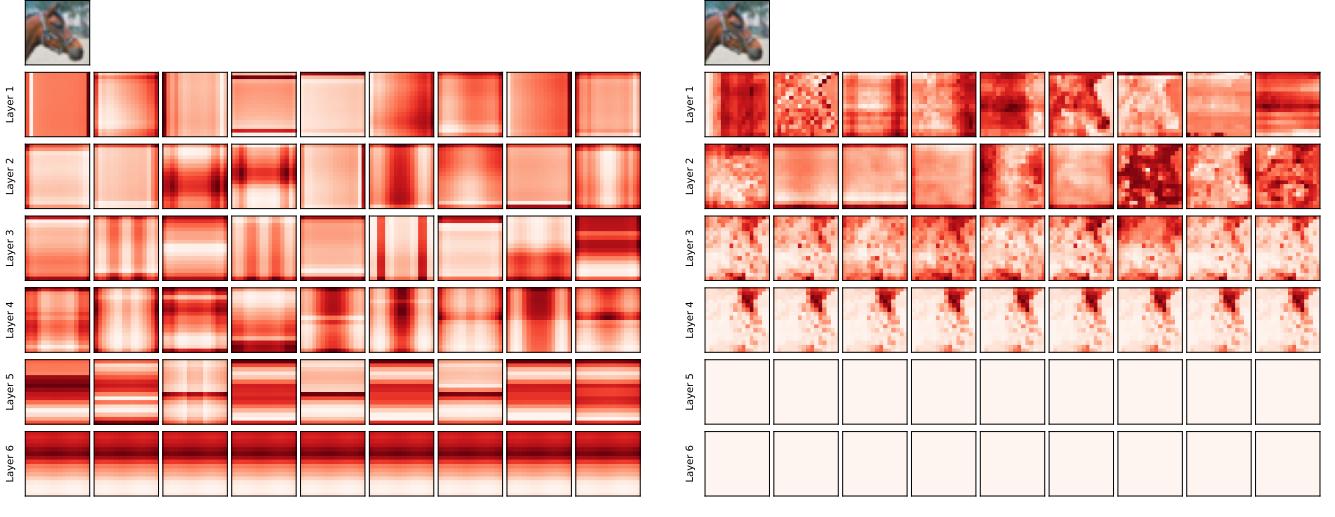


Figure 4: Attention probabilities for a model with 6 layers (rows) and 9 heads (columns) using learned relative positional encoding (with and without content). Attention maps are averaged over 50 test images to display head behavior and remove the dependence on the input content. The query pixel is in the center of the image.

152 3.3 Increasing the number of heads

153 As per the analogy derived between self-attention and convolutions, the number of heads is directly related to the kernel
 154 size in a convolution operation. Hence, we increase the number of heads from 9 to 16. It is important to note here that
 155 the unlike the general procedure of setting $D_h = D_{out}/N_h$ in transformer-based architectures, the paper suggests to
 156 concatenate heads of dimension $D_h = D_{out}$ since the effective number of learned filters is $\min(D_h, D_{out})$. Given the
 157 limited compute, we reduced D_{out} from 400 to 256 while increasing the number of heads to 16. As seen in Table. 2
 158 there seems to be no significant impact to the model's performance. However, the model takes longer to converge due to
 159 the increased number of parameters (Fig. 1b). We visualize the attention probabilities in Sections. B.4-B.6 (Appendix).



(a) Learned embedding w/o content

(b) Learned embedding w/ content

Figure 5: Average Attention visualization for a model with 6 layers (rows) and 9 heads (columns) using learned relative positional encoding (with and without content).

Model	Accuracy
Quadratic embedding	0.931
Learned embedding w/o content	0.912
Learned embedding w/ content	0.856

Table 2: Test accuracy on CIFAR-10; 16 heads

160 3.4 Additional Observations

161 3.4.1 Inductive biases in Transformers

162 As seen from the results mentioned in Table. 1, the fully attentional model utilizing learnable embeddings with image
 163 content performs poorly when compared to the other methods. As mentioned in the paper [3], transformers lack biases
 164 inherent to CNNs like translation equivariance and retention of 2D neighborhood structure, etc. Only the Multi-Layer
 165 Perceptron (MLP) layers used in these methods are local and translation equivariant, while the self attention layers are
 166 global. Even in the case of NLP, almost 75-90% of the predictions are correct even when the input words are randomly
 167 shuffled [15]. This implies that transformer-based methods do not sufficiently capture spatial information even with
 168 positional encodings and require a large amount of training data to do so. This could be the reason for improved
 169 performance in the case of Learned embedding w/o content and Quadratic embedding as the attention matrices are
 170 directly replaced with positional information.

171 3.4.2 Over-expressive power of attention matrices

172 The most important step of the self-attention operation is the generation of the attention matrix of size $N \times N$. In NLP,
 173 the value of N tends to be small (<100) in most cases as we are dealing with words in a sentence. On the contrary,
 174 images when broken down result in very long sequences of pixels, hence creating large attention matrices. Therefore,
 175 these attention matrices can be sparse and the model has a very high tendency of focusing on very high level information.
 176 This can lead to over-fitting and is talked about in the case of pointclouds⁶ where the number of points are very large
 177 (>1000). This has also been observed in our experiments. As seen in Figs. 3b, 4b, 5b, the attention heads in the last few
 178 2 layers are very sparse and do not capture any information. This can also be seen in the case of Quadratic encoding
 179 (Fig. 2), where certain attention heads focus on "non-intuitive" portions of the image (i.e) a thin strip of pixels or attend
 180 uniformly across a large patch of pixels. A simple and naive way to overcome this problem is to reduce the number of
 181 heads or layers but this is not an effective method as the model loses its capacity to learn strong features.

⁶https://github.com/juho-lee/set_transformer/issues/3#issuecomment-586711062

182 **4 Hierarchical Attention**

183 Given the problems described above, we need to propose a method which can avoid the over-expressive nature of
 184 independent attention heads while still being able to learn and derive strong features from the input image across layers.
 185 We now introduce a novel attention operation which we refer to as *Hierarchical Attention* (HA) operation. In the
 186 following sections, we explain the core idea behind the operation and perform detailed experiments to showcase its
 187 effectiveness.

188 **4.1 Methodology**

189 In most of the transformer-based methods, independent self-attention layers are stacked sequentially and the output
 190 from one is passed onto the next to derive the Q , K and V vectors. This allows each attention head to freely attend
 191 to specific features and derive a better representation. Deviating from these methods, the HA operation updates only
 192 the Q , V vectors while the K remains the same after each attention block. Further, the weights are shared across
 193 these attention layers inducing the transformer model to iteratively refine its representation across layers. This can
 194 be considered analogous to an unrolled recurrent neural network (RNN) as we are trying to sequentially improve the
 195 representation across layers based on the previous hidden state. This helps the model hierarchically learn complex
 196 features by focusing on corresponding portions of the K vector, and aggregating required information from the V vector.
 197 Figs. 6, 7 visualize the normal attention and the hierarchical attention operations respectively. For the sake of simplicity,
 198 we do not visualize the entire inner workings of the transformer like residual connections, layer normalization, etc. This
 199 is a very simple yet effective method which can be easily adapted to any existing attention-based network as described
 200 in the next section.

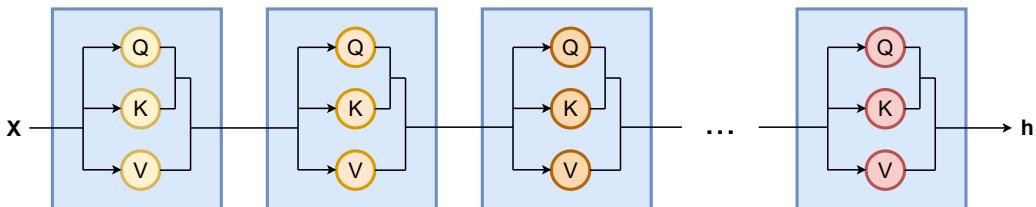


Figure 6: Normal Attention (Scaled Dot Product): The projection weights in each layer are different and independently learnt (different color). Q , K , V vectors are updated in each layer.

