

1 We thank the reviewers for the time and expertise they have invested in these reviews. We appreciate positive comments  
2 about the paper like connecting the existing meta-learning frameworks with unsupervised/self-supervised feature  
3 learning frameworks, clear presentation of the ideas, good control parameter experiments, and novel approach to the  
4 unsupervised meta-learning problem. Furthermore, the insightful comments we received point us towards several  
5 directions in which this work can be extended in the future. We focus on answering to the posed questions by the  
6 reviewers and provide feedback related to improvements suggestions.

7 Answers to reviewer #1.

8 **How much does hand-crafted knowledge play a role in the performance of the proposed method?** Domain  
9 knowledge, such as knowing what transformations of the image retain the classification, does play a role in the  
10 performance of the augmentation function. See left side of Table 1 where various hand-crafted augmentations performed  
11 at accuracy of 30.16%, 32.80%, 35.09% against the 24.17% baseline. It turned out that the learned auto-augment  
12 outperformed all the hand-crafted methods with 39.93%.

13 **In particular, a simple baseline would be to perform few-shot training on models trained with the augmentation  
14 methods proposed. That seems like a more informative/reasonable baseline than training from scratch.** We  
15 agree, we implemented this, and here are the results. For 5-way, 1-shot learning on Omniglot the accuracy was: training  
16 from scratch 52.5%, training from scratch with augmentation 55.8%, UMTRA 83.8%. For MiniImagenet the numbers  
17 were: from scratch without augmentation 27.6%, from scratch with augmentation 28.8%, UMTRA 39.93%.

18 **The impact of this work could be more substantial if the proposed method were generalized beyond data-  
19 augmentation to unsupervised learning in the few-shot setting using other proposed self-supervised techniques.**  
20 While the main body of the paper focuses on augmentation techniques, its application to the video domain in the  
21 supplemental material is using self-supervision: we randomly select 16 frames from one video for training and another  
22 16 frames from the same video for validation. For results, see table at supplemental material page 4.

23 **The approach of using Auto-Augment within UMTRA would suggest significant computational overhead for a  
24 new dataset:** This overhead indeed exists during the meta-training time. Note however, that this is an offline process  
25 which needs to be done once per domain. No augmentation is done during the target learning phase.

26 **How would augmentation work for non-image data?** The choice of augmentation is domain specific. For example,  
27 in our supplemental material, for videos we do not do augmentation and just pick two different parts of the video. Look  
28 at the Figure 5 at supplemental material.

29 Answers to reviewer #3

30 **that all samples are in a different class (...) how a person might test this assumption on any particular problem  
31 where classes are not previously articulated (as they are in imagenet). (...) How would we estimate if " $N \ll C$ "?  
32 (...) Does one need to fully articulate the task space – how do you project all possible arrangements of the data  
33 into separate classes down into some effective number of classes,  $c$ ?** Indeed, for the unlabeled datasets that we use  
34 for meta-learning, we do not know the exact value of  $C$ . This is ok, because UMTRA does not take  $C$  as a parameter.  
35 We only need a rough estimate of it to ensure that " $N \ll C$ " holds, we do not need to fully articulate the task space. For  
36 instance, if we download 10,000 random videos from YouTube, we can estimate that there will be at least hundreds of  
37 different activities in them.

38 Answers to reviewer #6

39 **I would run experiments on one other domain to assure read-  
40 ers of the generality of the system across domains with widely  
41 varying appearances within the dataset.**

Algorithm (N, K)	(5, 1)	(5, 5)	(5, 10)
Training from scratch	26.86	39.65	50.61
UMTRA	<b>33.43</b>	<b>50.19</b>	<b>58.84</b>

42 As requested, we run UMTRA on the CelebA dataset (which is unbalanced  
43 in terms of the number of examples in each class) for identity recognition and obtained the results shown in table above.  
44 Given 5 new identities and one image of each, UMTRA is able to learn the task better than training from scratch. The  
45 supplemental material also contains results on another domain, video classification.

46 Answers to reviewer #7

47 **The most important further contribution would be demonstrating that the augmentation approach can either  
48 be effectively automated in all cases, or that a rough hand-designed augmentation approach can be found that  
49 works in other domains.** We believe that some degree of domain knowledge is a necessary input into the design  
50 of the augmentation algorithm. For instance, we know that photos are crops of 2D projections of 3D scenes, which  
51 imply certain invariances that can be exploited by the augmentation. This augmentation can be then used across all  
52 scenarios involving photos. For instance, the CelebA results above were obtained with the same augmentation used for  
53 MiniImagenet.