

1 **Response Summary.** We truly appreciate reviewers’ valuable comments and positive feedback. We are encouraged
2 reviewers found the problems addressed by our paper are important (Reviewer 3, 4), our motivation is clear (Reviewer
3 1, 2, 3), and idea is novel (Reviewer 1, 3, 4). We are glad all reviewers found this paper is well written with excellent
4 reproducibility. One primary concern was the results on DARTS search space were not SOTA. Similar as DARTS and
5 other differentiable NAS, we also conducted more architecture search for high-performance models. We have released
6 a new architecture on GitHub (the anonymous link in the original submission file), with 2.50% test error on CIFAR10,
7 15.84% on CIFAR100, and 24.37% on ImageNet, which are SOTA results under DARTS experimental settings. Thanks
8 Reviewer 3 for the strong support, and we will provide a detailed item-wise response for other reviewers below.

9 **Responses to Review 1:**

10 **Q1:** As for exploration, there are some other traditional methods like BO. What the difference between BO and your
11 way is. **R1:** BO needs to build probabilistic models (e.g. GP) by training multiple architectures from scratch first,
12 while differentiable NAS only trains a supernet once. It is not intuitive to directly introduce the exploration in BO to
13 differentiable NAS, while our exploration could be easily applied to differentiable NAS to solve rich-get-richer problem.

14 **Q2:** Is it difficult to tune hyperparameters? **R2:** No. There are only two hyperparameters to be tuned, and the others are
15 default. The ablation studies in Sec 4.3, 4.4, and Appendix H show our model is robust to the two hyperparameters.

16 **Q3:** Some SOTA one-shot NAS methods [1,2] in the NAS-Bench-201 dataset missed. And the results on DARTS
17 search space are not the SOTA. **R3:** The two references are missing in the review. Our best single run can achieve
18 46.48% for ImageNet on NAS-Bench-201 as described in the bottom of Table 1 and line 246 in the original submission,
19 outperforming [1] (46.34% for ImageNet) provided by the reviewer. Please refer to **Resp. Summ.** and GitHub for the
20 results of our new high-performance model on the DARTS search space.

21 **Responses to Review 2:**

22 **Q1:** The paper is not very novel with limited contribution, the method mitigating catastrophic forgetting is not new and
23 similar with the previous EWC work. **R1:** This paper is the first paper introducing intelligent exploration into NAS,
24 through the probability density function based on a graph autoencoder. We believe it is novel in NAS and Reviewer 1
25 and 3 also agree on it. Besides, the regularization method to mitigate forgetting in our method is totally different from
26 EWC and WPL (EWC applied to one-shot NAS). EWC and WPL both calculate the joint posterior probability through
27 estimating the Fisher information matrix and assuming the previous models in optimal points, while the two conditions
28 hardly hold in differentiable NAS. We propose an architecture complementation scheme, and theoretically shows it
29 could optimize the joint posterior probability as EWC and WPL, without the assumption of the two conditions.

30 **Q2:** NAS-Bench is not a well-established benchmark that not many people are very familiar with. **R2:** The NAS-
31 Bench is a newly established benchmark with a much simpler search space, while the ground-truth test accuracy for all
32 candidates in the search space is reported, helping the NAS methods to conduct reproducible experiments with much less
33 computational requirements. Building a well-established benchmark is becoming a new interesting research direction in
34 NAS, and concurrent NAS-Bench 101, NAS-Bench 201, NAS-Bench 1Shot1, et. al., all help to relieve computational
35 requirements, and recent researches in NAS community prefer these benchmarks for enhancing reproducibility.

36 **Q3:** Require a stronger ImageNet result. **R3:** Please refer to **Resp. Summ.** and GitHub for our new SOTA results.

37 **Q4:** The author should not treat supernet training as a multi-task learning problem. The reviewer agrees that it is a
38 multi-model optimization problem, not a multi-task problem. **R4:** This paper focuses on the multi-model forgetting
39 problem in the supernet training. The multi-model forgetting in NAS is very related to catastrophic forgetting in multi-
40 task learning, as described in Sec. 2. To avoid confusion, we have rephrased "catastrophic forgetting" to "multi-model
41 forgetting" when describing the forgetting in NAS, and rephrased Sec.3.2.

42 **Q5:** In relieving forgetting, the selection of three models seems arbitrary and not intuitive. **R5:** It should be noted that
43 our architecture complementation scheme is to select specific architectures for regularization. We theoretically show our
44 method can optimize the joint posterior probability similar as EWC and WPL, but with less constraints. The ablation
45 study in Appendix H shows our method outperforms other naive schemes, including WPL, random selection, and so on.

46 **Q6:** The paper directly compares related works copied from the previous paper without hyperparameter tuning. **R6:**
47 The results of peer algorithms are from the original paper (NAS-Bench-201) since we adopt the same experimental
48 settings as that paper. Furthermore, our model outperforms most peer algorithms under all hyperparameter settings.

49 **Responses to Review 4:**

50 **Q1:** In Table 2, $\gamma = 0.2$ performs the worst with huge variance? Should the performance change smoothly with γ ?

51 **R1:** The performance should change smoothly with γ . As discussed in Sec. 4.3, a dynamic γ seems to achieve better
52 performance, and a small and static γ may lead to local optimal. As we conducted experiments with limited random
53 seeds, the outliers may greatly affect the statistical results. We have conducted experiments with more random seeds
54 and will remove outliers to obtain statistical results to avoid the effects of outliers in the final version.

55 **Q2:** The definition of complementary/orthogonal architecture? Should the union of α_{i-1} and α_i^c be the whole search
56 space, or just needs the union of α_{i-1} and α_i^c includes α_i ? **R2:** We define that α_m is orthogonal to α_n , so they do not
57 share parameters $\omega_m \cap \omega_n = \emptyset$. As to the complementary architecture, since we first select the α_{i-1} into the replay
58 buffer, the complementary architecture α_i^c is defined as $\omega_i \subseteq \{\omega_i^c \cup \omega_{i-1}\}$ that only needs the union of α_{i-1} and α_i^c
59 includes α_i , and α_i^c is also orthogonal to α_{i-1} .