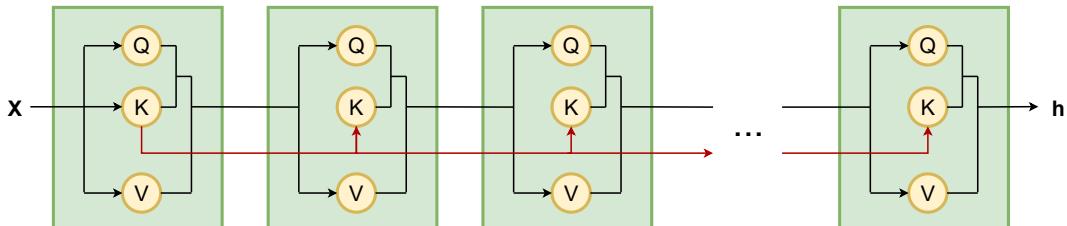


Figure 7: Hierarchical Attention (HA): The projection weights are shared across all the layers. Only Q , V vectors are updated while the K remain the same in each layer.

201 As mentioned earlier, there has been a lot of recent hype in replacing convolutions with attention layers. Hence, we
 202 choose two other popular and similar papers - "Exploring Self-attention for Image Recognition" [2], "An Image is
 203 worth 16x16 words: Transformers for Image Recognition at scale" [3] and apply the proposed HA operation to justify
 204 its effectiveness. For better understanding, we briefly introduce these papers in the following subsections.

205 **4.2 Pairwise and Patchwise self attention (SAN)**

206 Introduced by [2], pairwise self-attention is essentially a general representation of Eqn. 2. It is fundamentally a set
 207 operation, does not attach stationary weights to specific locations and is invariant to permutation and cardinality. The
 208 paper presents a number of variants of the pairwise attention that have greater expressive power than dot-product
 209 attention. Specifically, the weight computation does not collapse the channel dimension and allows the feature

210 aggregation to adapt to each channel. It can be mathematically formulated as follows:

$$y_i = \sum_{j \in R_i} \alpha(x_i, x_j) \odot \beta(x_j) \quad (5)$$

$$\alpha(x_i, x_j) = \gamma(\delta(x_i, x_j))$$

211 Here, i is the spatial index of feature vector x_i , $\delta(\cdot)$ is the relation function mapped onto another vector by $\gamma(\cdot)$. The
212 adaptive weight vectors $\alpha(x_i, x_j)$ aggregate feature vectors obtained from $\beta(\cdot)$. An important point to note here is that
213 δ can produce vectors of different dimensions when compared to β , allowing for a more expressive weight construction.

214 The patch-wise self attention is a variant of the pairwise operation, where x_i is replaced by a patch of feature vectors
215 $x_{R(i)}$ which allows the weight vector to incorporate information from all the feature vectors in the patch. The equations
216 are hence rewritten as:

$$y_i = \sum_{j \in R_i} \alpha(x_{R(i)})_j \odot \beta(x_j) \quad (6)$$

$$\alpha(x_{R(i)}) = \gamma(\delta(x_{R(i)}))$$

217 4.3 Vision Transformer

218 The Vision Transformer [3], has successfully shown that reliance on convolutions is no longer necessary and a pure
219 transformer outperforms all convolution based techniques significantly when pre-trained on large amounts of data. In
220 this method, an image is split into patches which is then projected onto another representation using a trainable layer
221 and then passed through a standard set of transformer operations as described in Eqn. 3.

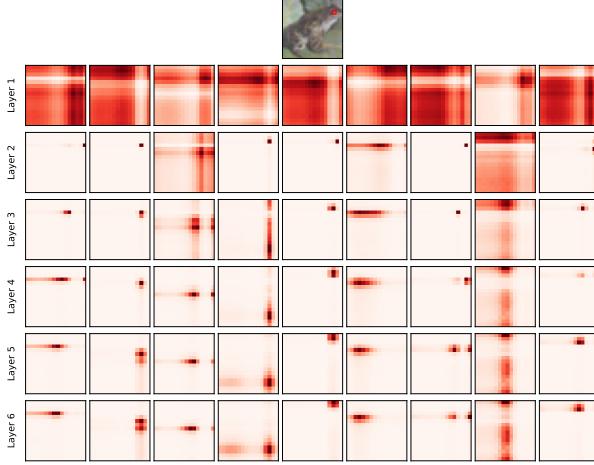
222 4.4 Results

223 To validate our intuition and to prove the effectiveness of the proposed method, we incorporate HA in all the methods
224 described above without modifying the overall structure of the architecture. Table 3 compares accuracy for each model
225 with its corresponding HA variant. We see significant improvements in the performance (at least 5% gain) in each
226 case while being able to reduce the number of parameters to almost 1/5th of the original model. As mentioned earlier,
227 transformers require sufficient training to perform equally well as convolution-based architectures. When pretrained on
228 large datasets (14M-300M images), transformer-based architectures achieve excellent performance and transfer to tasks
229 with fewer datapoints [3]. However, for all our experiments in this report we only focus on training these models from
230 scratch on the CIFAR10 dataset.

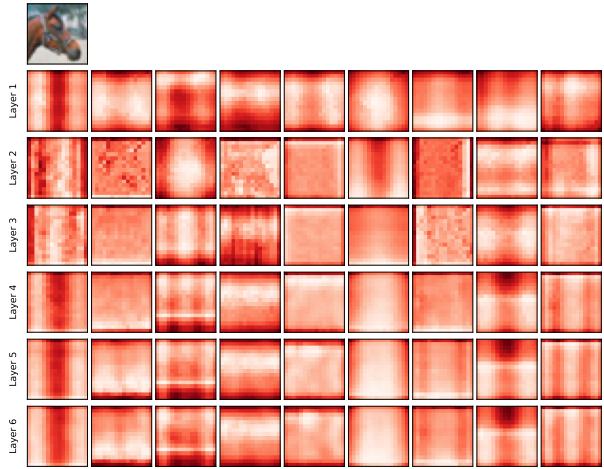
Model	Normal SA			Hierarchical SA		
	Accuracy	# of params	Wall time	Accuracy	# of params	Wall time
Learned embedding w/ content	0.871	29.5M	5.374	0.910	4.75M	4.703
SAN Pairwise (subtraction)	0.818	3.74M	16.586	0.853	0.35M	12.913
SAN Patchwise (subtraction)	0.857	3.74M	9.773	0.912	0.36M	7.419
Vision Transformer	0.843	37.4M	7.094	0.877	6.53M	6.302

Table 3: Comparison between models using normal SA and HA. Wall times are average inference times in milliseconds for the models over 300 iterations

231 In Fig. 8a, we visualize the attention probabilities for a given query pixel. The relationship between self-attention and
232 convolutions is striking as the model is attending to distinct pixels at a fixed shift from the query pixel reproducing the
233 receptive field of the convolution operation. The initial layers attend to local patterns while the deeper layers focus on
234 larger patterns positioned further away from the query pixel. Similarly, in Fig. 8b, the attention heads from the last two
235 layers are no longer sparse and help capture more information. Hence the visual correctness verifies the operation and
236 its increased performance. We also visualize the attention probabilities for SAN, VIT in Sections. B.7-B.8, B.9-B.10
237 (Appendix) respectively.



