1 Thank you all for the encouraging comments and helpful suggestions. Due to the space limit, we only respond to main
2 concerns, but will incorporate all comments in the final revision. Citations in this rebuttal refer to those in the main paper.

3 **Reviewer 1**
4 - **Bold text**: We bold the highest mean accuracy in all tables and show the bar of our model ($\sigma$-reptile) for readability
5 in Fig. 5. We will instead bold the statistically best models in the revision as you suggest.
6 - **When do we need long adaptation**: We will expand this in the revision. Generally speaking, long adaptation allows
7 higher model capacity for optimization-based meta learning [1-3]. The augmented Omniglot dataset was designed to
8 be challenging for short-horizon optimization, and MAML was found previously not to work for few-shot TTS [6].
9 - **Selecting modules according to $\sigma$ and setting to 0 in testing**: This is a very good question. Bayesian approaches
10 always allow some probability of selecting any module. It is common practice in the literature (e.g., Bayesian /
11 regularization-based feature selection [40] and model pruning) to choose features (modules in our case) with large
12 coefficients ($\sigma^2$). As explained in Sec D.2 we apply a weak regularizer on $\sigma^2$ so that the learned value for task
13 independent modules is close to zero. In practice, we often find that $\sigma^2$ of most modules is numerically indistinguishable
14 from 0, which reduces the gap between meta-training and meta-testing (where we set $\sigma = 0$). One exception is the
15 TTS experiment where some pruned modules have non-zero $\sigma^2$. All res-blocks have equal size in that case. There, we
16 choose the threshold with the best trade-off between pruned model size and validation performance. Your comment on
17 the relationship of $\sigma^2$ with the module size is however worth further investigation. We will add a discussion of this issue.
18 - **BatchNorm**: Each set of batch-norm parameters (offset, scale) has their own $\sigma$. Running stats are not parameters
19 but are computed from inputs. For consistency, all algorithms use transductive learning as in MAML and Reptile.
20 - **1-step MAML vs $\sigma$-Reptile**: MAML requires computing second-order gradients even with 1-inner step.

21 **Reviewer 2** We appreciate your assessment on the strengths of our work.

22 **Reviewer 3**
23 - **Choice of a Normal prior vs sparsity inducing prior**: The shrinkage prior is sparsity inducing for task parameters.
24 Given a weak prior on $\sigma^2$, the marginal prior distribution of task parameters is heavy-tailed with tied variance in a
25 group. Similarly, Laplace (spike-and-slab) priors is equivalent to a normal random variable with an exponential (mix
26 of point masses) prior on the variance. The estimate of $\sigma$ will shrink towards zero if the data does not support task
27 parameters deviating from the prior mean, and thus lead to sparse selection of parameter groups.
28 - **Showing benefit for some of the aspects that motivated this approach (interpretability, causality, transfer learn-
29 ing or domain adaptation)**: The cited line 30 is a general statement for the purpose of learning resuable modules in
30 the literature. We will clarify that we do not evaluate on all of these aspects. One of our main motivations is for saving
31 space in deployment by sharing the majority of parameters across tasks, which we demonstrate in all experiments.
32 Reviewer 1 also suggests that "*sharing parameters would also allow batching the evaluation of certain layers with
33 a GPU across tasks, which speeds up inference time*". Finally, the discovered middle layers in the WaveNet stack help
34 us better interpret the learned features in the different layers.
35 - **Bottleneck of hard-coded modules in few-shot domains**: This bottleneck may still exist in few-shot domains. Take
36 TTS for example, training a model from scratch typically requires at least 10K utterances. Learning a model with
37 $\sim 100$ utterances is a few-shot learning task. Nevertheless, the performance of recent works [5-7] with a hard-coded
38 task-specific component becomes saturated after around 10 utterances. Simply increasing the layer size of the task
39 specific component does not fix the problem.
40 - **Requiring large number of adaptation steps ($L$) in TTS**: We will clarify the reference to [6] of this requirement in
41 TTS. Generally speaking, allowing more adaptation steps increases the capacity of optimization-based methods [1-3].
42 - **Comparison to methods that learn learning rates**: We discuss the relationship of our method with those approaches
43 in Related Work and Appendix C. In Appendix C, we also compare with meta-SGD on a small synthetic example and
44 demonstrate the different behaviors for long adaptation. For real-data experiments, meta-SGD/-Curvature require back-
45 propagation through gradients and thus do not scale to long-adaptation regimes. Finally, SGD performs much worse than
46 Adam for the WaveNet model. It is not clear how to combine meta-SGD/-Curvature with Adam for a fair comparison.
47 - **Out-of-domain generalization in TTS**: We refer to the different size of the task dataset in meta-training and testing.
48 We will use a more accurate word in the revision.

49 **Reviewer 4**
50 - **Different learning rate in different approaches**: We run a random hyper-parameter search with 100 seeds for each
51 algorithm in each experiment and use the best found value from the validation task set. This is to ensure a fair comparison.
52 Particularly, we find that the reported hyper-parameters in the original papers (MAML, iMAML, Reptile) are not
53 optimal in our implementation. This small difference in learning rate makes essentially no difference in performance.
54 - **Small scale networks** The WaveNet model in the TTS experiment is a large network with 30 residual blocks. Every
55 block is a dilated and gated causal-convolutional net (see architecture in Figure 12 of Appendix E.2). There are 3M
56 parameters in total. In contrast to most other few-shot meta-learning works that use a pretrained residual network as
57 backbone and meta-learn only the last few layers, we meta-learn the entire CNN stack.