# Multi-Frame to Single-Frame:
# Knowledge Distillation for 3D Object Detection

Yue Wang[1], Alireza Fathi[2], Jiajun Wu[3],
Thomas Funkhouser[2], and Justin Solomon[1]

[1] MIT
{yuewangx,jsolomon}@mit.edu
[2] Google
{alirezafathi,tfunkhouser}@google.com
[3] Stanford
jiajunwu@cs.stanford.edu

**Abstract.** A common dilemma in 3D object detection for autonomous driving is that high-quality, dense point clouds are only available during training, but not testing. We use knowledge distillation to bridge the gap between a model trained on high-quality inputs at training time and another tested on low-quality inputs at inference time. In particular, we design a two-stage training pipeline for point cloud object detection. First, we train an object detection model on dense point clouds, which are generated from multiple frames using extra information only available at training time. Then, we train the model's identical counterpart on sparse single-frame point clouds with consistency regularization on features from both models. We show that this procedure improves performance on low-quality data during testing, without additional overhead.

## 1 Introduction

Recent advances in 3D object detection are beginning to yield practical vision systems for autonomous driving. Given the realities of data collection and allowable inference time in this application, however, most relevant 3D object detection methods take static one-frame point clouds as input, which can be extremely sparse. The goal of matching object detection rates achievable from dense point clouds constructed by fusing multiple frames remains elusive. On the other hand, aggregating point clouds from multiple frames is inherently hard in dynamical settings. For example, Figure 1 (b) shows a point cloud generated by combining 5 LiDAR scans, which exhibits significant motion blur; this blur—unavoidable in a moving environment—yields ambiguity in the object detection model.

Point cloud object detection training data often comes in *sequences*, with consistent labels that track moving objects between frames in each sequence; this extra information allows us to generate Figure 1 (c), which shows an aggregated point cloud after accounting for the motion of individual objects. This tracking information, however, is generally unavailable at testing time. This mismatch between training and testing data leads to an unavoidable dilemma: It is desirable

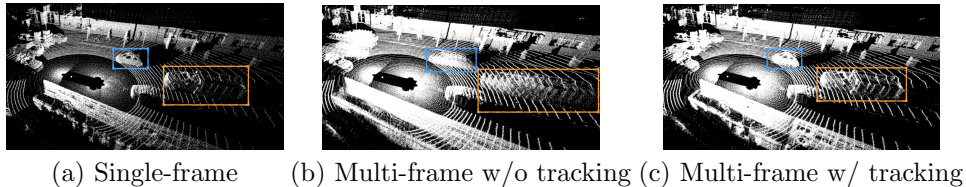(a) Single-frame      (b) Multi-frame w/o tracking (c) Multi-frame w/ tracking

Fig. 1: Visualizations of single/multi-frame point clouds. (a) single-frame point clouds; (b) multi-frame point clouds without bounding box tracking; (c) multi-frame point clouds with bounding box tracking. Two examples are labeled in blue and orange; the moving cars are extremely blurry in (b) while not in (c).

to build a model that gathers dense point clouds from multiple frames, but in the absence of fine-tuned registration procedures this data is only available at training time.

In this paper, we introduce a two-stage training pipeline to bridge the gap between training on dense point clouds and testing on sparse point clouds. In principle, our method is an intuitive extension and application of knowledge distillation [1,2] and learning using privileged information [13,6,11]. First, we train an object detection model on dense point clouds aggregated from multiple frames, using information available during training but not testing. Then, we train an identical model on single-frame sparse point clouds; we use the pre-trained model from the first stage to provide additional supervision for the second model. At test time, the second model makes predictions on single-frame point clouds, without extra inference overhead. Compared with a model trained solely on sparse point clouds, the proposed model benefits from distilled information gathered by the model trained on dense point clouds.

## 2   Method

We study how to use privileged information available at training time in the context of point cloud object detection. In contrast to existing works [11,4], which use additional image information and require extra annotations, our model takes advantage of the privileged information that is naturally available in a sequence of point clouds. Our model bridges the gap between the single-frame model and the (impractical) multi-frame model, while maintaining efficiency.