(a) Attention probabilities for a given query pixel. The query pixel (red square) is on the frog head.



(b) Average Attention visualization

Figure 8: Attention probabilities for a model with 6 layers (rows) and 9 heads (columns) using Hierarchical Learned 2D embedding w/ content

238 We summarize the hierarchical operation as follows:

- 239 1. Enable weight sharing between the layers of the model by reusing the K vector and updating the Q, V vectors
240 only. This helps the model progressively extract higher level features.
241 2. The method of progressive refinement ensures that the attention matrices do not get over expressive. This leads
242 to a significant improvement over the corresponding non-hierarchical cases.
243 3. The total number of parameters in the model is independent of the number of layers and this property helps
244 significantly reduce the number of parameters when compared to the non-hierarchical versions. Also, this
245 helps make deeper models without worrying about memory constraints.
246 4. Even if the model is provided more layers than are necessary, every layer learns to attend to a different
247 pattern. The new features learnt at every layer add on to those learnt by the previous ones, which provides its
248 characteristic hierarchical nature. By visualizing the attention scores on a test image, we obtain convincing
249 evidence to support this hypothesis.

250 **5 Conclusion**

251 In this report, we study the application of self-attention for image recognition, specifically image classification. We
252 validate the original paper's claims by performing detailed experiments on the CIFAR10 dataset. We discuss two
253 recent papers which propose stand-alone attention based models for image classification. Finally, we introduce a novel
254 Hierarchical Attention operation which significantly improves performance across various models while also reducing
255 the number of parameters. The preliminary results raises various questions: do we actually need multiple independent
256 layers in large transformers? Does this improved performance also translate to large datasets and across various other
257 image recognition tasks like object detection and image segmentation? We would like to answer all these questions and
258 provide a more rigorous understanding of the proposed method in the future.

259 **References**

- 260 [1] Jean-Baptiste Cordonnier, Andreas Loukas, and Martin Jaggi. On the relationship between self-attention and
261 convolutional layers. *CoRR*, abs/1911.03584, 2019.
- 262 [2] Hengshuang Zhao, Jiaya Jia, and Vladlen Koltun. Exploring self-attention for image recognition, 2020.
- 263 [3] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner,
264 Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An
265 image is worth 16x16 words: Transformers for image recognition at scale, 2020.
- 266 [4] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In
267 Yoshua Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015,*
268 *San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.
- 269 [5] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott E. Reed, Dragomir Anguelov, Dumitru Erhan,
270 Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. *CoRR*, abs/1409.4842, 2014.
- 271 [6] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *CoRR*,
272 abs/1512.03385, 2015.
- 273 [7] Jie Hu, Li Shen, Samuel Albanie, Gang Sun, and Enhua Wu. Squeeze-and-excitation networks, 2019.
- 274 [8] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and
275 Illia Polosukhin. Attention is all you need, 2017.
- 276 [9] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko.
277 End-to-end object detection with transformers, 2020.
- 278 [10] Xiaolong Wang, Ross B. Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. *CoRR*,
279 abs/1711.07971, 2017.
- 280 [11] Prajit Ramachandran, Niki Parmar, Ashish Vaswani, Irwan Bello, Anselm Levskaya, and Jonathon Shlens.
281 Stand-alone self-attention in vision models. *CoRR*, abs/1906.05909, 2019.
- 282 [12] Huiyu Wang, Yukun Zhu, Bradley Green, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. Axial-deeplab:
283 Stand-alone axial-attention for panoptic segmentation, 2020.
- 284 [13] Zihang Dai, Zhilin Yang, Yiming Yang, Jaime G. Carbonell, Quoc V. Le, and Ruslan Salakhutdinov. Transformer-
285 xl: Attentive language models beyond a fixed-length context. *CoRR*, abs/1901.02860, 2019.
- 286 [14] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, 2009.
- 287 [15] Thang M. Pham, Trung Bui, Long Mai, and Anh Nguyen. Out of order: How important is the sequential order of
288 words in a sentence in natural language understanding tasks?, 2020.

289 **Appendix**290 **A Hyperparameters**

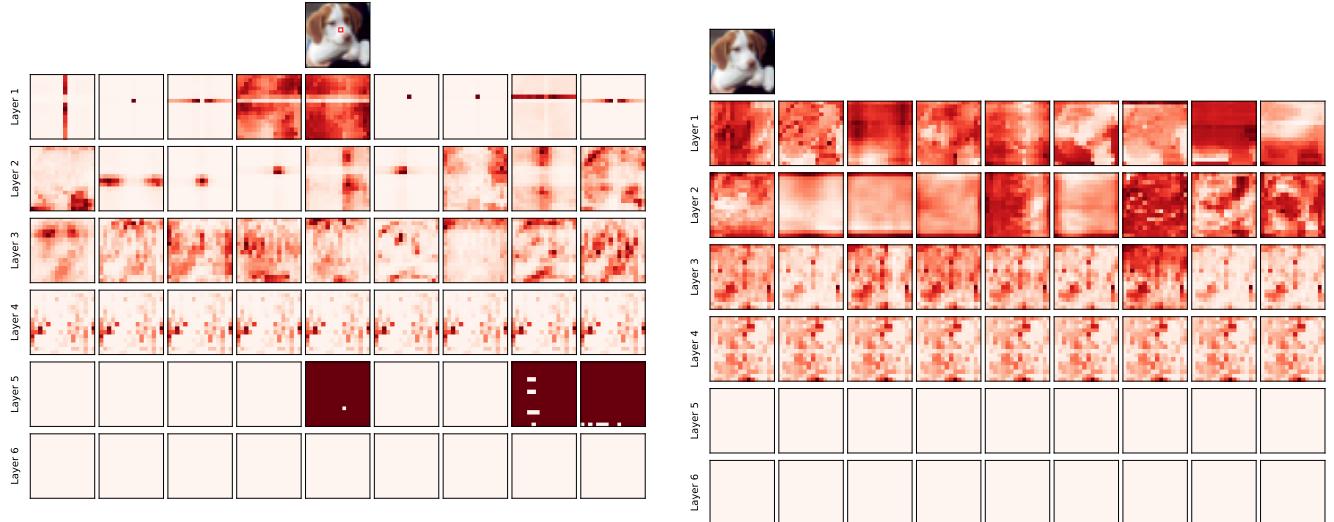
Model	Epochs	Mini batch size	Optimizer		Scheduler	
			Type	LR	Type	Warm up/step epochs
ResNet18	300	100	SGD	0.1	Cosine	10
Quadratic emb.	300	100	SGD	0.1	Cosine	15
Learned emb. w/o content	300	100	SGD	0.1	Cosine	15
Learned emb. w/ content (SA/HA)	300	100	SGD	0.1	Cosine	15
SAN Pairwise (SA)	300	256	SGD	0.1	Multistep	100, 200, 250
SAN Patchwise (SA)	300	256	SGD	0.1	Multistep	100, 200, 250
Vision Transformer (SA)	100	32	SGD	0.01	Cosine	5
SAN Pairwise (HA)	300	40	SGD	0.1	Cosine	10
SAN Patchwise (HA)	300	100	SGD	0.1	Cosine	10
Vision Transformer (HA)	100	32	SGD	0.01	Cosine	5

Table 4: Hyperparameter configuration for all experiments. SA refers to normal self-attention and HA refers to Hierarchical Attention. The momentum and weight decay for SGD were set to 0.9 and 0.0001 respectively for all experiments.

