
TOWARDS REAL-TIME 3D OBJECT DETECTION FOR AUTONOMOUS VEHICLES WITH PRUNING SEARCH

Pu Zhao¹ Wei Niu² Geng Yuan¹ Yuxuan Cai¹ Bin Ren² Yanzhi Wang¹ Xue Lin¹

ABSTRACT

In autonomous driving, 3D object detection is essential as it provides basic knowledge about the environment. However, as deep learning based 3D detection methods are usually computation intensive, it is challenging to support real-time 3D detection on edge-computing devices with limited computation and memory resources. To facilitate this, we propose a compiler-aware pruning search framework, to achieve real-time inference of 3D object detection on the resource-limited mobile devices. Specifically, a generator is applied to sample better pruning proposals in the search space, and an evaluator is adopted to evaluate the sampled pruning proposal performance with Bayesian optimization. We demonstrate that the pruning search framework can achieve real-time 3D object detection on mobile (Samsung Galaxy S20 phone) with state-of-the-art detection performance.

1 INTRODUCTION

As the rapid development of autonomous vehicles to self-drive without human intervention, object detection (especially 3D detection to deal with LiDAR data) serves as a fundamental prerequisite for autonomous navigation. 3D detection can extract desirable knowledge about its environment from 3D point clouds of LiDAR sensors, thus enabling high-level computations and optimizations for auto-driving.

It is essential to implement real-time 3D object detection on autonomous vehicles. However, the current deep neural networks (DNNs) based 3D object detectors usually cost tremendous memory and computation resources, leading to difficulties for real-time implementation, especially on autonomous vehicles with limited hardware resource.

To reduce the DNN model size and computations, DNN weight pruning (Guo et al., 2016; Wen et al., 2016) has shown great advantages to remove redundancy in the model, therefore reducing storage/computation cost and accelerating inference. There are *unstructured pruning* scheme (Han et al., 2015; Liu et al., 2018b) to remove arbitrary weight, *coarse-grained structured pruning* scheme (He et al., 2017; Liu et al., 2018b; Wen et al., 2016) to eliminate whole filters/channels, and *fine-grained structured pruning* (Ma et al., 2020; Niu et al., 2020) to assign different pruning patterns to convolutional (CONV) kernels. Compared with unstructured pruning, structured pruning can achieve higher

hardware parallelism and inference acceleration, assisted by compiler-level code generation and optimization techniques (Niu et al., 2020), with competitive detection performance.

Though compiler optimization can support various structured pruning schemes with notable mobile acceleration performance, different sparsity schemes lead to different accuracy and acceleration performance with compiler optimization. For the specific 3D detection problem, it is still questionable to adopt which sparsity scheme with which pruning rate to satisfy the accuracy and real-time requirements. To find the pruning solution, motivated by the idea of Neural Architecture Search (NAS) (Liu et al., 2018a; Zoph & Le, 2017), we propose a compiler-aware pruning search framework to automatically determine the pruning scheme and pruning rate for each individual layer. The *objective* is to maximize accuracy with an inference speed/latency constraint on the target mobile device. Different from previous work with fixed pruning scheme for all layers, our work can have different pruning schemes and rates for different layers in the model. We summarize our contribution as follows,

- We incorporate DNN latency constraint into pruning search to satisfy a predefined real-time requirement.
- Our framework configures various pruning schemes and pruning rates for various layers, different from previous works with fixed pruning scheme for all layers.
- We adopt an ensemble of neural predictors and Bayesian optimization (BO) to reduce the number of evaluated proposals, leading to less searching efforts.
- We are the first to achieve (close-to) real-time (98ms)

¹Northeastern University, MA ²William & Mary, VA. Correspondence to: Pu Zhao <zhao.pu@northeastern.edu>, Xue Lin <xue.lin@northeastern.edu>.

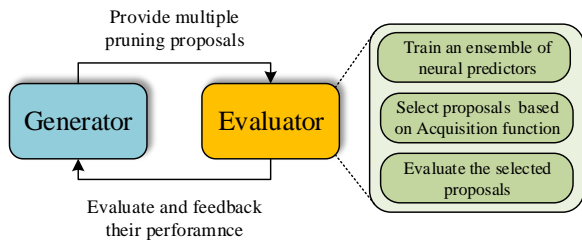


Figure 1. Automatic network pruning search framework

3D detection with PointPillars, on an off-the-shelf mobile phone with minor (or no) accuracy loss.