Consider a sequence of point clouds $\{\mathcal{X}_0, \mathcal{X}_1, \ldots, \mathcal{X}_t, \ldots, \mathcal{X}_T\}$ up to frame $T$ with per-frame pose $\{\mathcal{T}_0, \mathcal{T}_1, \ldots, \mathcal{T}_t, \ldots, \mathcal{T}_T\}$. Each frame is a point cloud $\mathcal{X}_t = \{\boldsymbol{x}_t^0, \boldsymbol{x}_t^1, \ldots, \boldsymbol{x}_t^i, \ldots, \boldsymbol{x}_t^N\} \subset \mathbb{R}^3$ with additional point-wise features $\mathcal{F}_t = \{\boldsymbol{f}_t^0, \boldsymbol{f}_t^1, \ldots, \boldsymbol{f}_t^i, \ldots, \boldsymbol{f}_t^N\} \subset \mathbb{R}^c$, such as reflectance and material. For simplicity of notation, we assume the same number of points $N$ in each point cloud.

In each frame $t$, a set of ground-truth bounding boxes $\{\boldsymbol{b}_t^0, \boldsymbol{b}_t^1, \ldots, \boldsymbol{b}_t^k, \ldots, \boldsymbol{b}_t^K\}$ is provided, where each bounding box is parameterized by its center location $\boldsymbol{p}_t^k$, its size $\boldsymbol{s}_t^k$, and its heading angle $\boldsymbol{\theta}_t^k$. For ease of presentation, assume $\{\boldsymbol{b}_0^k, \boldsymbol{b}_1^k, \ldots, \boldsymbol{b}_t^k \ldots, \boldsymbol{b}_T^k\}$ are bounding boxes of the $k$-th object over all frames;
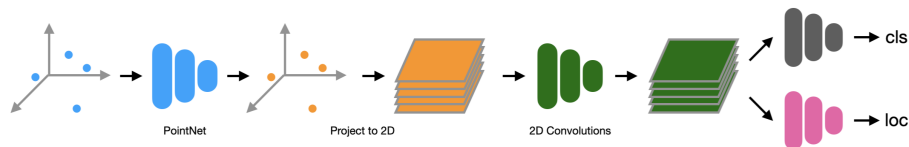
Fig. 2: Our object detection model: first, point clouds are encoded using a lightweight PointNet; then, points are projected into BEV 2D grid, followed by additional 2D convolutional layers for further feature embedding; finally, a classification branch and localization branch predict the existence of bounding boxes per BEV pixel and bounding box parameters, respectively. cls: classification branch; loc: localization branch.

that is, we assume each frame in the sequence is labeled *consistently*. We can take advantage of these bounding box annotations to generate dense point clouds as in Figure 1 (c), but this is only possible at training time.

At test time, since we can no longer access the bounding box annotations, the inputs to the model are either single-frame point clouds (Figure 1 (a)) or multi-frame point clouds (Figure 1 (b)) without accounting for moving objects. So the model trained on the desired point clouds (Figure 1 (c)) can never be deployed for inference. To address this difference between training and testing, we propose a two-stage training pipeline. First, we produce dense point clouds by registering multiple frames and train an object detection model on dense point clouds. Then, we distill the previous model to a new model by leveraging both dense point clouds and sparse point clouds. Finally, we use the new model for testing on sparse point clouds without extra cost.

## 3    Experiments

First, we introduce the dataset, metrics, implementation details, and optimization details in §3.1. Then, we show performance of the proposed method for vehicle and pedestrian detection on the Waymo Open Dataset [12] in §3.2.

### 3.1    Details

*Metrics.* We use birds-eye view (BEV) and 3D mean average precision (mAP) as metrics. The intersection-over-union (IoU) threshold is 0.7 for vehicles and 0.5 for pedestrians. In addition to overall scores, we provide breakdowns based on the distances between the origin and ground-truth boxes: 0m-30m, 30m-50m, and 50m-infinity (Inf). Following existing works  [5,9,14], we evaluate our method on bounding boxes that have more than 5 points.

