

1 We thank all reviewers for their encouraging and constructive feedback.

2 **Rev#1:** We thank the reviewer for their suggestion to tease apart the impact of architectural decisions, have a clearer
3 elaboration of the fusion task, and include full neural network model formulation. We will include additional details
4 and ablation studies in the revised version to further clarify each of these aspects.

5 **[Fusion priority]** As mentioned in Sec 3 of the paper, each node is associated with a fusion priority score (ranging
6 from 1 to F , where F is a hyperparameter). The scores determine the order in which nodes are fused.

7 **[Better figure]** The reviewer is correct that each Transformer layer in Fig.3 only outputs a feature embedding that is
8 used by the next layer, and only the final layer computes actions. In Eq.2, the actions are produced by the final layer in
9 Fig.3. However, the temporal dependencies between actions in multiple RL steps are carried out by state embedding
10 such that the state of the previous step is used as an input to the network in the next step. Fig.3 is an illustration of
11 a single task, while Fig.4 demonstrates the extension of the same framework for multiple tasks. We add a few more
12 Transformer layers with residual connections as task heads. The yellow part in Fig.4 is the same as Fig.3.

13 **[Why not use the first layer of the TRF to replace the modulation layer?]** Having a separate modulation layer
14 provides additional parameterization to the network and hence better orchestration of the training of the parameters. It
15 empirically results in better generalization.

16 **Rev#2:** We thank the reviewer’s insightful feedback on providing more details and potential open sourcing for
17 reproducibility. We will add more detailed explanations about GNN’s limitations on tracking global node dependencies.

18 **[Reproducibility]** We made an extensive effort to document the details of our network architecture, the input graphs,
19 and other hyperparameters in Sec. 4 and Supp.A. We plan to incorporate further details to enhance reproducibility and
20 open source the GO framework along with the performance model. We would also like to emphasize that GO, as a
21 general algorithm, can be applied to problems at different layers of the compilation stack, even on different platforms.
22 We have successfully applied GO to two other graph optimization problems at different stages in the compiler stack –
23 tensor layout optimization and operator fusion in XLA (as opposed to at the TF level discussed in this paper).

24 **[Can GO be trained online?]** Yes, GO can be trained in an on-line scenario similar to Decima. Specifically, GO-one
25 is quite similar to an online training method where state transitions and rewards can be collected by interacting with the
26 compiler. However, as pointed out, an online RL can make poor decisions in early stages of training and the quality
27 of training is subjective to input data distribution. This is particular challenging for a compilation problem, where a
28 program needs to be compiled within a short time. One possible approach is to pre-train the model like GO-finetune,
29 then apply online training to continuously improve the policy.

30 **Rev#3:** We appreciate the reviewer’s detailed feedback and will address the questions in the final version.

31 **[Reproducibility and validation]** Please see our response to Rev#2. The accuracy of the performance model has been
32 validated against true runtime measurements (on actual hardware) on several industry standard models, including MLP,
33 CNN, RNN, LSTM, Transformer, BERT, etc. We will provide more validation results in the final paper.

34 **[Table1: is there a methodology for these yes/no distinctions?]** We empirically observed these distinctions. For
35 example, a vanilla Transformer network will be OOM for problems larger than 10k nodes, and a vanilla LSTM cannot
36 generate placement decisions for input sequence longer than a few hundreds. Ablation study in Supp.C.1 addresses
37 some of the distinctions.

38 **[Source for the claim the workloads are realistic]** We used important workloads that are widely used in real
39 applications or are incorporated in industry standard benchmarks like MLPerf. These include ResNet, InceptionNet,
40 WaveNet, Transformer-XL, AmoebaNet, NMT, and RNNLM.

41 **[SA convergence time]** SA takes around 24 hours to converge in our evaluated tasks. For a smaller graph, GO takes
42 less than an hour to converge. For a large graph over 10k nodes, GO can take 1-4 hours to converge.

43 **[Better figure]** We acknowledge that Figure 3 can be improved, and we will incorporate the suggestions from the
44 reviewer to the final version.

45 **[More details of the performance model]** Exact features used in node embedding (as presented in Fig.3 dashed box):
46 tensor shape, op type, and adjacency matrix. Fusion is optimized using a priority-based node traversal. Please also see
47 reponse to Rev#1 [Fusion priority].

48 **[Is this function $f^{(l+1)}$ different each iteration?]** Yes. l indexes the layer of the graph neural net and different layers
49 of the graph neural net have different f .

50 **Rev#4:** As mentioned in response to Rev#2, we plan to open source the GO framework along with the performance
51 model, and will include details about this in the final version.