291 **B Attention Visualization**

292 We present more examples for visualising the attention probabilities in various models.

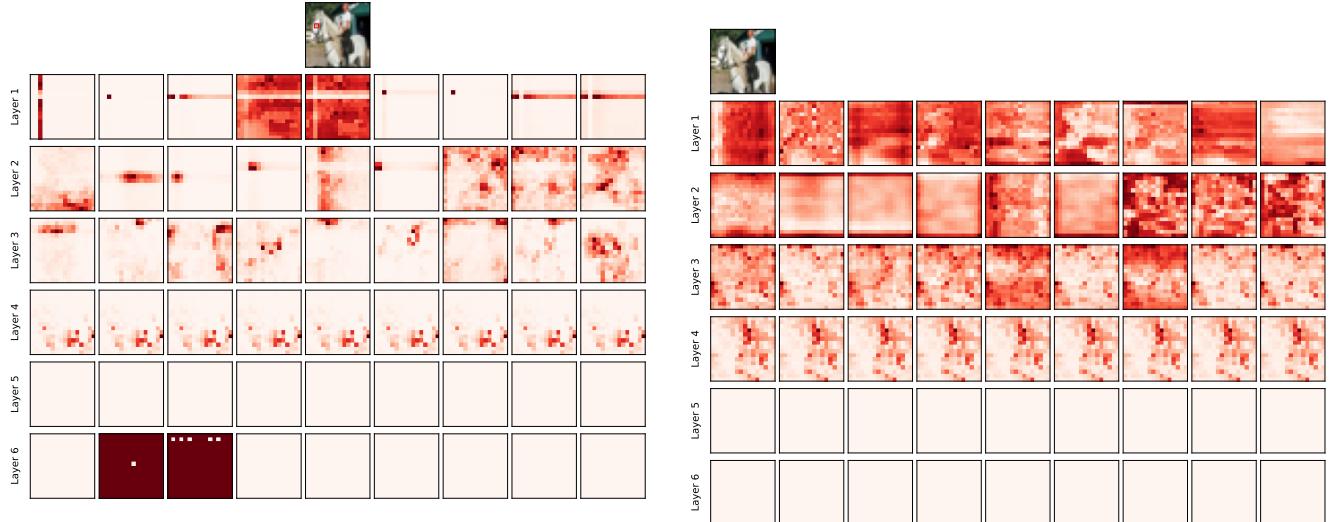
293 **B.1 Learned embedding w/ content; 9 heads**



(a) Attention probabilities for a given query pixel. The query pixel (red square) is on the dog head.

(b) Average Attention visualization

Figure 9: Attention probabilities for a model with 6 layers (rows) and 9 heads (columns) using learned embedding w/ content.

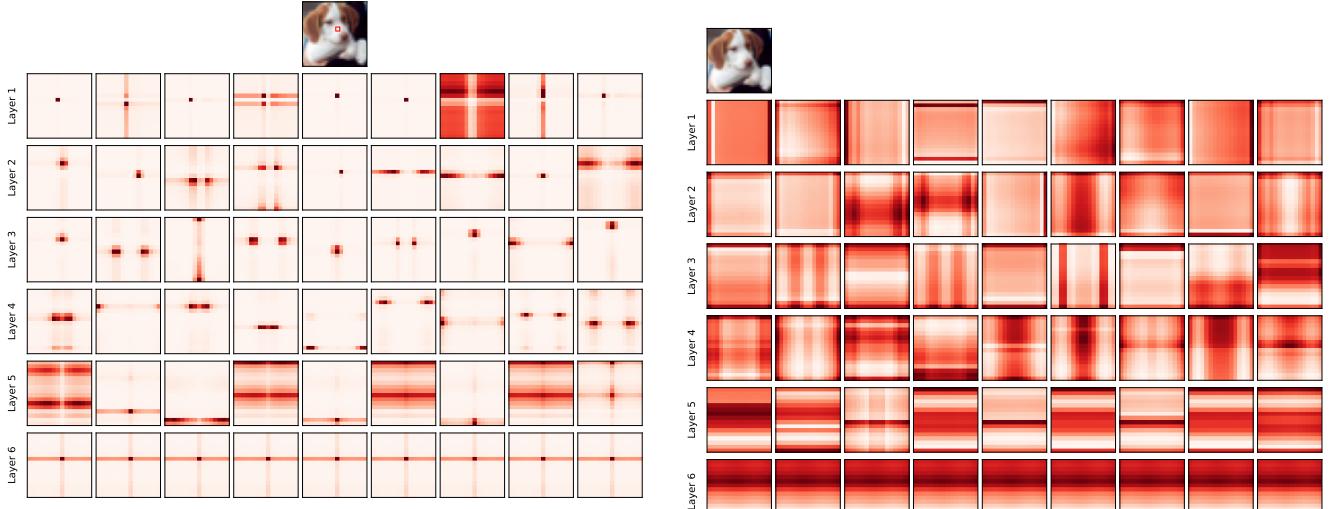


(a) Attention probabilities for a given query pixel. The query pixel (red square) is on the horse head.

(b) Average Attention visualization

Figure 10: Attention probabilities for a model with 6 layers (rows) and 9 heads (columns) using learned embedding w/ content.

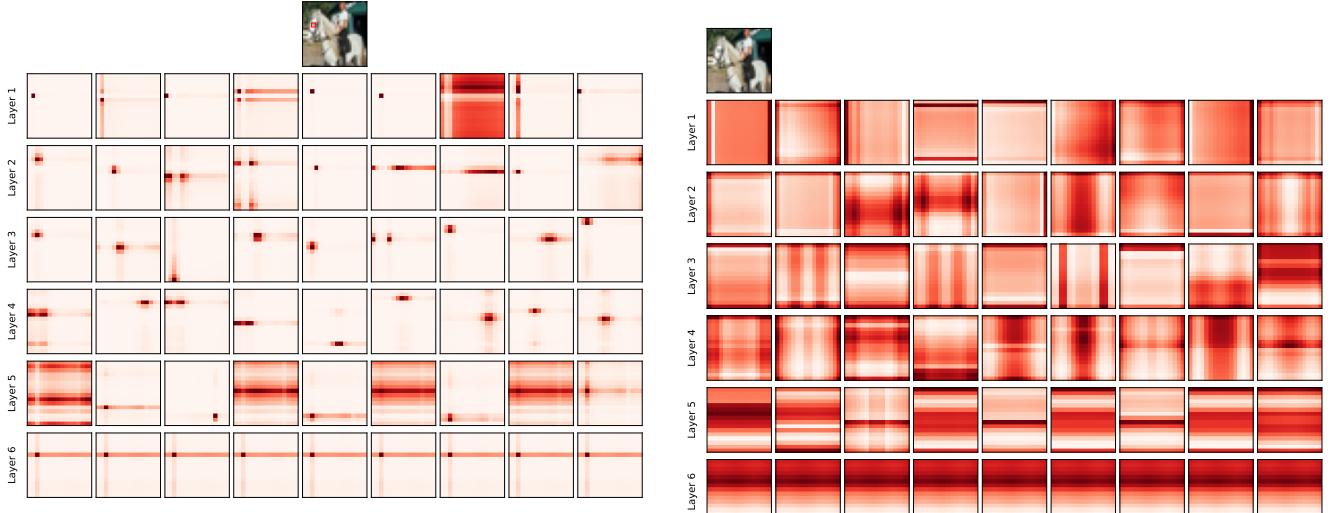
294 B.2 Learned embedding w/o content; 9 heads



(a) Attention probabilities for a given query pixel. The query pixel (red square) is on the dog head.

(b) Average Attention visualization

Figure 11: Attention probabilities for a model with 6 layers (rows) and 9 heads (columns) using learned embedding w/o content.

