

1 A Experiment Details

2 A.1 Hyperparameters

3 In this section, we list out all the selected hyperparameters in our experiments for reproducibility in
4 Table 1 and Table 2. For selecting the number of layers in the GS and the segmented transformer
5 networks. We did a hyperparameter sweeping of [2, 4, 6, 8] network layers for both GS and the
6 segmented transformer network and find the optimal hyperparameters as presented in below table.

Table 1: Hyperparameters for Policy Network. *gs_layers*: GraphSAGE layers, *gs_knn*: GraphSAGE maximum neighbors, *trf_d_model*: Dimension of the segmented transformer model, *trf_n_head*: Number of attention heads, *trf_layers*: Number of transformer layers, *trf_d_heads*: Dimension of each attention head, *trf_d_inner*: Dimension of inner hidden size in positionwise feed-forward.

Parameters	Value	Parameters	Value
<i>gs_layers</i>	4	<i>gs_dim</i>	128
<i>gs_knn</i>	5	<i>trf_layers</i>	4
<i>trf_d_model</i>	128	<i>trf_n_head</i>	3
<i>trf_d_head</i>	15	<i>trf_d_inner</i>	512

Table 2: Hyperparameters for PPO.

Parameters	Value	Parameters	Value
<i>learning rate</i>	0.5	<i>num of rollouts</i>	800
<i>minibatches</i>	40	<i>epochs</i>	20
<i>epsilon</i>	0.2	<i>entropy</i>	0.5
<i>optimizer</i>	Adam		

7 A.2 Input Graphs

8 We used important workloads that are widely used in real applications or are incorporated in industry
9 standard benchmarks like MLPerf. These include ResNet, InceptionNet, WaveNet, Transformer-XL,
10 AmoebaNet, NMT, and RNNLM. In this section, we give a detailed explanation on the selected
11 models and hyperparameters.

12 A.2.1 Inception-V3

13 Inception-V3 [6] is a multi-branch convolutional network used for a variety of computer vision
14 tasks, including classification, recognition, or generation. The network consists of blocks made of
15 multiple branches of convolutional and pooling operations. Within a block, the branches of ops can
16 be executed in parallel. However, the model is mostly sequential as the outputs of each block are
17 concatenated together to form the input to the next block. We use a batch size of 64. The TensorFlow
18 graph of this model contains 24,713 operations.

19 A.2.2 AmoebaNet

20 AmoebaNet [5] is an automatically designed neural network that yields SOTA performance on
21 ImageNet. Similar to Inception-V3, it contains Inception-like blocks called cells, which receives a
22 direct input from the previous cell and a skip input from the cell before it. The network is made of
23 redundant cells stacked together, therefore is more modular than Inception-V3. We use a batch size
24 of 64. The TensorFlow graphs contains 9,430 operations.

25 A.2.3 RNNLM

26 Recurrent Neural Network Language Model [10, 4] is made of many LSTM cells organized in a
27 grid structure. The processing of each LSTM cell only depends on the results of 2 other cells (from
28 the previous layer, and from the previous time step), which make the concurrent execution of many
29 LSTM cells possible given enough hardware resources. We use batch size 64 and a hidden size of

30 2048. The corresponding TensorFlow graph contains 9,021 operations for a 2-layer model. The
31 number of ops grow roughly proportional with the number of layers.

32 **A.2.4 GNMT**

33 Neural Machine Translation with attention mechanism [1, 9] has an architecture similar to that of
34 RNNLM, but its many hidden states make it far more computationally expensive than RNNLM. To
35 reduce the training time, prior work [9] propose placing each LSTM layer, as well as the attention and
36 the softmax layer, on a separate device. This strategy demonstrates early success in human placement,
37 we show that GO can find significantly better placements. We use batch size 64. The original 2-layer
38 encoder-decoder consisting of 28,044 operations. An extended 4-layer version consisting of 46,600
39 operations, An even larger 8-layer version consisting of 83,712 operations.