2 BACKGROUND AND RELATED WORK

2.1 3D Object Detection

3D object detection detects objects with point clouds from LiDAR sensors. PointPillars (Lang et al., 2019) is a popular 3D detection method with three main stages: (1) A feature encoder network to convert a point cloud to a sparse pseudo-image; (2) a 2D CONV backbone to transform the pseudo-image into high-level representation; and (3) a detection head to regress 3D boxes. Besides PointPillars, there are various 3D detection methods such as SECOND (Yan et al., 2018) and Point-GNN (Shi & Rajkumar, 2020). We mainly focus on PointPillars as we found that PointPillars is the only one which can run on mobile while others are not available on mobile since their special structures to deal with sparse data are not supported by mobile compiler.

2.2 Weight Pruning Schemes

Previous weight pruning work can be categorized into: unstructured pruning (Guo et al., 2016; Han et al., 2015; Mao et al., 2017), coarse-grained structured pruning (He et al., 2017; Luo et al., 2017; Wen et al., 2016), and fine-grained structured pruning including pattern (Ma et al., 2020; Niu et al., 2020) and block (Dong et al., 2020) pruning.

Unstructured pruning (Han et al., 2015; Mao et al., 2017) removes weights at arbitrary positions, leading to irregular sparse weight matrix with indices, incurring damages to the parallel implementations and acceleration performance on hardware. Different from unstructured pruning, coarse-grained structured pruning (He et al., 2017; Liu et al., 2018b) removes the whole filters/channels to maintain model structure with high regularity for efficient hardware parallel implementation, at the cost of certain obvious accuracy degradation. To overcome the disadvantages, fine-grained structured pruning (Ma et al., 2020; Niu et al., 2020) follows a pruning pattern (chosen from a predefined library) to prune each CONV kernel, where the predefined patterns have been optimized with compiler optimizations for mobile acceleration. Fine-grained structured pruning can achieve

Table 1. Search space for each DNN layer

Pruning scheme	{Filter, Pattern-based, Block-based}
Pruning rate	{1×, 2×, 3×, 5×, 7×, 10×, 15×}

high accuracy due to the flexibility with different patterns, and high hardware parallelism (and mobile acceleration) with compiler-based code generation and optimization.

3 AUTOMATIC NEURAL PRUNING SEARCH

We show the framework in Fig. 1, consisting of two basic components: a *generator* and an *evaluator*. Given the search space, the generator first generates or samples various *pruning proposals*. Then the evaluator evaluates their detection accuracy and speed performance, and feeds them back to the generator. Next the generator samples new pruning proposals based on existing proposals’ performance. After iterations, the framework can obtain the final pruning proposal with satisfying detection and speed performance.

In each iteration, the evaluator first trains an ensemble of neural predictors and then selects proposals based on their acquisition function values enabled by the predictor ensemble. Next the selected proposals are evaluated to obtain their performance while the rest unselected proposals are not evaluated, thus reducing evaluation time and efforts.

After the framework finishes and outputs a final pruning proposal, we further apply ADMM pruning (Zhang et al., 2018) to perform an enhanced pruning following the best proposal. Compared with the simple magnitude pruning (Han et al., 2015) method applied during evaluation for time-saving, ADMM usually outperforms magnitude pruning in terms of accuracy with an increased complexity, that is why we only adopt it for the final proposal.

3.1 Generator

The generator samples *pruning proposals* from the search space. Each pruning proposal is a directed graph consisting of the layer-wise pruning scheme and pruning rate. For example, it has 20 nodes for a 10-layer DNN model.

3.1.1 Proposal Formulation (Search Space)

Each pruning proposal contains the pruning scheme and pruning rate for each layer of the model, as shown in Tab. 1.

Per-layer pruning schemes: The generator can choose from filter (channel) pruning (Zhuang et al., 2018), pattern-based pruning (Ma et al., 2020) and block-based pruning (Dong et al., 2020) for each layer. As different layers may have different best-suited pruning schemes, the generator can choose different pruning schemes for different layers, also supported by our compiler code generation.

Per-layer pruning rate: The pruning rate is the rate between the number of original parameters and that of left parameters after pruning. We can choose from the list $\{1\times, 2\times, 3\times, 5\times, 7\times, 10\times, 15\times\}$, where $1\times$ means the layer is not pruned (i.e., bypassing this layer).