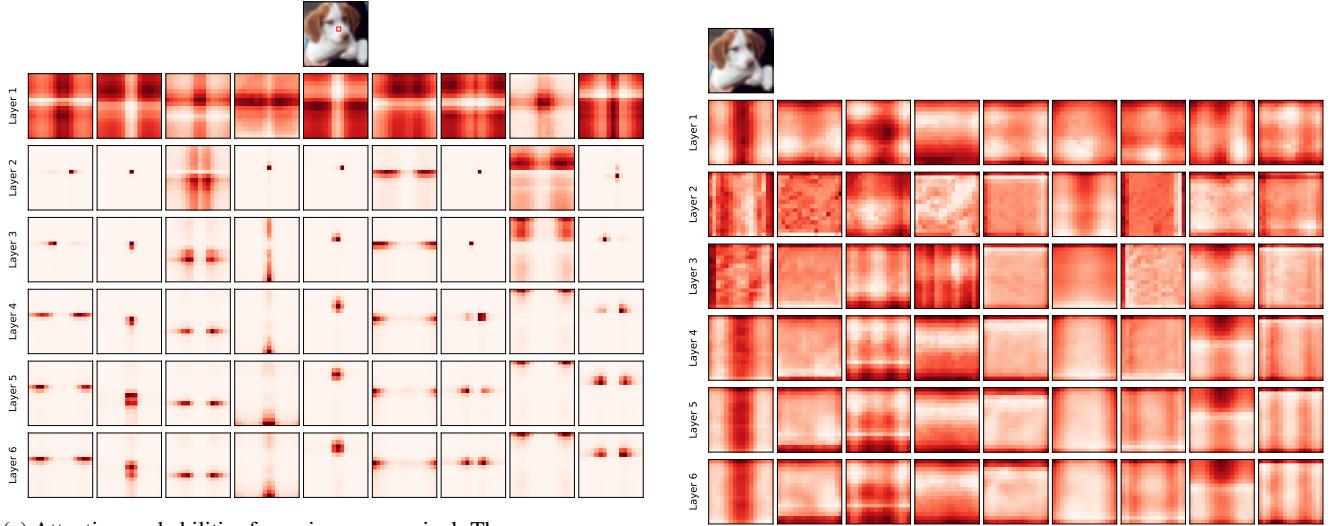


(a) Attention probabilities for a given query pixel. The query pixel (red square) is on the horse head.

(b) Average Attention visualization

Figure 12: Attention probabilities for a model with 6 layers (rows) and 9 heads (columns) using learned embedding w/o content.

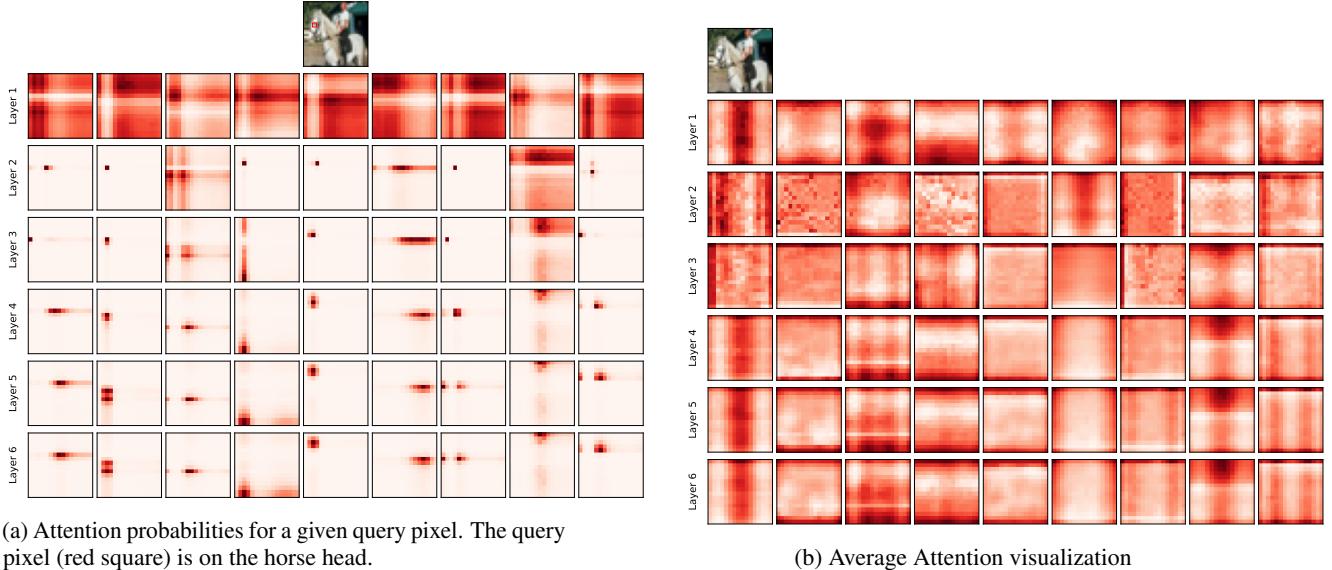
295 **B.3 Hierarchical Learned embedding w/ content; 9 heads**



(a) Attention probabilities for a given query pixel. The query pixel (red square) is on the dog head.

(b) Average Attention visualization

Figure 13: Attention probabilities for a model with 6 layers (rows) and 9 heads (columns) using hierarchical learned embedding w/ content.

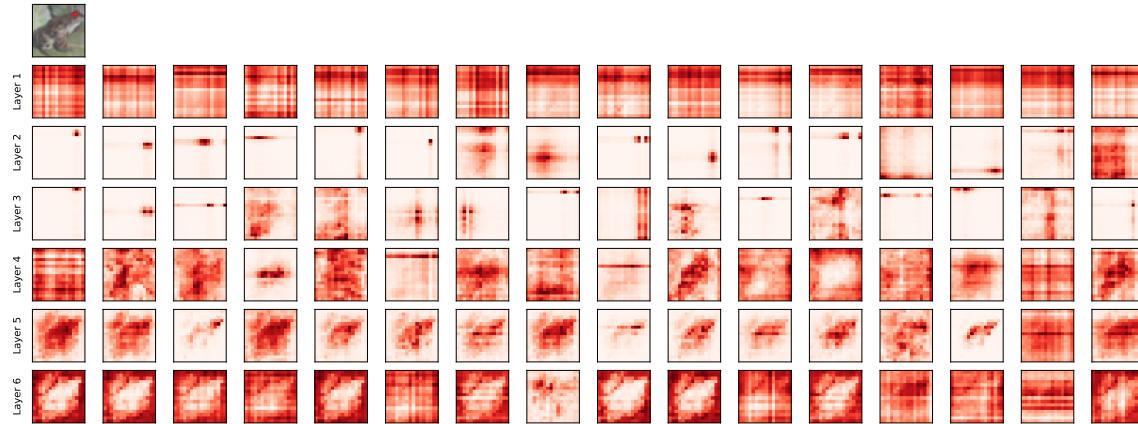


(a) Attention probabilities for a given query pixel. The query pixel (red square) is on the horse head.

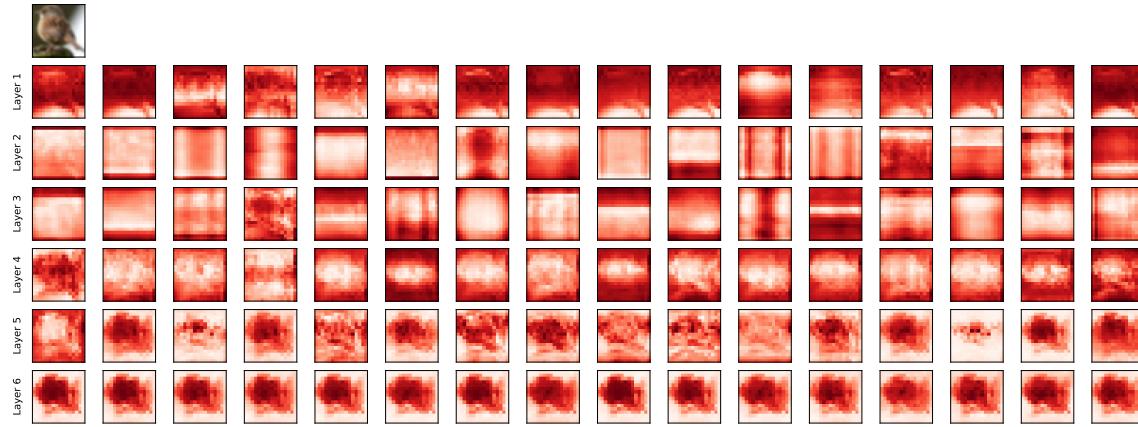
(b) Average Attention visualization

Figure 14: Attention probabilities for a model with 6 layers (rows) and 9 heads (columns) using hierarchical learned embedding w/ content.

