1 **Rebuttal for ID 41**. We would like to thank the reviewers for their time and thoughtful comments.

2 **[R1]** *"The goal of the replacing convolutions with local self-attention is a bit in contradiction..."* Self-attention has
3 several advantages over convolutions, even when restricted locally.

- 4 • In contrast with convolutions where each position shares the same kernel weights, multi-head self-attention
  5 generates local kernels that can have different weights per position due to the computation depending on
  6 content-content interactions.

- 7 • Additionally, local self-attention is more parameter efficient than convolutions: using a $7 \times 7$ local self-attention
  8 layer outperforms using a $3 \times 3$ convolution while having $3\times$ fewer parameters. Furthermore, the $7 \times 7$ local
  9 self-attention layer has $2.4\times$ fewer FLOPs than a $3 \times 3$ convolutional layer.

10 **[R1]** *"it seems that the parameters of the convolutional layer have been replaced by the parameters ... $W_q$, $W_k$ and*
11 $W_v$.*"* We view the attention mechanism as a method for manufacturing convolutional kernels based on the content of
12 a given location. In some sense, this is a relaxation of locally-connected layers by not requiring that the kernels be
13 identical across spatial locations. The relaxation goes further by allowing the weights themselves to depend dynamically
14 on the content of each image.

15 **[R1]** *"Why are positional features important for self-attention"* Without positional information, attention will not be
16 sensitive to the ordering of the pixels because it will only use content-content interactions. Convolutions implicitly
17 carry a relative positional encoding by having weights that depend on relative distance.

18 **[R1]** *"One could consider that using only the positional interaction is a degenerated form of convolution"* We agree that
19 the importance of the content-relative interaction is surprising and concur in its similarity to traditional convolution, but
20 expect that in future investigations on more challenging tasks than classification the relative importance of content-
21 content interactions will increase.

22 **[R1]** *"throrough comparison of CNN and the proposed self-attention from a computational point of view [...] and the*
23 *expected behaviour and properties.* For kernel size $k$, channels $d$, convolution cost scales as $k^2 d^2$ FLOPs per position
24 with $k^2 d^2$ parameters, while self-attention cost scales as $3d^2 + k^2 d + kd$ FLOPs per position with $d^2$ parameters.

25 **[R3]** *"Why on Table 1 for ResNet-50 is full attention better than convolution-stem + attention?"* In the cases where full
26 attention outperforms convolutional-stem with attention, the difference is small ($\leq 0.2\%$) and can likely be explained
27 by variance in training runs. In the final version, we will add error bars to capture the variance.

28 **[R3]** *"...enlarging the spatial extent $k$ in attention improves performance but plateaus off at $11 \times 11$..."* We will
29 experiment with larger $k$ in the final version. We suspect that the effect of changing $k$ is task dependent.

30 **[R3]** *"What if you have binary/illusory/sketch images where you may need attention in the first place?"* While this
31 work is focused on demonstrating the attention can be used as a fundamental primitive for building vision models,
32 studying the performance on different input domains is an exciting future direction, as is understanding the relative
33 merits of convolution and attention beyond standard classification and detection tasks. Other study directions include
34 benchmarking performance of convolutional vs. attentional models on transfer and self-supervised settings.

35 **[R3]** *"grammatical errors and typos. Also somehow most of the references were missing in the paper."* We apologize
36 for accidentally clipping 4 pages of references section and any grammatical errors. We have addressed all of these
37 issues and will restore the references in the revised manuscript.

38 **[R4]** *"downsampling is carried out with average pooling with stride 2... instead of increasing the stride of the self-*
39 *attention layer"* We tried this downsampling approach in early experimentation and found it slightly underperforms
40 compared to average pooling. However, this experiment was conducted on a preliminary architecture, so we plan on
41 running experiments to benchmark this conceptually simpler approach on our final architecture.

42 **[R4]** *"not clear what self-attention can learn with respect to convolution, and what would happen with deeper models"*
43 We agree that a more rigorous study of the modeling and optimization capabilities of attention and convolution would
44 be illuminating. We leave this to future work. However, one clear difference is attention can generate a different kernel
45 per position based on content, while convolution uses the same kernel for every position.