3.1.2 Proposal Updating

The generator keeps a record of all evaluated pruning proposals with their evaluation performance. To generate new pruning proposals, it mutates the best proposals in the records by randomly changing one pruning scheme or one pruning rate of one layer. More specifically, it first selects K proposals with the highest evaluation performance, and mutates each of them iteratively until it gets C new proposals.

3.1.3 Proposal encoding

As pruning proposals are basically graphs, special attention is required for the proposal representation. Different from traditional representations with an adjacency matrix for graphs, we adopt the pruning encoding to encode each proposal with a vector of binary values. There is a binary feature for each possible node in each layer, denoting whether the node (pruning scheme or pruning rate of certain layer) is adopted or not. To encode a proposal, we simply check which pruning scheme or rate for each layer is applied, and set the corresponding features to 1s. This simple proposal encoding can help with proposal evaluation.

3.2 Evaluator

The evaluator needs to evaluate proposal performance. We define the performance measurement (reward) as:

$$m = V - \alpha \cdot \max(0, r - R), \quad (1)$$

where V is the validation mean average precision (mAP) of the model, r is the model inference latency, which is actually measured on a mobile device with compiler code optimization and generation for inference acceleration¹. R is the threshold for the real-time requirement. Generally, satisfying real-time requirement ($r < R$) with high mAP leads to high m . Otherwise if mAP is low or the real-time requirement is violated, m is small.

3.2.1 Fast Evaluation with BO

As it incurs large time cost to evaluate the performance of each pruning proposal (including pruning and retraining the model with multiple epochs), we adopt Bayesian optimization (BO) (Chen et al., 2018) to accelerate evaluation. The generator provides C pruning proposals, and the evaluator first use BO to select B proposal with potentially better performance. Next the evaluator measure the performance

¹The compiler code generation can be performed in parallel with the pruning proposal evaluation to obtain its speed.

Algorithm 1 Evaluation with predictor ensemble & BO

Input: Observation data \mathcal{D} , BO batch size B , BO acquisition function $\phi(\cdot)$
Output: The best pruning proposal g
for steps do
 Generate a pool of candidate pruning proposals \mathcal{G}_c ;
 Train an ensemble of neural predictors with \mathcal{D} ;
 Select $\{\hat{g}^i\}_{i=1}^B = \arg \max_{g \in \mathcal{G}_c} \phi(g)$;
 Evaluate the proposal and obtain reward $\{r^i\}_{i=1}^B$ of $\{\hat{g}^i\}_{i=1}^B$;
 $\mathcal{D} \leftarrow \mathcal{D} \cup (\{\hat{g}^i\}_{i=1}^B, \{r^i\}_{i=1}^B)$;
end for

of the selected proposals by pruning following each selected proposal to obtain its accurate accuracy and speed performance, while the rest unselected proposals are not evaluated. Thus, the number of actual evaluated proposals is reduced to save the overall pruning evaluation efforts.

In general, there are two main components in BO including training an ensemble of neural predictors and selecting proposal based on acquisition function values enabled by the predictor ensemble. To make use of BO, the ensemble of neural predictors provides an accuracy prediction with its corresponding uncertainty estimate for an unseen pruning proposal. Then BO is able to choose the proposal maximizing the acquisition function. We show the full algorithm in Algorithm 1 and specify the two components next.

3.2.2 Ensemble of Neural Predictors

We use a neural network repeatedly trained on the current set of evaluated pruning proposals with their evaluation performance as a neural predictor to predict the reward (incorporating the accuracy and speed performance) of unseen pruning proposals. The neural network is a sequential fully-connected network with 8 layers of width 30 trained by the Adam optimizer with a learning rate of 0.01. Note that it does not cost much predictor training efforts due to their simple architectures and parallel training.

For the loss function in neural predictors, mean absolute percentage error (MAPE) is adopted as it can give a higher weight to pruning proposals with higher performance:

$$\mathcal{L}(m_{\text{pred}}, m_{\text{true}}) = \frac{1}{n} \sum_{i=1}^n \left| \frac{m_{\text{pred}}^{(i)} - m_{\text{UB}}}{m_{\text{true}}^{(i)} - m_{\text{UB}}} - 1 \right|, \quad (2)$$

where $m_{\text{pred}}^{(i)}$ and $m_{\text{true}}^{(i)}$ are the predicted and true values of the reward for the i -th proposal in a batch, and m_{UB} is a global upper bound on the maximum true reward.