40 **A.2.5 Transformer-XL**

41 Transformer-XL [2] is an modified version of Transformer [8] that supports segment-level recurrence
42 and a novel positional encoding scheme. This innovation enables learning dependency that is 80%
43 longer than RNNs, and 450% longer than vanilla Transformers. We use a transformer-XL with batch
44 size of 64, sequence length of 256, segment length of 64, model hidden dimension of 500 and feed
45 forward hidden dimension of 1000, 10 heads, and head dimension of 50. The 2-layer Transformer-XL
46 contains 2,618 operations. The number of ops grow roughly proportional with the number of layers.

47 **A.2.6 WaveNet**

48 WaveNet [7] is a generative model for speech synthesis. The model is fully probabilistic and
49 autoregressive, with the predictive distribution for each audio sample conditioned on all previous ones.
50 Architecturally, WaveNet uses causal convolutions with dilations, to obtain a large receptive field.
51 We use a WaveNet model with batch size 64 and a receptive field size of 2048 (9-layers per stack).
52 An 5-stack WaveNet contains 4,374 operations and a 10-stack WaveNet contains 8,516 operations.

53 **A.3 Performance Model**

54 An analytical performance model [3] based on the roofline model is used in this work. Specifically, the
55 execution time of a kernel is estimated using analytically modeled flops and bytes of memory accessed,
56 together with GPU's peak achievable FLOPS and memory bandwidth. The kernel launching overhead
57 is not considered. For data transfer between devices, the transfer time is estimated using the tensor
58 size and bandwidth. A virtual scheduler is developed to handle the dependency among different ops
59 in the graph within and across devices with several available scheduling policies including the priority
60 based one used in this work. The accuracy of the performance model has been validated against true
61 runtime measurements (on actual hardware) on several industry standard models, including MLP,
62 CNN, RNN, LSTM, Transformer, BERT, etc.

63 **B Ops Scheduling**

64 **B.1 Default Scheduling is Sub-optimal When Coupled with Device Placement**

65 Many of the large machine learning models require partitioning of the dataflow graph for the model
66 over a set of heterogeneous devices (like CPUs and GPUs). This partitioning is necessary to overcome
67 the memory limitations of a single device while also reducing the step time (i.e., execution time of
68 a single training step) by exploiting the inherent parallelism in the model architecture. However,
69 the best partitioning of the graph over a set of devices is a complex non-differentiable function of
70 the graph structure and computation capabilities of the devices. Furthermore, even for a seemingly
71 good partitioning and placement of the graph, poor choices in scheduling operations can often lead
72 to near-sequential execution of the partitioned blocks and therefore poor runtimes and low device
73 utilization. For instance, consider the timelines for a single step of Transformer XL model with
74 identical partitioning and placement shown in Figure 1. While the first timeline schedules operations
75 first-in-first-out (FIFO), the second timeline uses a human expert provided schedule. By carefully
76 choosing the ordering in which operations are scheduled, we can reduce step time (as illustrated