*Implementation details.* As shown in Figure 2, the object detection model is a variant of PointPillars [5]: first, points in 3D are projected to BEV, with a lightweight PointNet [10] to aggregate features in each BEV pixel; then, three convolutional blocks are employed to further embed the BEV features; finally, we

| Method | BEV mAP (IoU=0.7) | | | | 3D mAP (IoU=0.7) | | | |
|---|---|---|---|---|---|---|---|---|
| | Overall | 0 - 30m | 30 - 50m | 50m - Inf | Overall | 0 - 30m | 30 - 50m | 50m - Inf |
| StarNet [9] | - | - | - | - | 53.7 | - | - | - |
| LaserNet† [8] | 71.57 | 92.94 | 74.92 | 48.87 | 55.1 | 84.9 | 53.11 | 23.92 |
| PointPillars∗ [5] | 75.57 | 92.1 | 74.06 | 55.47 | 56.62 | 81.01 | 51.75 | 27.94 |
| MVF [14] | 80.4 | 93.59 | 79.21 | 63.09 | 62.93 | 86.3 | 60.2 | 36.02 |
| Baseline (1 frame) | 83.16 | 94.2 | 81.22 | 64.43 | 63.09 | 84.73 | 58.55 | 33.55 |
| + Distillation (1 frame) | **84.79** | **94.84** | **82.69** | **68.15** | **66.6** | **87.56** | **60.97** | **38.94** |
| Improvements | **+1.63** | **+0.64** | **+1.47** | **+3.72** | **+3.51** | **+2.83** | **+2.42** | **+5.39** |
| Oracle (5 frames) ‡ | 86.31 | 95.27 | 84.46 | 70.7 | 68.79 | 87.98 | 65.09 | 41.61 |

Table 1: Results on vehicle. †: our implementation. ∗: re-implemented by [14]. ‡: the model trained and tested on dense point clouds as in Figure 1 (c). ‡ is not a realistic setup, but we provide results to serve as an oracle.

use two branches to predict bounding box existence and bounding box parameters, respectively. The frame rate at test time is 24 FPS on par with PointPillars [5].

*Optimization.* We use Adam [3] to train both the multi-frame and the single-frame model. For both models, the initial learning rate is $3 \times 10^{-4}$. We then linearly increase it to $3 \times 10^{-3}$ in the first 5 epochs. Finally it is decreased to $3 \times 10^{-6}$ using cosine scheduling [7]. The model is trained for 75 epochs at each stage. The weight $\lambda$ of feature consistency is 0.1 for vehicles and 0.01 for pedestrians. We use 5 frames for the multi-frame model. Our batch size is 128. We train our models on TPUv3. Training the multi-frame model is about 2 times slower than training the single-frame model.

### 3.2 Results on the Waymo Open Dataset

In Table 1 and Table 2, we mainly compare the model trained with our two-stage strategy to a single-frame model trained from scratch. For the sake of completeness, we also provide results of StarNet [9], LaserNet [8], PointPillars [5], and MVF [14]; compared to them, even the single-frame model achieves better performance. The model trained on 5 frames significantly outperforms others, including the single-frame model, but this is an unrealistic setup for testing—points from multiple frames are aggregated by using ground-truth tracking information. The single-frame model with distillation improves over its vanilla counterpart in all metrics, which indicates the model indeed learns additional information with the supervision from a multi-frame model.

## References

1. Buciluă, C., Caruana, R., Niculescu-Mizil, A.: Model compression. In: Special Interest Group on Knowledge Discovery and Data Mining (SIGKDD) (2006)

