

1 **Author response for “Incremental Few-Shot Learning with Attention Attractor Networks”**

2 We thank all reviewers for their time and insightful comments. We address individual comments below.

3 **To R1 on the normalization constant:** Thank you for pointing this out. In our experiments we actually used the  
4 same normalization constant (same softmaxes) for both support set and query set, both having  $W_a$ . Predictions for the  
5 base classes are simply ignored when calculating the loss for the support set. We will fix the paper to clarify this.

6 **To R1 on *mini-* vs. *tiered-ImageNet*:** *tiered-ImageNet* contains many more images, so this allows the network to be  
7 pretrained with a better feature representation, and we are also able to fit it with a larger and deeper network (ResNet-18  
8 vs. ResNet-10). The challenge in tiered ImageNet lies in the split on higher level class categories, so there is a domain  
9 shift in training, validation and testing.

10 **To R1 on computation overhead:** Thank you for asking. We are not directly inverting the matrix but computing the  
11 inverse matrix vector product. Second, we are using truncated RBP, which computes a low-rank approximation of the  
12 matrix inverse, i.e. doing Jacobian transpose vector product for a fixed number of steps (see Line 15 in Alg. 1). In the  
13 experiments, we used 20 steps, so this has the same time complexity as truncated BPTT for 20 steps, while RBP saves  
14 memory by not having to store intermediate activations. We will add these clarifications in the paper.

15 **To R1 on code release:** We will release the code that can reproduce the experiments if the paper gets accepted. The  
16 release is expected to happen before the conference.

17 **To R2 on catastrophic forgetting literature:** Thank you for providing the references. We will cite them and include  
18 them in our discussion of related work. [1] uses the diagonal approximation of the Fisher matrix to prevent the parameter  
19 from drifting too far. Since our backbone network is frozen, this wouldn’t be an issue in our setting, however, it could  
20 be complementary to our method if we choose to finetune the backbone network as well. [2] stores a few data points of  
21 old tasks in the “Coreset” so that they can estimate the variational distribution, whereas in our case, the regularizer is  
22 learned through meta-learning.

23 **To R2 on disjoint classes & data augmentation:** Data augmentation is only applied during pretraining, and for  
24 simplicity not applied during meta-learning. To ensure that in meta-learning novel classes do not overlap with base  
25 classes, we mask out the base classes that are used in the few-shot episode.

26 **To R2 on  $\theta_E$  in Figure 1:** Thanks for the suggestion. We will modify Figure 1 to match with the notations. The  
27 meta parameters  $\theta_E$  is enclosed in the teal part of the figure, which is learned during meta-learning stage, i.e. for mini  
28 ImageNet, there are 64 training classes, and the query set will be 59+5 classes, where there are 5 “fake” novel classes.

29 **To R2 on other notation issues:** Thanks for pointing them out. We will adopt your suggestions and change the  
30 notation to ensure consistency.

31 **To R2 on  $\Delta_a$  and  $\Delta_b$ :** These two measure how much drop in accuracy is caused by jointly classifying base and novel  
32 classes. We first record the accuracy by classifying base classes alone (i.e. standard classification problem), and novel  
33 classes alone (i.e. standard few-shot problem), and then observe the accuracy of base and novel classes in a mixed query  
34 set. Lastly, we take the difference of these two accuracies to be the drop in performance ( $\Delta_a$  and  $\Delta_b$ ).

35 **To R3 on motivation of attention mechanism:** The motivation is to assess the similarity between the novel classes  
36 and the base classes, so that during the learning of novel classes, the system learns to differentiate from the base classes  
37 and prevent interference. The attention vector is used to retrieve an attractor to regularize the episodic training.

38 **To R3 on clarity:** Could you comment on the specific places that you feel are not clearly written? We will revise  
39 these sections in the next version.

40 **To R3 on OptimizerStep:** Line 12-17 is the part of the algorithm that computes the gradients of the meta-parameters  
41  $\theta_E$ . “OptimizerStep” is an external function that performs optimization given some gradient direction (e.g. gradient  
42 descent, momentum, Adam, etc.).

43 **To R3 on T-BPTT:** The reason why T-BPTT fails is that it only optimizes for a short horizon (e.g. 20-100 steps of  
44 gradient descent). When we train the episodic objective till convergence, T-BPTT is no longer trained for that, so the  
45 accuracy goes down for longer iterations.