296 **B.4 Learned embedding w/ content; 16 heads**



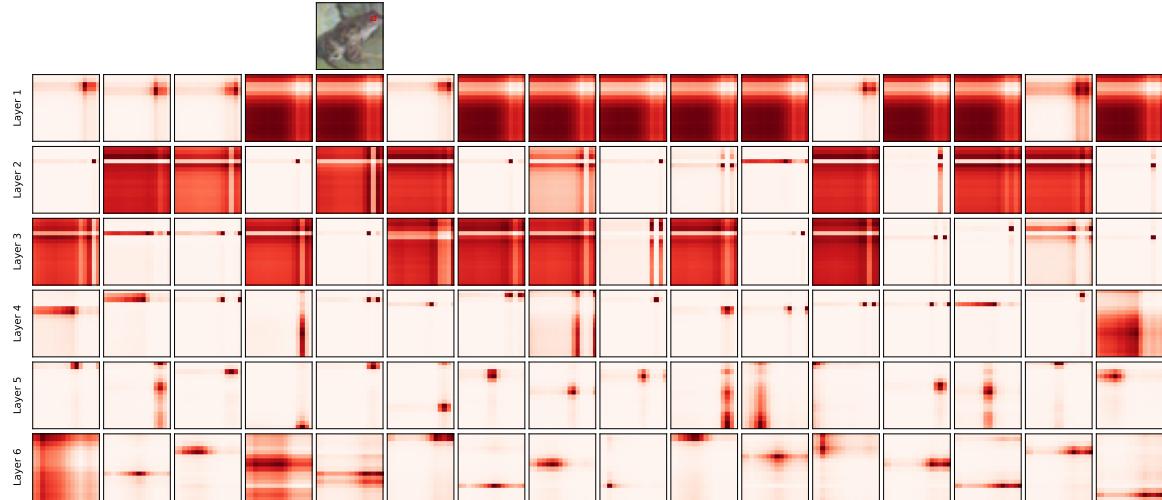
(a) Attention probabilities for a given query pixel. The query pixel (red square) is on the frog head.



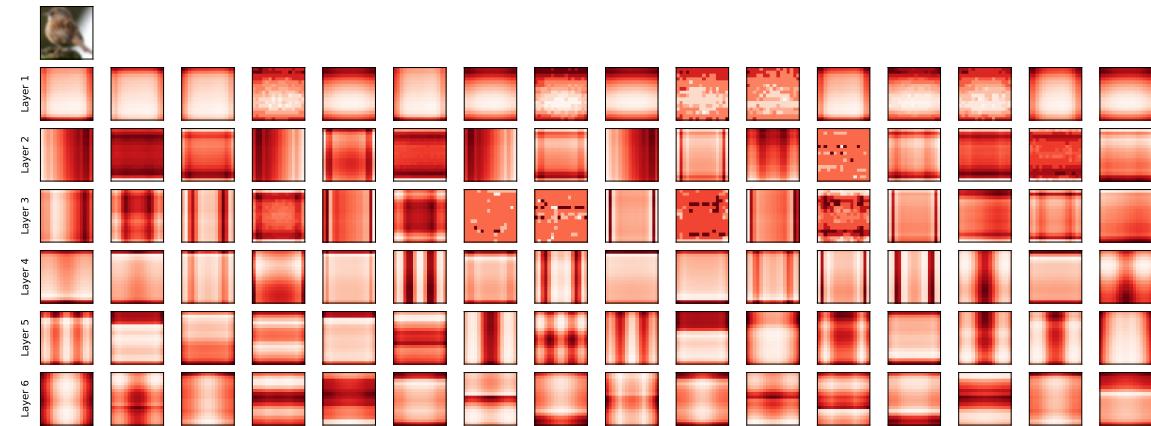
(b) Average Attention visualization

Figure 15: Attention probabilities for a model with 6 layers (rows) and 16 heads (columns) using learned embedding w/ content.

297 **B.5 Learned embedding w/o content; 16 heads**



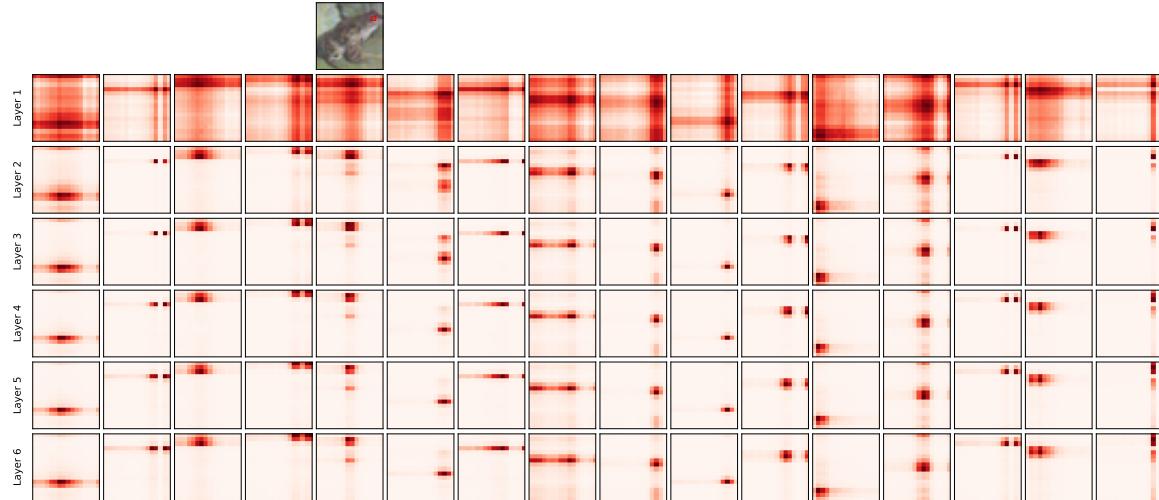
(a) Attention probabilities for a given query pixel. The query pixel (red square) is on the frog head.



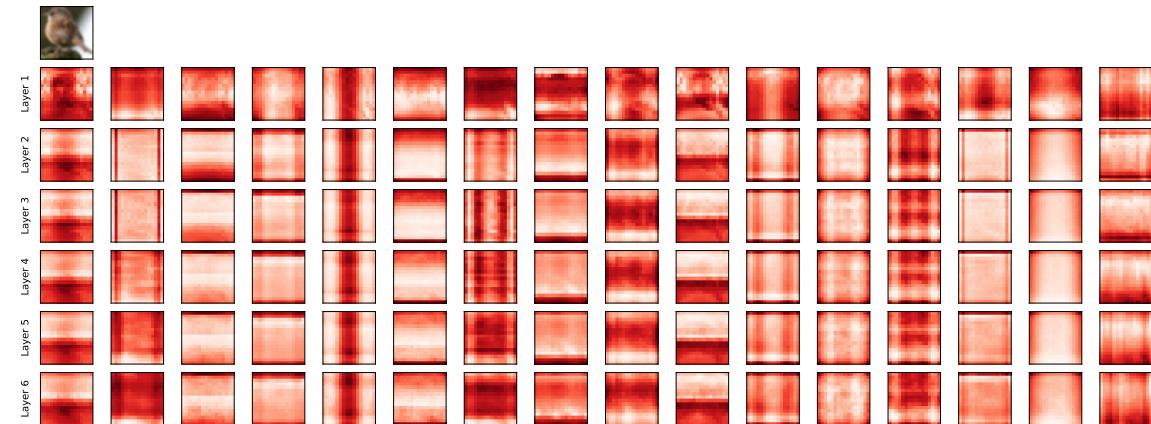
(b) Average Attention visualization

Figure 16: Attention probabilities for a model with 6 layers (rows) and 16 heads (columns) using learned embedding w/o content.

