

1 We thank the reviewers for their helpful comments. The code and models will be open-sourced.

2 **Reviewer 1**

3 **Inference speed test:** Along with peak RAM, we report inference FLOPs for all the models. We see a consistent
4 reduction in FLOPs per inference using RNNPool. Resource-constrained devices like microcontrollers, the key focus
5 of our work, do not have parallelism – so the reduction in FLOPs translates reasonably accurately to faster inference; we
6 will add inference numbers on Cortex M-class devices in the manuscript. On hardware with parallelism, the application
7 of RNNPool operator to multiple patches in an activation map can be fully parallelized. For face detection on 320×240
8 image, the RNNPool layer has about 4800 patches that can be executed in parallel – sufficient for many cloud hardware.

9 **Inappropriate comparison of experiments:** The main goal of Table 1 (that includes DenseNet121 on ImageNet-10)
10 is to show that RNNPool works well with a variety of architectures. While we agree that DenseNet121 result might
11 be an artifact of a big model applied to small ImageNet-10 dataset, we observe a similar trend for smaller models
12 like MobileNet also (that are in fact more critical for the resource-constrained settings that we study). We observe a
13 similar trend for ImageNet-1K as well (Table 3) but focus mainly on MobileNets due to resource-constrained devices
14 setting. Table 2 also shows that pooling/downsampling operators significantly affect accuracy in both DenseNet121 and
15 MobileNetV2. Finally, the rationale behind ImageNet-10 can be found in Appendix A.1.

16 **Reviewer 2**

17 **FLOPs of ReNet:** In Table 2 second block, we compare the accuracy of various downsampling methods while keeping
18 the working RAM same for each. Therefore, ReNet and RNNPool were placed in the same position and have the same
19 shape of input-output feature maps. So when the output map from RNNPool has $4h_2$ channels, the hidden state size of
20 the RNNs in ReNet was taken to be $2h_2$ to ensure an output map of $4h_2$ channels. And even then ReNet’s accuracy is
21 significantly worse than RNNPool on multiple tasks. Furthermore, in ReNet each input patch is flattened and passed as
22 a single input to an RNN which leads to a huge matrix multiplication operation. Finally, ReNet authors suggest using
23 GRU or LSTM as the RNN unit; we use GRU as it is more efficient. RNNPool uses FastGRNN as the RNN unit which
24 also contributes to the reduction in FLOPs.

25 **Similarity to ReNet:** We respectfully disagree with the reviewer. The global and local contexts do not depend on
26 overlapping or non-overlapping patches but are captured by the constituent operations in a layer. As mentioned in
27 L106-123, ReNet applies each horizontal and vertical RNN sweep on the *entire* image in *one pass* and the resulting
28 hidden states provide the output feature map. Even if ReNet uses overlapping patches, it will still apply both RNNs
29 in *one pass* over the image. In contrast, RNNPool is applied *patchwise* (it is never applied to the entire image in one
30 pass) and only the *final* states from each RNN run are used for output feature map. So, the intuition and goal of the two
31 methods seem quite different. We do present a comparison with ReNet (seen as a downsampling operator) in Table 2 on
32 two different tasks and with three different base models across 2 datasets.

33 **Comparison to DNN compression methods:** RNNPool is a complementary technique to model compression, i.e., we
34 can use model compression (e.g. pruning, quantization) on RNNPool based models as well to get additional reduction
35 in working RAM, model size & inference FLOPs. To illustrate this, we created the RNNPool-Face-Quant (Table 4)
36 model using a standard quantization based compression method that maintains accuracy with about $4\times$ compression.

37 **Reviewer 3**

38 **Quantization methods:** See response to R2; we will add references to the manuscript. As mentioned above, RNNPool
39 is complementary to and supports quantization (and other standard model compression techniques).

40 **Training from scratch:** RNNPool is a new building block (pooling operator with parameters) for CNNs. As is the
41 case with any architecture involving a new building block, training the model end-to-end provides the best performance.
42 XNOR-Nets are also trained from scratch because of the formulation and could be coupled with RNNPool based
43 models if needed during end-to-end training.

44 **Training of RNNPool & underlying difficulty:** The instances of RNNPool used in our models typically require a
45 sequence length of 8 for the RNNs, and thus are reasonably stable to train. We use the same hyperparameters as the base
46 CNN models to train the RNNPool variants. We do observe that the RNNPool based variants require more epochs to
47 train than the base CNN models. But, overall, the training complexity of RNNPool based models is similar to the base
48 CNN models. We will add training accuracy vs epochs plot in the appendix.

49 **Reviewer 4**

50 Thank you for the positive feedback. The pointer to SqueezeNAS and making a stronger case for RNNPool with the
51 argument is helpful. We will include them in the next revision and fix the citation to "Seeing AI".

52 **Deployment on a microcontroller:** We will open-source the code to deploy RNNPool based models on ARM Cortex
53 M-class microcontrollers, and add some of the details in the appendix.