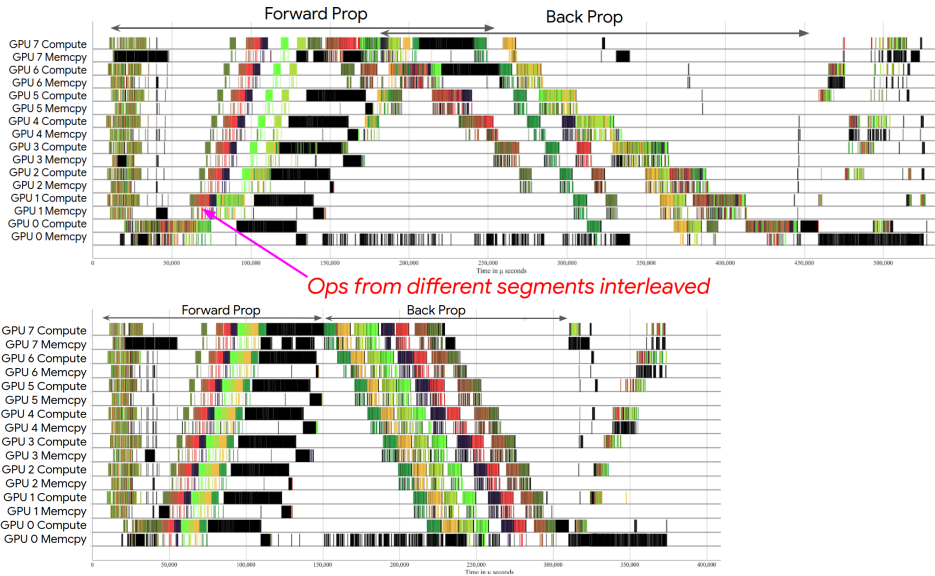


Figure 1: Effect of human expert scheduling on timelines for Transformer-XL workload. Step time is reduced by 25.5% from (a) to (b) as shown by the shorter timeline (compressed x-axis). Higher device utilization is observed in (b), with less idle time (as represented by white space). One layer is placed on each device (each row), with colors denoting segments. Each successive device can start on a segment only after it is completed on the previous device.

77 through a shorter timeline) as well as increase device utilization (as illustrated through less white
 78 space in the timelines).

79 B.2 Ops Scheduling for Individual Graphs

80 **Baseline – Heuristics.** We use *First-In-First Out (FIFO)*, which is the default scheduling policy in
 81 TensorFlow, and a static *Fanout heuristic*, which assigns priorities based on the number of dependent
 82 operations as baselines. Fanout heuristic aims to be a simplification of heuristics such as critical-path-
 83 first, and aims to capture the sentiment of prioritizing operations that other ops depend on. In our
 84 experiments, Fanout heuristic utilizes the priority scheduler implementation while FIFO numbers are
 85 from the default scheduler.

86 **Baseline – Human expert optimization.** To show the potential for scheduling, we manually added
 87 artificial control dependencies to our workloads. This required delving into timelines, diagnosing
 88 scheduling issues, and identifying which of the 49K nodes to add dependencies between. After
 89 adding 21 dependencies, we reduced the step time of an 8-layer Transformer-XL model by 25.5%
 90 as shown in Figure 1. With FIFO scheduling, we note that operations from different segments are
 91 interleaved on each device, which we infer to be the cause of unnecessary waiting on subsequent
 92 devices.

93 **Setup.** We evaluate our approach on the workloads RNNLM [11, 4] and Transformer-XL [2]. These
 94 workloads are used for language modeling and represent complex state-of-the-art models with grid-
 95 like dependencies that offer opportunities for model parallelism. We use a manual placement of
 96 assigning one layer to each GPU, and set the number of layers and thus dictating the size of the
 97 model to be equal to the number of GPUs. We use a batch size of 64 for training, and a GPU cluster
 98 topology consisting of a single machine with 8 Nvidia Tesla P100s, each with 16 GB of memory and
 99 NVLink interconnect.

100 We ran simulated annealing twice with different random seeds but same starting point, with 5000
 101 iterations each. We trained the RL model with a replay buffer and a rollout of 400 actions per iteration.
 102 To speed up the search, we use tf-sim¹ to estimate the step time. It considers operations’ running
 103 time characteristics and data transfer times on target hardware.

¹TensorFlow Performance Simulator. Publication in progress.

Table 3: Effect of scheduling priorities on reductions in step times (relative to default FIFO policy).

Workload	Layers	FH	SA	RL
RNNLM	2	-0.84%	2.96%	8.77%
RNNLM	4	-0.74%	12.01%	7.85%
RNNLM	8	-0.59%	14.65%	4.30%
Transformer-XL	2	-4.86%	6.30%	14.46%
Transformer-XL	4	-3.17%	9.27%	26.98%
Transformer-XL	8	-0.05%	19.27%	10.67%

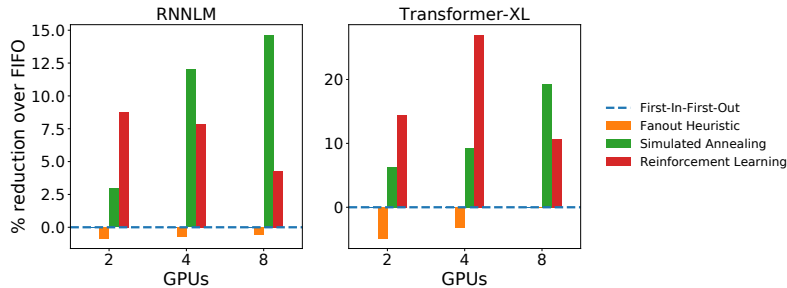


Figure 2: Effect of scheduling priorities on reductions in step times.