To incorporate BO, it also needs an uncertainty estimate for the prediction. So we adopt an ensemble of neural predictors to provide the uncertainty estimate. More specifically, we train P neural predictors using different random weight initializations and training data orders. Then for any proposal, we can obtain the mean and standard deviation of these P predictions. More specifically, we train an ensemble

Table 2. Comparison of various pruning methods for PointPillars

Methods (grid size)	Para. #	Comp. # (MACs)	Speed (ms)	Car 3D detection		
				Easy	Moderate	Hard
PointPillars (0.16)	5.8M	60G	553	85.16	74.39	69.42
Filter (0.16)	1.1M	10.8G	178	80.63	67.51	65.28
Pattern (0.16)	1.1M	10.7G	225	83.64	74.30	68.42
Block (0.16)	1.1M	10.7G	268	82.86	75.43	69.71
Ours (0.16)	1.1M	10.7G	193	85.52	76.69	70.10
PointPillars (0.24)	5.4M	28G	253	84.24	75.28	68.46
Filter (0.24)	0.8M	4.0G	82	81.36	68.06	65.77
Pattern (0.24)	0.8M	3.9G	116	82.16	73.93	68.25
Block (0.24)	0.8M	4.0G	140	83.69	74.09	68.06
Ours (0.24)	0.8M	3.9G	98	85.38	75.72	68.53

of P predictive models, $\{f_p\}_{p=1}^P$, where $f_p : A \rightarrow \mathbb{R}$ with a pruning proposal g as input and the predicted reward as output. The mean prediction and its deviation are given by,

$$\hat{f}(g) = \frac{1}{P} \sum_{p=1}^P f_p(g), \text{ and } \hat{\sigma}(g) = \sqrt{\frac{\sum_{p=1}^P (f_p(g) - \hat{f}(g))^2}{P-1}}. \quad (3)$$

3.2.3 Selection with Acquisition Function

After training an ensemble of neural predictors, we can obtain the acquisition function value for proposals in the pool and select a small part of proposals with largest acquisition function values. We choose upper confidence bound (UCB) (Srinivas et al., 2010) as the acquisition function, i.e.,

$$\phi_{\text{UCB}}(g) = \hat{f}(g) + \beta \hat{\sigma}(g) \quad (4)$$

where the tradeoff parameter β is set to 0.5.

3.2.4 Evaluation with Magnitude Pruning

After selecting pruning proposals from the pool, the evaluator uses magnitude based pruning framework (Han et al., 2015) (with two steps including pruning and retraining) to perform the actual pruning and obtain its evaluation performance for the proposal. Besides, once the pruning step is finished, the speed measurement on a mobile device can be performed in parallel with the accuracy measurement.

4 EXPERIMENTAL RESULTS

4.1 Experiment Setup

We focus on 3D detection and employ the PointPillars (Lang et al., 2019) as starting point and test on KITTI dataset (Geiger et al., 2012). We use 40 GPUs for parallel training and pruning search and it takes about 6 days to find the best pruning proposal. In Eq. (1), we set α to 0.01 and the inference time is measured in milliseconds. The pool size C is set to 50 and the Bayesian batch size B is set to 10. We test the speed performance on the mobile GPU (Qualcomm Adreno 640) of a Samsung Galaxy S20 smartphone.

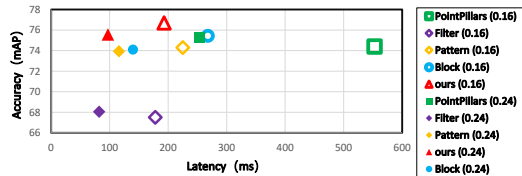


Figure 2. Comparison with other methods

4.2 Performance on 3D Object Detection

As shown in Tab. 2 and Fig. 2, we compare the performance of the original unpruned PointPillars model and the model derived by our method and other pruning methods with different grid sizes (0.16m and 0.24m). We set the threshold of the real-time requirements to 200ms for 0.16m grid size and 100ms for 0.24m. For the grid size, as large grid size leads to small pseudo-image input size for the model, the 0.24m grid size has a smaller parameter and computation numbers, and a faster inference speed on mobile GPUs, compared with the 0.16m grid size.

