

1 **All reviewers:** *Improvement over SOTA:* Our work 1) improves results by the union of narrow-space search, 2) accounts
2 for constraints, and 3) improves accuracy for the given constraints. The reviewers agree with the novelty of point
3 1. Even though reference work solves points 2 and 3 and we reuse a variant of a genetic algorithm inside the inner
4 optimization loop of our automatic workflow (L163), we claim the uniqueness of our overall approach and results.
5 More precisely, we provide a low effort solution to a NAS problem by decomposing it into three steps that can be
6 solved 1) analytically, 2) by calibration measurement, and 3) by training models. We stress run times inequalities; 1
7 analytical evaluation (100ths per sec) \ll 1 calibration measurement (1sec – 100sec) \ll 1 model training (>1h even
8 on P100 GPUs). Solving the NAS for one specific constraint, the left parts depicted in Fig. 3 can be evaluated efficiently
9 (especially, the genetic algorithm does not involve model training). We assess 10'000 of models before training, where
10 only a small subset of those is required to be trained.

11 *Novelty:* First, we extend the search range to small models (we cover five orders of magnitudes in the number of
12 parameters) whereas reference work reports around a reference point or on a limited dynamic range. Second, we
13 evaluate quantized models on the front of optimal models. Even though quantization and reduced precision models
14 are established, key reference results are obtained on selected operation points. In contrast, our evaluation reveals the
15 effects of trade-offs between quantization and NAS results for smaller networks.

16 *Significance:* In our opinion, the significance and innovation should not be solely judged on the subparts of contributions
17 on their own but rather by the overall impact of the approach. Our method is simple and clean. It has been designed for
18 practical industrial applications, targeting industrial experts with limited AI background and/or time to invest in AI
19 expertise. We provide a flexible and economical design strategy that helps to reduce development time to design AI
20 models for edge devices and impact industrial solutions.

21 **Reviewer 1:** The manual and automatic workflow are independent approaches to solve the same NAS problem and they
22 do not interact. Fig. 4 (right), is the analog plot of Fig. 1 (right), where the “random configuration” shows the statistics
23 of networks that are sampled in the full space, the “automatic search” are the results of running and concatenating the
24 genetic algorithm as described in L181+. We will better clarify this in the revised version.

25 The closest related work is MnasNet [46] that follows a traditional NAS approach, where the feedback loop includes
26 the expensive training step reporting a run time of 4.5 days on 64 TPUv2 devices. That is equivalent to about 640
27 GPUh/TPUv2 pod*24h/d*4.5d = 69'120 GPU hours. In comparison, our approach (100 selected models) leads to about
28 200 GPU hours (similar model size as ref.). Accounting for the explainable difference of the dataset 24x (1.2M/50k
29 samples) but also for the training 1/20x (they: 5 epochs, we: 100 epochs) our NAS approach remains about two orders
30 of magnitude more resource-efficient! 2) Results: all MnasNet models exceed 3.9M parameters and their latency
31 domain spans one order of magnitude between 20ms and 140ms (MnasNet, Fig. 5). In contrast, we include the 1k up to
32 1M parameters range that is essential to reduce the per IoT device cost and we cover two orders of magnitude on a
33 much weaker device in terms of latency time.

34 **Reviewer 2:** We agree, manual and automatic comparisons are challenging to compare due to the expert knowledge. In
35 our results, we assume a best-effort approach with an experienced data scientist. All decisions taken by the expert are
36 listed in the appendix. We stress, that for a fair comparison, the narrow-search space is common for both approaches.
37 Knowing Fig. 6 implies you have already spent the experimenting runtime, so if your intention is to automate and
38 shorten the overall development time you aim to not depend on the “human in the loop” at that stage. If you still
39 do, most probably you want to extend the search space to account for new design patterns, if you do so, for a fair
40 comparison, you should rerun both workflows on the new space. If your initial assumption was right, and you indeed
41 have a stronger search space, we expect to improve results in both cases.

42 **Reviewer 3:** The number of sampled narrow spaces that are needed to deliver a good solution depends directly on the
43 quality of the search space. If it would be a priori clear, or there is strong evidence that one space is considerably better
44 than others, it would be enough to perform all considerations on that space only. However, since we do not have such
45 evidence prior to run the experiments, we propose the aggregation of multiple subspaces. Since each definition involves
46 paper and pen definitions and related implementations the number of considered narrow spaces must be reasonably
47 small. We conducted all experiments with five narrow-spaces.

48 We do not provide theoretical guarantees. In contrast, we adopt the in-the-field practice of empirically reporting results.
49 Our opinion is that the main contribution of our work is formulating the NAS problem in a decomposed way that allows
50 constraint evaluations and HW measurements to rule out candidates before they are required to be trained.

51 We conducted the key results on CIFAR-10 since we know that this is a well-established dataset with many reference
52 results, especially experiments conducted in IoT settings on non-standard hardware. We hope that this choice allows
53 a common ground between HW developers, the deep learning community, and the IoT industry to better judge our
54 achievements. Fig. 8 of our paper includes results where we applied our search strategy to additional thirteen image
55 classification tasks, as detailed in the appendix.