104 **SA and RL priority scheduling reduce step time.** Table 3 and Figure 2 show step time reductions
 105 from our simulated runs, where the step time is the time needed to train one minibatch. We observe
 106 that simulated annealing (SA) performs strongly, achieving the most reduction in step time for half of
 107 the workloads. SA was difficult to tune to get good results, especially in designing the action space
 108 and finding a set of hyperparameters that could work well for all workloads. Given that we did not do
 109 such tuning for the RL model, its performance should be considered positively given that it was a
 110 more automated solution and still outperformed SA for half the workloads in addition to achieving
 111 up to 26.98% over FIFO, which was greater than that of SA’s 19.27% over FIFO.

112 From examining timelines such as those in Figure 1, we expected the potential amount of scheduling
 113 badness (and thus maximum possible speedup from scheduling) to grow near linearly with the number
 114 of devices. This is consistent with the trend of the maximum achieved speedup scaling with the
 115 number of GPUs. We observed that the Fanout heuristic performs worse than FIFO scheduling,
 116 underlining the difficulty of designing a heuristic that works well across all real workloads.

117 C Ablation Study

118 C.1 Model Architecture

119 As we discussed model architecture alternatives and explained how we designed the network archi-
 120 tecture in Section 4.2 in the main paper, we provide more empirical results compared to a decision
 121 network based on MLPs and a decision network based on Graph Attention Networks (GATs). The
 122 MLPs is combined with GraphSAGE to provide generalization while the GATs supports general-
 123 ization naturally. We don’t compare with an LSTM or a vanilla Transformer model as they cannot
 124 scale to problem size interesting enough. According to Figure 3, for many bigger models (e.g.
 125 AmoebaNet, GNMT, 8-layer RNNLM and Transformer-XL), GATs failed to generate any valid graph
 126 optimizations. For the workloads they can generate valid optimizations, GATs is on average 8%
 127 worse and MLP is on average 11% worse in the optimized graph run time, compared to GO-one. The
 128 results also shows that the global attention brings in about 11% additional performance, compared to
 129 non-attention based decision network.

130 C.2 Fine-tuning Results

131 We also evaluate a fine-tuning strategy by pre-training the graph embedding and placement network
 132 and fine-tuning the network on the down stream tasks. The difference here compared to the main

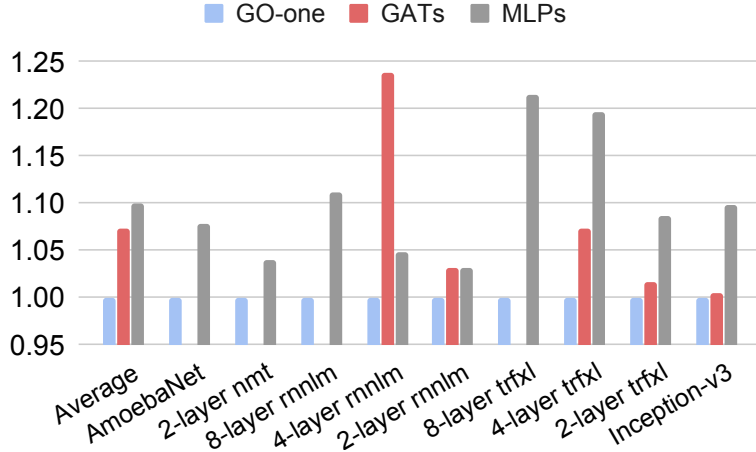


Figure 3: Lrnt Policy Comparison with Alternative Decision Networks.

