

1 We would like to thank all the reviewers for careful reading our paper, and we sincerely appreciate the constructive
2 comments. We will individually address the questions raised by each reviewer.

3 **Reviewer 1 Transferability to diverse datasets.** We conducted
4 additional experiments on the SVHN dataset shown in the Table to
5 further demonstrate our algorithm. We used the same Meta network
6 prior previously trained on the multi-scale Imagenet datasets and
7 quickly adapted the meta network to SVHN for two epochs in less
8 than an hour. The networks were trained for 160 epochs with cutout
9 augmentation. We also trained the CIFAR10 specialized network

Architecture	Top-1 Err	Params(M)
WideResnet [2]	1.30 ± 0.03	11.7
MetaQNN [3]	2.24	9.81
DARTS (CIFAR10 Searched)	2.09	3.3
BASE (Multitask Prior)	2.13	3.22
BASE (Imagenet32 Tuned)	2.07	3.29
BASE (SVHN Tuned)	2.01	3.22

10 from DARTS. The adapted network achieves the best performance in our experiments and has comparable performance
11 to other work for the model size. It is difficult to compare to [1] since they don't explicitly state their final performance
12 on the SVHN dataset outside of small graphs and search smaller networks. Roughly, based on the figures, it appears to
13 achieve about 97% accuracy in a 2 GPU hours and around 97.3% accuracy in 20 GPU hours.

14 **Reviewer 2 Novelty of the Bayesian View.** Our Bayesian view of architecture search is fundamentally different
15 from SNAS [4] since our approach has the ability to naturally extend to meta architecture search while SNAS cannot.
16 Specifically, we can transfer prior information among tasks which is critical for meta architecture search. SNAS's
17 formulation which is motivated from the policy gradient in reinforcement learning, has no prior structure and can only
18 be applied to a single task. Compared with SNAS, the networks found by BASE achieved comparable performance on
19 CIFAR10 and higher accuracy on Imagenet with meta architecture search.

20 **Optimization Embedding.** Optimization embedding is used for parameterization of the posterior distribution of the
21 convolutional weights and architecture parameters in Eq. (11), ϕ_D^t and ψ_D^t . In the few-shot learning experiments, we
22 fully unroll the optimization embedding while in the other experiments, we use the finite difference approximation
23 corresponding to line 10 of algorithm 1. Compared to the vanilla independent posterior parametrization with an arbitrary
24 neural network, optimization embedding allows us to model the dependencies between D and W explicitly which
25 allows to better adaptation to new datasets.

26 **Training Time.** As with most work in the field of architecture search, our method focuses on architecture search over
27 model training. However, with our low search times, we do recognize the increasing importance of training time, and
28 we will add a discussion to our final paper. Our Bayesian formulation does naturally allow transfer of weights from the
29 meta-network if we don't transfer to larger networks which could potentially eliminate the additional training time.

30 **Transferability of Meta-Weights.** We did not transfer the posterior of weights to the final networks in our experiments.
31 Figure 4, however, does demonstrate the fast adaptation of both weights and architecture where our meta-network
32 achieves 75% accuracy on the unseen CIFAR10 in one epoch based on the posterior of both weights and architecture.

33 **Clarification about few-shot learning results.** Indeed, MAML and REPTILE use the same model architectures.
34 However, they are fundamentally different algorithms for Meta-learning: MAML uses the full gradient and REPTILE
35 uses the finite difference approximation, leading to different performances.

36 To be clear, DARTS, BASE (Softmax), and BASE (Gumbel) are the algorithms used to generate network architectures.
37 They are *not* few-shot classification algorithms and in Table. 3, the actual few-shot learning with which the models are
38 trained is conducted using MAML. Therefore, it only makes sense to directly compare against the MAML baseline.
39 The REPTILE baseline is included only for completeness. It is true that the models used in MAML were much smaller,
40 and we will update the paper. However, we chose to use our larger search space so we could directly compare to the
41 the CIFAR10 optimized network architecture from DARTS. This network, though much larger, significantly beat the
42 original MAML paper results with the same algorithm. BASE was able to find networks in the same search space which
43 achieved even better results. We will make this point clearer in our final paper.

44 **Error Rates.** It is definitely possible that the cell space and scaling may complicate the results, but this is still a
45 common methodology and shouldn't invalidate our work. I believe the particular part you quoted from our paper was in
46 reference to our found Imagenet-optimized architecture which achieved fairly high performance, but we will revise
47 those sections to clarify our claims.

48 **Clarity of Figures.** We will be expanding the related work, clarifying the notation, and fixing those typos. To address
49 the issue in Figure 1, the M in that figure should be L following the notation in the rest of the paper which corresponds
50 to the number of choices of operations for each layer type.

51 **Reviewer 3** Thank you for your appreciation and suggestions for paper refinement. We will polish the the final
52 version of our paper following your comments.

53 [1] Wistuba et al. Inductive Transfer for Neural Architecture Optimization. *arXiv:1903.03536*, 2019.

54 [2] DeVries et al. Improved regularization of convolutional neural networks with cutout. *arXiv:1708.04552*, 2017.

55 [3] Baker et al. Designing neural network architectures using reinforcement learning. *ICLR* 2017.

56 [4] Xie et al. SNAS: stochastic neural architecture search. *ICLR* 2019.