298 **B.6 Hierarchical Learned embedding w/ content; 16 heads**



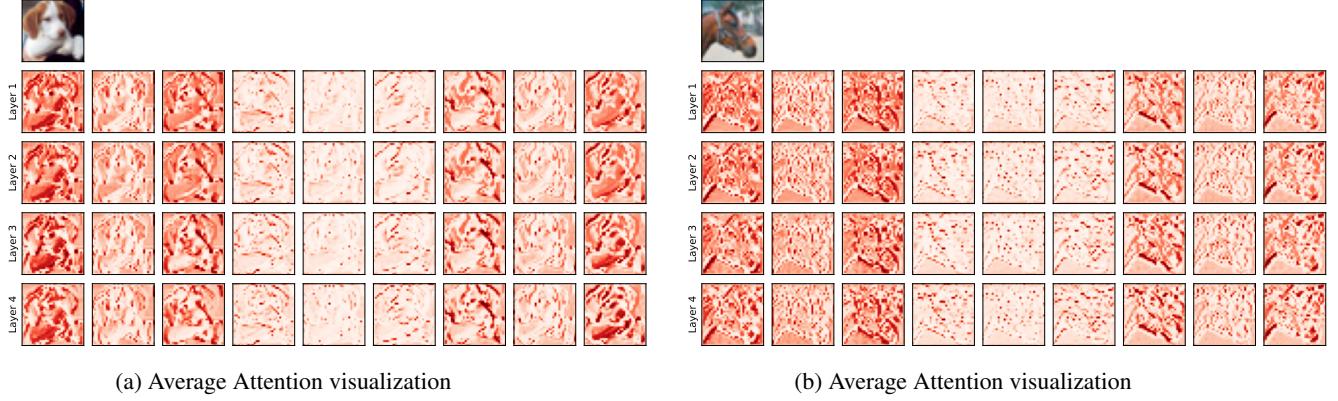
(a) Attention probabilities for a given query pixel. The query pixel (red square) is on the frog head.



(b) Average Attention visualization

Figure 17: Attention probabilities for a model with 6 layers (rows) and 16 heads (columns) using hierarchical learned embedding w/ content.

299 **B.7 Hierarchical SAN Pairwise**

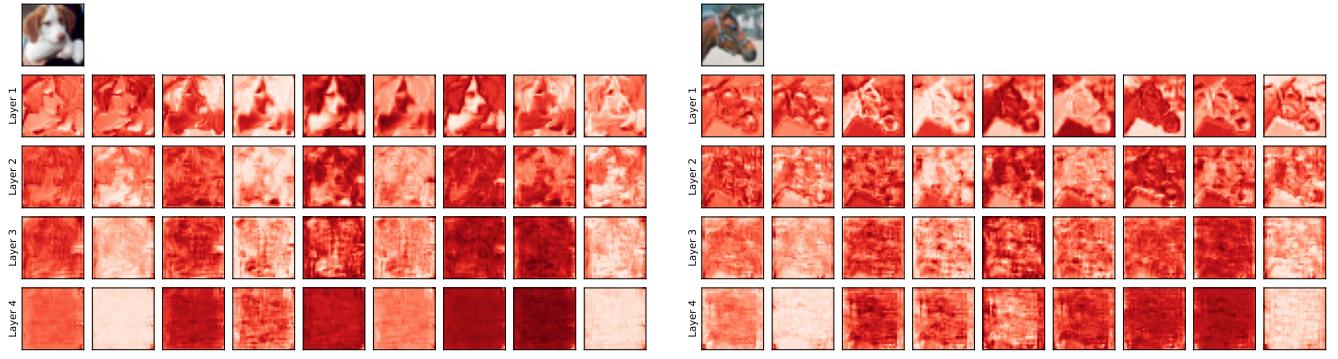


(a) Average Attention visualization

(b) Average Attention visualization

Figure 18: Attention probabilities for a model with 4 layers (rows) and 9 heads (columns) using hierarchical SAN Pairwise.