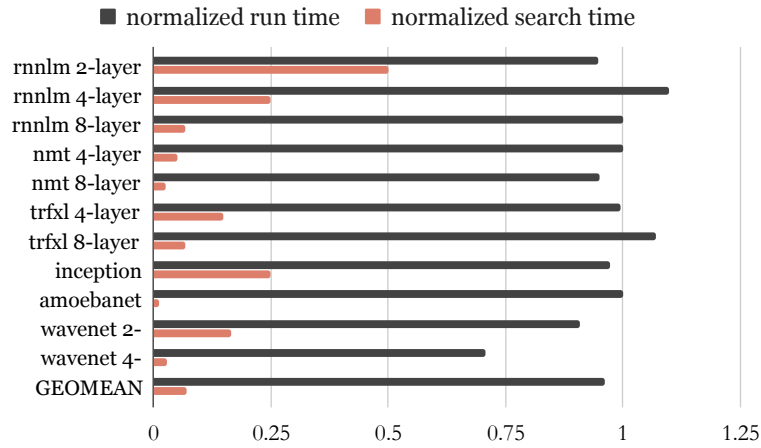


Figure 4: Normalized run time (step time for the generated placement) and normalized training time (search time) for fine-tuning. Time is normalized to GO-one.

133 paper Section 5.1 "Generalization Across Graphs" is that we also include the target graphs in the
 134 pre-training dataset. When GO-batch is used as a pre-training strategy, the graph embedding and
 135 placement network assimilate meaningful graph representations and placement policies from a wide
 136 set of graphs, thus can be used as a strong baseline network for fine-tuning on downstream tasks. We
 137 compare the generated placement run time and the placement search time, normalized to GO-one. We
 138 find that fine-tuning further reduces the the placed graph run time by an average of 5% and placement
 139 search time by an average of 86%, compared to GO-one.

140 References

- 141 [1] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly
 142 learning to align and translate. In *ICLR 2015*. 2015.
- 143 [2] Zihang Dai, Zhilin Yang, Yiming Yang, Jaime G. Carbonell, Quoc V. Le, and Ruslan Salakhut-
 144 dinov. Transformer-xl: Attentive language models beyond a fixed-length context. *ACL*, 2019.
- 145 [3] Google. TF-Sim: TensorFlow Performance Simulator, under submission. 2020.
- 146 [4] Rafal Jozefowicz, Oriol Vinyals, Mike Schuster, Noam Shazeer, and Yonghui Wu. Exploring
 147 the limits of language modeling. *arXiv preprint arXiv:1602.02410*, 2016.
- 148 [5] Esteban Real, Alok Aggarwal, Yanping Huang, and Quoc V. Le. Regularized evolution for
 149 image classifier architecture search. *CoRR*, abs/1802.01548, 2018.

- 150 [6] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna.
151 Rethinking the inception architecture for computer vision. *CoRR*, abs/1512.00567, 2015.
- 152 [7] Aäron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex
153 Graves, Nal Kalchbrenner, Andrew W. Senior, and Koray Kavukcuoglu. Wavenet: A generative
154 model for raw audio. *CoRR*, abs/1609.03499, 2016.
- 155 [8] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
156 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. V. Luxburg,
157 S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural
158 Information Processing Systems 30*, pages 5998–6008. Curran Associates, Inc., 2017.
- 159 [9] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang
160 Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah,
161 Melvin Johnson, Xiaobing Liu, Lukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo,
162 Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason
163 Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, and Jeffrey
164 Dean. Google’s neural machine translation system: Bridging the gap between human and
165 machine translation. *CoRR*, abs/1609.08144, 2016.
- 166 [10] Wojciech Zaremba, Ilya Sutskever, and Oriol Vinyals. Recurrent neural network regularization.
167 *CoRR*, abs/1409.2329, 2014.
- 168 [11] Wojciech Zaremba, Ilya Sutskever, and Oriol Vinyals. Recurrent neural network regularization.
169 *arXiv preprint arXiv:1409.2329*, 2014.