| Method | BEV mAP (IoU=0.5) | | | | 3D mAP (IoU=0.5) | | | |
|---|---|---|---|---|---|---|---|---|
| | Overall | 0 - 30m | 30 - 50m | 50m - Inf | Overall | 0 - 30m | 30 - 50m | 50m - Inf |
| StarNet [9] | - | - | - | - | 66.8 | - | - | - |
| LaserNet† [8] | 70.01 | 78.24 | 69.47 | 52.68 | 63.4 | 73.47 | 61.55 | 42.69 |
| PointPillars∗ [5] | 68.7 | 75.0 | 66.6 | 58.7 | 60.0 | 68.9 | 57.6 | 46.0 |
| MVF [14] | 74.38 | 80.01 | 72.98 | 62.51 | 65.33 | 72.51 | 63.35 | 50.62 |
| Baseline (1 frame) | 74.48 | 81.72 | 72.9 | 58.77 | 68.32 | 78.12 | 64.85 | 50.21 |
| + Distillation (1 frame) | **76.3** | **83.03** | **74.49** | **62.18** | **69.9** | **78.83** | **66.96** | **52.73** |
| Improvements | **+1.82** | **+1.31** | **+1.59** | **+3.41** | **+1.58** | **+0.71** | **+2.11** | **+2.52** |
| Oracle (5 frames) ‡ | 80.54 | 86.57 | 78.13 | 69.08 | 74.68 | 82.31 | 71.45 | 61.85 |

Table 2: Results on pedestrian. †: our implementation. ∗: re-implemented by [14]. ‡: the model trained and tested on dense point clouds as in Figure 1 (c). ‡ is not a realistic setup, but we provide results to serve as an oracle.

2. Hinton, G., Vinyals, O., Dean, J.: Distilling the knowledge in a neural network. In: NeurIPS Deep Learning and Representation Learning Workshop (2015)
3. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: The International Conference on Learning Representations (ICLR) (2014)
4. Lambert, J., Sener, O., Savarese, S.: Deep learning under privileged information using heteroscedastic dropout. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2018)
5. Lang, A.H., Vora, S., Caesar, H., Zhou, L., Yang, J., Beijbom, O.: Pointpillars: Fast encoders for object detection from point clouds. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2019)
6. Lopez-Paz, D., Bottou, L., Schölkopf, B., Vapnik, V.: Unifying distillation and privileged information. In: International Conference on Learning Representations (ICLR) (2016)
7. Loshchilov, I., Hutter, F.: Sgdr: Stochastic gradient descent with warm restarts. In: The International Conference on Learning Representations (ICLR) (2017)
8. Meyer, G.P., Laddha, A., Kee, E., Vallespi-Gonzalez, C., Wellington, C.K.: Lasernet: An efficient probabilistic 3d object detector for autonomous driving. In: The Conference on Computer Vision and Pattern Recognition (CVPR) (2019)
9. Ngiam, J., Caine, B., Han, W., Yang, B., Chai, Y., Sun, P., Zhou, Y., Yi, X., Alsharif, O., Nguyen, P., Chen, Z., Shlens, J., Vasudevan, V.: Starnet: Targeted computation for object detection in point clouds. ArXiv **abs/1908.11069** (2019)
10. Qi, C.R., Su, H., Mo, K., Guibas, L.J.: Pointnet: Deep learning on point sets for 3d classification and segmentation. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2017)
11. Su, J.C., Maji, S.: Adapting models to signal degradation using distillation. In: British Machine Vision Conference (BMVC) (2017)
12. Sun, P., Kretzschmar, H., Dotiwalla, X., Chouard, A., Patnaik, V., Tak-wing Tsui, P., Guo, J.C.Y., Zhou, Y., Chai, Y., Caine, B., Vasudevan, V., Han, W., Ngiam, J., Zhao, H., Timofeev, A., Ettinger, S., Krivokon, M., Gao, A., Joshi, A., Zhang, Y., Shlens, J., Chen, Z.F., Anguelov, D.: Scalability in perception for autonomous driving: Waymo open dataset. ArXiv **abs/1912.04838** (2019)

13. Vapnik, V., Vashist, A.: A new learning paradigm: Learning using privileged information. Neural Networks (2009)
14. Zhou, Y., Sun, P., Zhang, Y., Anguelov, D., Gao, J., Ouyang, T., Guo, J., Ngiam, J., Vasudevan, V.: End-to-end multi-view fusion for 3d object detection in lidar point clouds. In: The Conference on Robot Learning (CoRL) (2019)