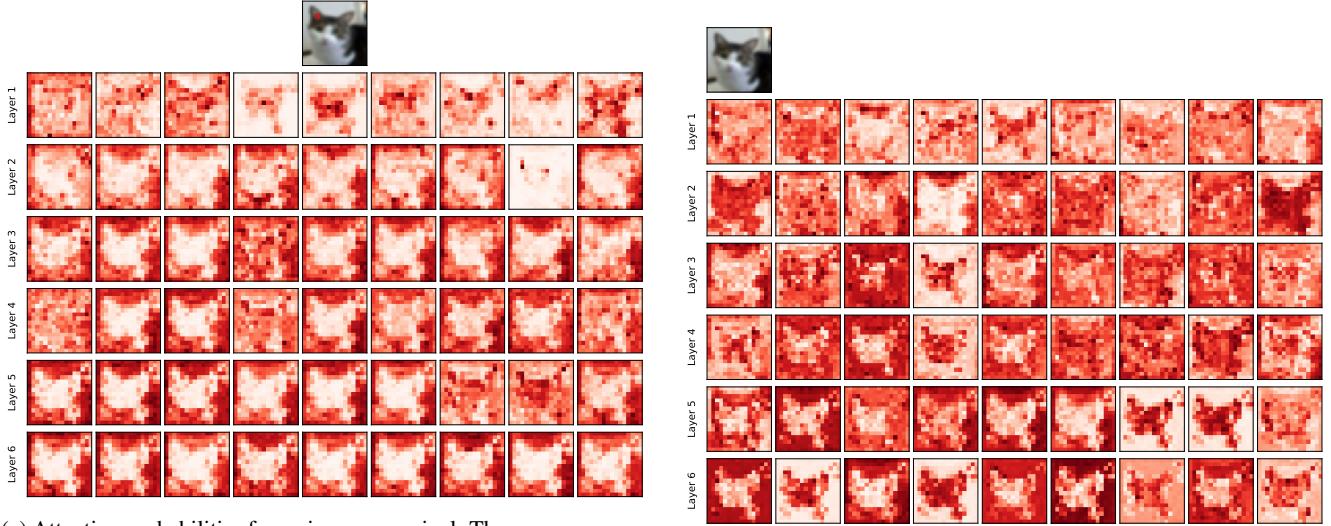
300 **B.8 Hierarchical SAN Patchwise**



(a) Average Attention visualization

(b) Average Attention visualization

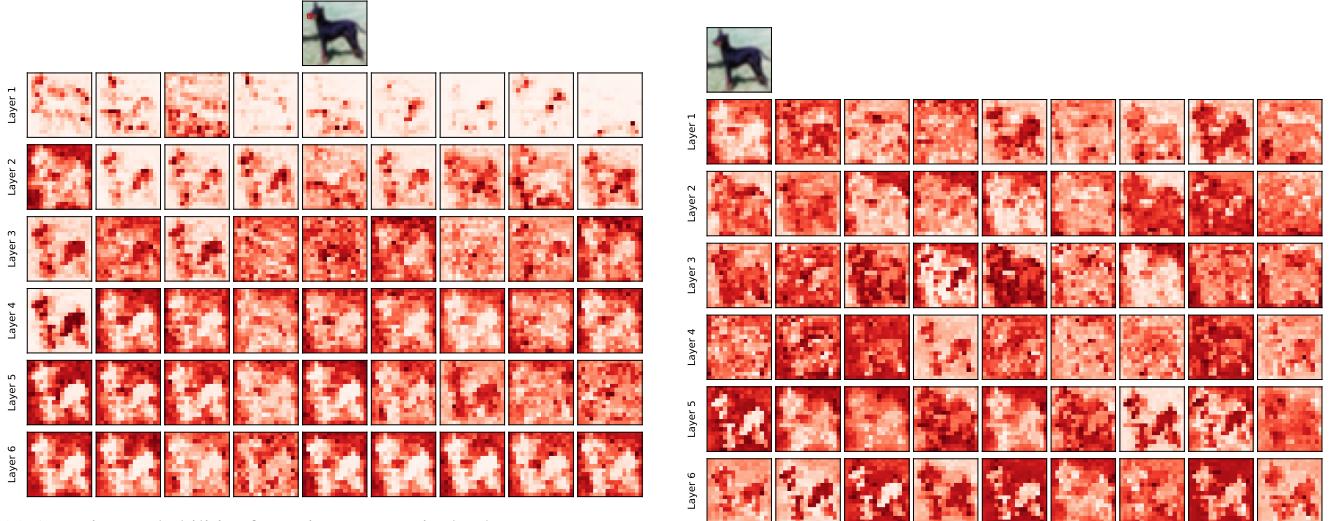
Figure 19: Attention probabilities for a model with 4 layers (rows) and 9 heads (columns) using hierarchical SAN Patchwise.



(a) Attention probabilities for a given query pixel. The query (red square) is on the cat's ear

(b) Average Attention visualization

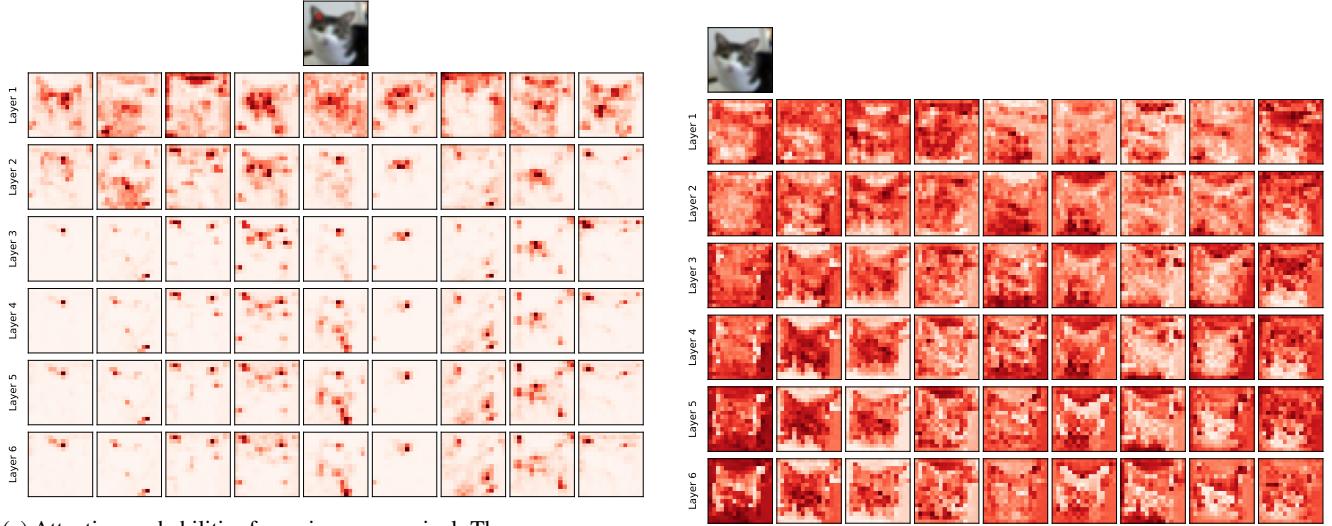
Figure 20: Attention probabilities for a model with 6 layers (rows) and 9 heads (columns) using ViT with patch size 2×2 .



(a) Attention probabilities for a given query pixel. The query (red square) is on the dog's snout

(b) Average Attention visualization

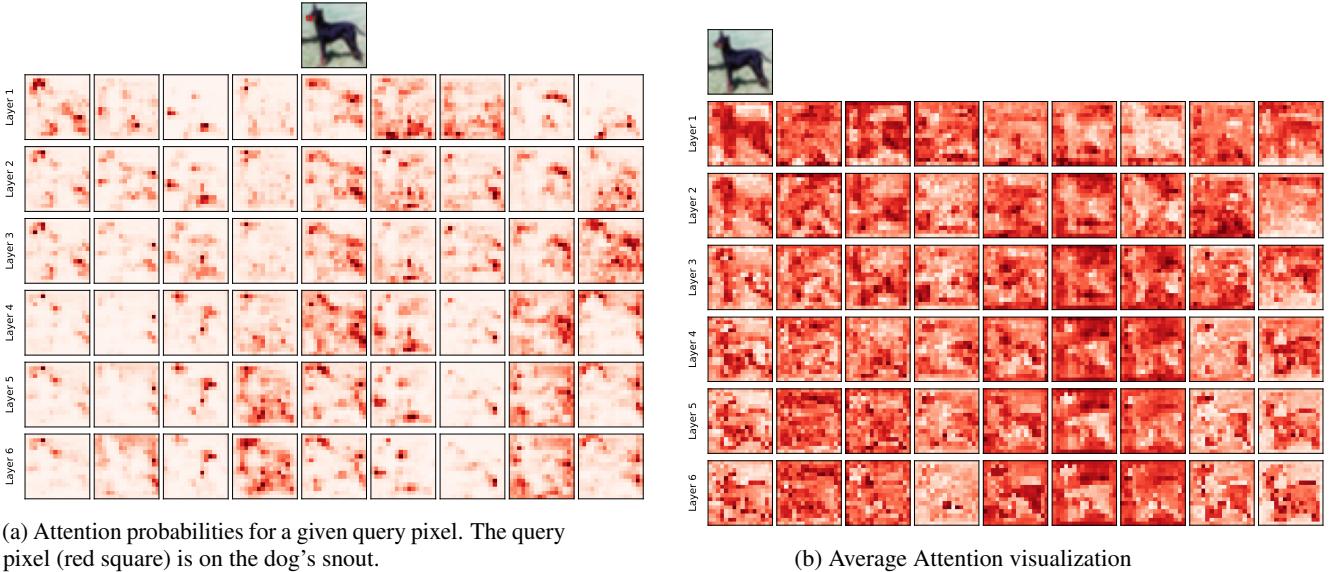
Figure 21: Attention probabilities for a model with 6 layers (rows) and 9 heads (columns) using ViT with patch size 2×2 .



(a) Attention probabilities for a given query pixel. The query pixel (red square) is on the cat's ear.

(b) Average Attention visualization

Figure 22: Attention probabilities for a model with 6 layers (rows) and 9 heads (columns) using hierarchical ViT with patch size 2×2 .



(a) Attention probabilities for a given query pixel. The query pixel (red square) is on the dog's snout.

(b) Average Attention visualization

Figure 23: Attention probabilities for a model with 6 layers (rows) and 16 heads (columns) using hierarchical ViT with patch size 2×2 .