

1 We thank all the reviewers for their constructive comments. Below are detailed responses.

2 **R1&R3: Co-design process elaboration.** We provide a simple pseudo-  
 3 code in Alg. A due to space limit. We will provide details in the final draft.

**Algorithm A.** The co-design process.

```

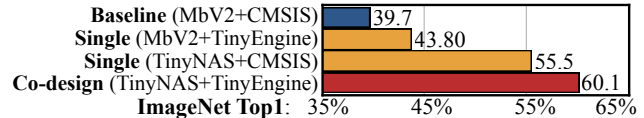
4 # TinyNAS: sample a DNN arch
5 for arch in arch_space:
6 # TinyEngine: find a good schedule
7 for schedule in schedule_space:
8 # check if satisfy mem. constraints
9 if can_fit_memory(arch, schedule):
10 # eval acc. and update best arch
11 acc = get_valid_acc(arch)
12 best_acc = max(best_acc, acc)
13 break

```

4 **R1: More deployment devices and tasks.** MCUNet generalizes well  
 5 across different MCU devices with different capacities: we show the ImageNet  
 6 top-1 accuracy on F746 (320kB SRAM, 1MB Flash) and H743 (512kB  
 7 SRAM, 2MB Flash) in Table A, MCUNet consistently outperforms the base-  
 8 line by a large margin (up to 20.4%). MCUNet also generalizes beyond  
 9 classification to detection. On PASCAL VOC with YOLO, MCUNet signifi-  
 10 cantly improves the mAP from 31.6% to 51.4% on H743. To the best of our  
 11 knowledge, this is the first large-scale object detection experiment on tiny MCU devices.

	ImgNet(F746)	ImgNet(H743)	VOC(H743)
MbV2+CMSIS	39.7%	53.8%	31.6%
MCUNet	<b>60.1%</b>	<b>65.1%</b>	<b>51.4%</b>

**Table A.** MCUNet shows consistent improvement across different devices (F746, H743) and tasks (classification, detection).



**Figure A.** MCUNet’s co-design scheme outperforms single-design ones on ImageNet classification.

12 **R1: Improvements from co-design over single-design.**

13 We showed the advantage of the co-design scheme in Table  
 14 2 of the original paper, where co-design achieves 4.6% higher accuracy compared to the best single-design result. We highlight the advantage of the co-design scheme in Figure A. We will make it more clear in the final draft.

15 **R1: Whether the overall network topology brings major improvement.** Yes, considering the overall network  
 16 topology enables specialized im2col, specialized loop tiling and unrolling strategies, which accounts for 49% of the  
 17 overall performance boost achieved by TinyEngine.

18 **R2: Why the auto-tuning in TVM fails to work on MCUs.** MicroTVM’s auto-tuning is based on a pre-defined  
 19 implementation template. However, the template does not include our advanced optimizations, e.g., scheduling memory  
 20 according to the overall network topology. Therefore, auto-tuning cannot match our speedup and memory reduction.

21 **R4: Contributions of TinyNAS.** We would like to clarify that TinyNAS is novel for the “actual NAS procedure”.  
 22 TinyML on MCU is a very new area; existing NAS methods *cannot* fit the tight memory constraints. TinyNAS is the  
 23 *first* NAS algorithm to enable large-scale deep learning on MCU devices. Since there is no carefully-tweaked design  
 24 space like those for mobile phones, we have to start from a huge search space so that it is likely to contain a good  
 25 model for various MCUs. The space needs to cover not only the micro-level architecture designs (e.g., kernel size,  
 26 expansion ratio) but also the macro-level designs like input resolution and channel widths (Section 3.1). Existing NAS  
 27 methods fail to achieve good performance on the huge space (Table 5 in original paper), since the large space makes  
 28 weight-sharing difficult and leads to *low* sample efficiency due to the *sparse* search reward. Our TinyNAS overcomes  
 29 the search inefficiency with a two-stage search algorithm. The first stage is to shrink/prune the huge search space to a  
 30 smaller sub-space, so the reward is no longer sparse, and the sample efficiency is improved. The second stage is to  
 31 perform micro-level optimization in the pruned sub-space. Both stages are the “actual NAS procedure”; they work  
 32 jointly with TinyEngine to achieve a decent performance, and should not be considered separately.

33 **R4: Comparison to budget-aware NAS methods.** TinyNAS argues that a two-stage  
 34 algorithm that gradually narrows down the search space is important to avoid the sparse  
 35 search award. Therefore, a fair comparison needs to start from the same full space. We  
 36 modify existing NAS methods to use the same search space under the same memory  
 37 constraint as ours. Compared to Single Path One-Shot NAS (SPOS) [17] and Once-For-  
 38 All (OFA) [5] on ImageNet-100 (ImgN<sub>100</sub>), TinyNAS outperforms both SOTA methods  
 39 (Table B), which verifies the advantage of our two-stage search mechanism. Other NAS  
 40 methods (e.g., [6, 44]) cannot handle the macro-level architecture like backbone channel  
 41 widths like ours. Therefore, we scale their channels&resolutions to fit the same memory

**Table B.** Compare NAS.

Method	ImgN <sub>100</sub>	ImgN <sub>1k</sub>
MnasNet [14]	-	51.8%
FBNet [44]	-	50.6%
Proxyless [6]	-	54.4%
SPOS [17]	75.6%	53.6%
OFA [5]	77.0%	54.0%
TinyNAS	<b>78.7%</b>	<b>60.1%</b>

42 budget of STM32F746 (320kB). Under the same MobileNet-v2 search space, TinyNAS shows significant advantage on  
 43 ImageNet (ImgN<sub>1k</sub>) with up to **9.5%** better top-1 accuracy, which verifies memory-awareness is important for TinyML.

44 **R4: Existing NAS methods that optimize memory footprint.** The two papers provided by the reviewer do not  
 45 optimize the working memory footprint. MorphNet [Gordon et al., 2018] only considers FLOPs and model size as  
 46 constraints. Though [Veniat et al., 2018] mentions “memory consumption cost”, it actually refers to model size but not  
 47 activation memory, which is the bottleneck. Neither explored memory-bounded NAS at tiny MCU scale (<1MB).

48 **R4: Details about the experimental protocol.** Many experimental protocol details are provided in Section 4.1 and  
 49 Section G of the supplementary (e.g., datasets, momentum, weight decay, training epochs). We will add more details to  
 50 the main paper in the final version to help reproduction.

51 **R4: Limited space for NAS.** Both stages are the actual NAS procedure to search a good model from a huge search  
 52 space. Therefore, we have dedicated a considerable amount of space for the NAS procedure. Due to the space limit, we  
 53 put some of the details of the second stage in the supplementary. We will add it to the main paper in the final version.