

1 We thank all reviewers for the thoughtful comments and constructive suggestions to improve our paper. In general, all
 2 reviewers find our general message: “*the model learned in each task is itself part of the inductive bias*” convincing.
 3 The core idea of incorporating model complexity into task embedding is “*well motivated and interesting*” (R4), “*novel*
 4 *and interesting, and generally applicable to many meta-learning models*” (R5), and “*is interesting and well-illustrated*”
 5 (R6). The only reservation shared by all reviewers is that experiments are not sufficient to support this claim.

6 Here, we address the major concern raised by all three reviewers — **generalization to more baselines**. We conducted
 7 additional experiments on two competitive baselines with large backbone feature extractors. To summarize, MATE
 8 brings consistent improvements by exploiting model information in task representations, which confirms our original
 9 finding. We plan to try more baselines and report in the final version. We also provide details about **meta-testing**
 10 **protocol** (R4), discuss **the gain brought by MATE** (R4, R5) and **the choice of FiLM layer conditioning** (R5).

11 ▷ **Applying MATE to more baselines**
 12 (R4, R5, R6). Per all your suggestions,

13 we conducted experiments on two more
 14 baselines. Due to the limited rebuttal time
 15 window and the well-known difficulty in
 16 finding suitable, reproducible implemen-
 17 tation of SOTA meta learning works, we
 18 turn to two baselines that have been com-
 19 pared in this paper, namely, Prototypical

Model	Backbone	5-way 1-shot	5-way 5-shot
MetaOptNet [20]	ResNet-12	72.00 ± 0.70%	84.20 ± 0.50%
MetaOptNet + MATE	ResNet-12	72.30 ± 0.70%	85.20 ± 0.40%
ProtoNets [43]	ResNet-12	71.35 ± 0.73%	84.07 ± 0.51%
ProtoNets + MATE	ResNet-12	71.49 ± 0.70%	84.71 ± 0.50%
R2D2 [6]	ResNet-12	72.51 ± 0.72%	84.60 ± 0.50%
R2D2 + MATE	ResNet-12	72.59 ± 0.70%	85.04 ± 0.50%

20 Networks [43] and R2D2 [6], but use larger convolutional backbones. We limit the experiments on CIFAR-FS, and
 21 will include miniImageNet results on the new baselines. The results are shown in the above table. Although ProtoNets
 22 and R2D2 are somehow old, we would still like to justify that comparing on these two are meaningful and can help to
 23 corroborate the generality of MATE framework. It is known that the original ProtoNets and R2D2 have much lower
 24 performance than more recent works, e.g. they are 12.2% and 4.8% lower in 5-way 5-shot accuracy on CIFAR-FS
 25 compared to MetaOptNet [20], respectively. However, once we try replace the backbone feature extractor with the
 26 same ResNet-12 used in MetaOptNet, ProtoNets and R2D2 both show competitive results, and especially R2D2
 27 already performs better than MetaOptNet just by ensuring a fair backbone. Then, MATE can still consistently provide
 28 improvements to both (enhanced) baselines: 1) applying MATE to ProtoNets+ResNet12 yields +0.64% 5-shot accuracy
 29 and slightly better 1-shot accuracy (+0.14%); 2) applying MATE to R2D2+ResNet12 yields +0.44% 5-shot accuracy
 30 improvement and similar 1-shot accuracy (+0.08%). These results are hence consistent with our original finding that
 31 MATE brings more benefits to 5-shot accuracy than to 1-shot, which is reasonable because we can obtain more accurate
 32 information about data distribution on the task with more data and thus task representation of higher quality.

33 ▷ **Protocol used for meta-testing** (R4). During the meta-testing stage, we sample 1,000 episodes (Section 3.2) from
 34 the meta-testing set following either 5-way 1-shot or 5-way 5-shot settings. The query set in each meta-testing episode
 35 contains 15 query images over which we calculate the meta-testing accuracy. We then report the average accuracy and
 36 standard deviation of the accuracies over the 1,000 meta-testing episodes. Due to large amount of testing episodes used,
 37 the standard deviation of the accuracies is sometimes very close. We confirm that the numbers reported in the tables in
 38 this paper are all correct. We would like to thank R4 for pointing out the ambiguity of the testing protocol. We will
 39 clarify this and add more details of the experiments in the final version.

40 ▷ **Limited performance gain** (R4, R5). Firstly, we thank R4 for the appreciation of the improvement brought by
 41 MATE, which we think can be further strengthened by the additional experiments we just conducted. Secondly, we
 42 humbly clarify that we calculate the accuracy standard deviation over 1,000 meta-testing tasks instead of the confidence
 43 interval. Hence, the accuracy improvement over 0.5% can show consistent improvement over a large sample of tasks.
 44 We’d also like to emphasize that incorporating model information into task embedding does help with and improve the
 45 performance, which is supported by the comparison of FiLM+KME and FiLM+SVM in Table 3 (2nd and 3rd rows).

46 ▷ **Conditioning FiLM layers on ω** (R5)? If we understand correctly, R5 suggests to condition FiLM layers on the
 47 optimal parameters learned by SVM (Section 3.1), instead of the model-aware task features proposed in this paper. We
 48 think this question can be answered well by humbly reminding R5 of the connection of our proposed method with kernel
 49 mean embedding (KME) [28], as we described before Section 3.1. In Eq. (1), if we ignore the model information by
 50 taking $f_M(x) \equiv 1$, Eq. (1) reduces to KME. Further, if ϕ corresponds to the canonical feature map of the characteristic
 51 kernel, the map defined by Eq. (1) is injective, i.e., the representation $\Phi(T)$ captures all information about the task T
 52 [10, Lemma 2]. Therefore, conditioning FiLM layers on the model-aware task feature defined in Eq. (1), which is very
 53 likely to contain most of information on the task, could possibly make the FiLM easier to train and, more importantly,
 54 more interpretable. We plan to conduct comparison experiments and report related results in final version.