For the same grid size, compared with the original unpruned PointPillars model, we observe that our method can significantly reduce the number of parameters and computations, achieving state-of-the-art detection performance while satisfying the real-time requirement. The accuracy of our method is even higher than the unpruned model, demonstrating that the unpruned model may suffer from over-fitting problem and removing the redundancy can help with its accuracy.

We also compare with other pruning methods for the same grid size. For other pruning methods, the same pruning scheme is applied to all layers and the pruning rate is set to the same with the overall pruning ratio of our pruned model (80% for grid size 0.16m and 86% for 0.24m). As observed, the proposed method achieves the best detection performance with highest accuracy compared with other methods with the same pruning scheme for every layer, demonstrating the advantages of using different pruning scheme for different layers. We notice that although filter pruning can be faster than our method, it suffers from an obvious degradation on the detection performance.

For the speed performance, with grid size 0.24m, the proposed method only needs 98ms to process one frame on mobile devices with the highest accuracy, demonstrating its superior performance to achieve (close-to) real-time inference on mobile with state-of-the-art detection performance.

5 CONCLUSION

We propose pruning search to configure the pruning scheme and rate for each layer with real-time inference requirement. Our experiments demonstrate that the proposed method achieves (close-to) real-time (98ms) 3D object detection on a mobile phone with minor (or no) accuracy loss.

REFERENCES

- Chen, Y., Huang, A., et al. Bayesian optimization in alphago. *arXiv:1812.06855*, 2018.
- Dong, P., Wang, S., et al. Rtmobile: Beyond real-time mobile acceleration of rnns for speech recognition. *arXiv:2002.11474*, 2020.
- Geiger, A., Lenz, P., and Urtasun, R. Are we ready for autonomous driving? the kitti vision benchmark suite. In *CVPR*, 2012.
- Guo, Y., Yao, A., and Chen, Y. Dynamic network surgery for efficient dnns. In *NeurIPS*, pp. 1379–1387, 2016.
- Han, S., Pool, J., et al. Learning both weights and connections for efficient neural network. In *NeurIPS*, pp. 1135–1143, 2015.
- He, Y., Zhang, X., and Sun, J. Channel pruning for accelerating very deep neural networks. In *ICCV*, pp. 1389–1397, 2017.
- Lang, A. H., Vora, S., et al. Pointpillars: Fast encoders for object detection from point clouds. In *CVPR*, pp. 12697–12705, 2019.
- Liu, H., Simonyan, K., and Yang, Y. Darts: Differentiable architecture search. *arXiv preprint arXiv:1806.09055*, 2018a.
- Liu, Z., Sun, M., Zhou, T., Huang, G., and Darrell, T. Rethinking the value of network pruning. *arXiv preprint arXiv:1810.05270*, 2018b.
- Luo, J.-H., Wu, J., and Lin, W. Thinet: A filter level pruning method for deep neural network compression. In *ICCV*, pp. 5058–5066, 2017.
- Ma, X. et al. Pconv: The missing but desirable sparsity in dnn weight pruning for real-time execution on mobile devices. In *AAAI*, 2020.
- Mao, H., Han, S., et al. Exploring the regularity of sparse structure in convolutional neural networks. *arXiv:1705.08922*, 2017.
- Niu, W. et al. Patdnn: Achieving real-time dnn execution on mobile devices with pattern-based weight pruning. *arXiv:2001.00138*, 2020.
- Shi, W. and Rajkumar, R. Point-gnn: Graph neural network for 3d object detection in a point cloud. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 1711–1719, 2020.
- Srinivas, N., Krause, A., et al. Gaussian process optimization in the bandit setting: No regret and experimental design. In *ICML*, 2010. ISBN 9781605589077.
- Wen, W., Wu, C., et al. Learning structured sparsity in deep neural networks. In *NeurIPS*, pp. 2074–2082, 2016.
- Yan, Y., Mao, Y., and Li, B. Second: Sparsely embedded convolutional detection. *Sensors*, 18(10):3337, 2018.
- Zhang, T., Ye, S., et al. Systematic weight pruning of dnns using alternating direction method of multipliers. *ECCV*, 2018.
- Zhuang, Z., Tan, M., et al. Discrimination-aware channel pruning for deep neural networks. In *NeurIPS*, pp. 875–886, 2018.
- Zoph, B. and Le, Q. V. Neural architecture search with reinforcement learning. In *ICLR*, 2017.