

1 Thank you for the insightful comments and the opportunity to follow up.

2 **R1, R2, R3: Comparing Nimble with TensorRT, TVM, TensorFlow(XLA).** TensorRT and TVM employ graph  
3 optimizations (e.g., aggressive operator fusion) and kernel selection/tuning, which are orthogonal to our idea. We  
4 applied operator fusion (less aggressive than TensorRT’s) and basic kernel selection for Conv op (use either cuDNN or  
5 PyTorch’s native implementation) to Nimble and measure its performance. Figure 1 shows that Nimble outperforms all  
6 cases except MobileNet V2 (TVM). The reason is that TVM spends more than a day in tuning Conv kernels of a model,  
7 and such tuning happens to be remarkably effective on MobileNet V2, finding more efficient kernels compared to those  
8 of cuDNN and PyTorch. Note that TensorRT and TVM do not support training for now. We will add these results in our  
9 revised paper, along with detailed discussions on related works including TorchScript and TFRT.

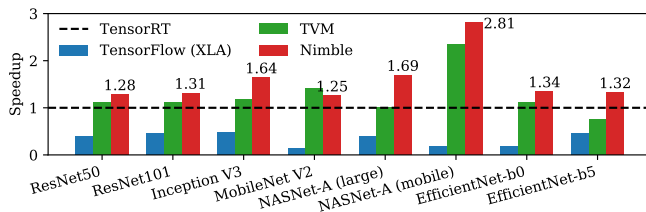


Figure 1: Speedup compared to TensorRT on inference workloads (batch size 1) using V100.

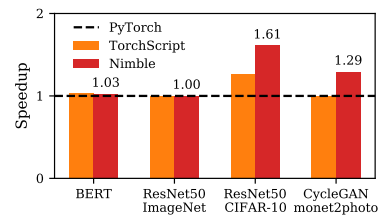


Figure 2: Speedup compared to PyTorch on training using V100.

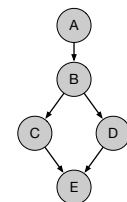


Figure 3: An example DAG.

10 **R2: Clarifying our contributions.** Nimble is the first work to automatically avoid framework overheads and aggres-  
11 sively parallelize GPU kernels using multiple streams for static DL models. Nimble introduces ahead-of-time (AoT)  
12 preparation to avoid framework overheads (details discussed below). This AoT preparation can be done quickly, and  
13 Nimble does not experience the overheads when executing the DL model afterwards. Furthermore, this opens up an  
14 opportunity to more efficiently utilize multiple GPU streams (as discussed in the first paragraph of §3.2). Nimble  
15 proposes a new multi-stream algorithm that maximizes parallelism and minimizes the number of synchronizations  
16 across streams. Nimble automates applying these techniques by capturing the GPU kernel call trace and running only  
17 GPU kernel calls for each execution, without redesigning the framework runtime. Nimble also uses CUDAGraph to  
18 reduce the number of GPU kernel launches. As a result, Nimble exhibits significant inference (and training) speedup on  
19 various models. Moreover, Nimble is easy to use; a user just needs to wrap a PyTorch model in a Nimble object (two  
20 additional lines), and use the Nimble object for inference or training.

21 Nimble’s approach is different from prior approaches that identify framework overheads and remove such overheads.  
22 Instead, Nimble captures the core DL computations (i.e., GPU kernel call trace) to run and prepares an environment for  
23 executing the captured trace for a new input.

24 **R2: Sources of framework overhead.** As discussed in the paper, the framework overhead is incurred not only by  
25 well-known sources like memory allocation but also by other sources such as inferring the output shape, dispatching  
26 appropriate GPU kernels, and preparing GPU kernel arguments. Existing approaches (e.g., memory preallocation)  
27 are limited to specific sources of overhead. Redesigning the framework to remove all sources is very challenging. As  
28 described above, we present a solution to avoid all overhead sources without rewriting the framework.

29 **R1, R2: Large training workloads.** Figure 2 shows Nimble’s performance when training larger models. We use batch  
30 sizes of 32, 64, and 1 for BERT, ResNet50, and CycleGAN respectively. As shown in the result of BERT and ResNet50  
31 (ImageNet), framework overhead is less pronounced when a model mostly consists of large kernels (kernels with large  
32 amounts of computation), leading to limited performance improvement in Nimble. We will include these results to  
33 show the limitations of Nimble in our revised paper. Nonetheless, there exist other important cases where the model  
34 contains small kernels. For example, training classification models on the CIFAR dataset or training typical GAN  
35 models generally involve small kernels, hence Nimble achieves training speedup on such models.

36 **R2: Generality of Nimble.** Nimble supports static DL models, and is not applicable to dynamic models. Yet, we  
37 believe that Nimble covers a wide range of models and has practical, real-world impacts; for instance, TensorRT is  
38 widely deployed in production albeit its limited applicability.

39 **R3: Stream capture / Initialization cost / Comments on §3.2 and Proof.** We capture all operations of the model at  
40 once. The mean and maximum AoT preparation time for the models in Figure 1 are 0.35 s and 1.07 s (NASNet-A  
41 (large)), respectively. We will describe the algorithm and proof more clearly in our revised paper.

42 **R3: Simplifying the stream assignment algorithm.** To our understanding, we cannot omit the process of constructing  
43 a bipartite graph. We describe an example in Figure 3. Since every path from A to E includes the edge (A, B), the  
44 maximum flow of graph is trivially 1, and does not give useful information for the stream assignment of the graph. We  
45 greatly appreciate the suggestions.