

# SUIT: Learning Significance-guided Information for 3D Temporal Detection

Zheyuan Zhou<sup>1</sup>, Jiachen Lu<sup>1</sup>, Yihan Zeng<sup>2</sup>, Hang Xu<sup>2</sup>, Li Zhang<sup>1†</sup>

**Abstract**—3D object detection from LiDAR point cloud is of critical importance for autonomous driving and robotics. While sequential point cloud has the potential to enhance 3D perception through temporal information, utilizing these temporal features effectively and efficiently remains a challenging problem. Based on the observation that the foreground information is sparsely distributed in LiDAR scenes, we believe sufficient knowledge can be provided by sparse format rather than dense maps. To this end, we propose to learn Significance-guided Information for 3D Temporal detection (SUIT), which simplifies temporal information as sparse features for information fusion across frames. Specifically, we first introduce a significant sampling mechanism that extracts information-rich yet sparse features based on predicted object centroids. On top of that, we present an explicit geometric transformation learning technique, which learns the object-centric transformations among sparse features across frames. We evaluate our method on large-scale nuScenes and Waymo dataset, where our SUIT not only significantly reduces the memory and computation cost of temporal fusion, but also performs well over the state-of-the-art baselines.

## I. INTRODUCTION

3D object detection with LiDAR point cloud is a fundamental task in autonomous driving [1], [2], [3], [4], [5], [6], [7], [8], [9], [10], [11]. Since the LiDAR sensor collects data continuously, the sequential point cloud frames can provide additional information compared to a sparse single frame, which encourages the development of multi-frame detectors [12], [13], [14], [15]. However, the introduction of temporal clues inevitably increases the memory burden and computation cost when storing and fusion the information cross frames. Besides, the temporal information of sequential data is unaligned due to the rapid movement of objects in scenes cross time, which results in invalid fusion. To tackle those issues, it requires a rethink of the temporal information learning for an effective yet efficient multi-frame detector.

Previous works have shown success in improving multi-frame fusion for 3D object detection. Direct point-level fusion methods, such as concatenating multi-frame point clouds, have been explored with positive results [16], [17], [2]. Meanwhile, feature-level fusion methods, which leverage the Bird’s-eye-view feature maps of multiple frames, have also achieved considerable improvements [18], [19], [20], [12], [13], [14]. Recently, those feature-level fusion methods

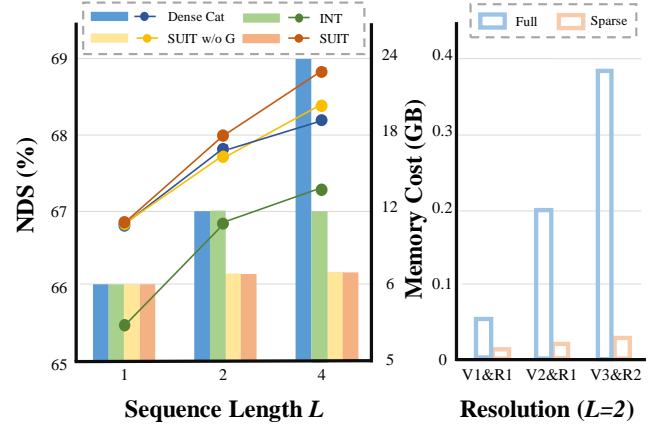


Fig. 1. Comparison of performance and memory cost. Our SUIT not only reduces the memory cost of temporal fusion with sparse features over different sequence lengths and resolutions but also achieves better temporal feature alignment with geometry transformation (G). V1, V2, V3 represent voxel size of  $[0.2m, 0.2m]$ ,  $[0.1m, 0.1m]$ ,  $[0.075m, 0.075m]$ . R1, R2 represent point cloud range of  $[-51.2m, 51.2m]$  and  $[-54.0m, 54.0m]$ .

have focused on learning temporal correlation between dense features using LSTM [21], GRU [19], or Transformer [20].

Despite the performance boost, current methods can not escape from utilizing the whole BEV feature maps at different times, resulting in high memory cost and computation waste by the repeated information across frames. Previous work [22] has proved that using range images only can achieve remarkable performance. This observation inspired us that the whole dense BEV feature map is redundant as historical information since only the foreground areas swept by the LiDAR have clues that can be effectively resolved and exploited by the detector. Furthermore, the point clouds are subject to discontinuities and inconsistencies during the acquisition process, leaving a huge gap for improvement between the alignment of multi-frame features. Though [14] have proposed a way to mitigate the long-sequence memory usage, it still requires a dense feature from the past times. Besides, it adopts simple concatenation to fuse temporal fusion without consideration of the explicit exploration of the unalignment of information between frames. By contrast, our method only exploits the slight sparse feature rather than complete BEV feature maps. Therefore, memory consumption can be drastically whittled down whilst preserving most of the object information. Based on these significant features

<sup>1</sup>Zheyuan Zhou (zheyuanzhou20@fudan.edu.cn), Jiachen Lu and Li Zhang (lizhangfd@fudan.edu.cn) are with Fudan University.

<sup>2</sup>Yihan Zeng and Hang Xu are with Huawei Noah’s Ark Lab.

<sup>†</sup>Li Zhang is the corresponding author with School of Data Science, Fudan University.

with object-dense information, we can explicitly extract their object-centric correlation.

Under these observations, we propose a multi-frame 3D detector: **Significance-guided 3D Temporal Detector (SUIT)**. SUIT introduces two components: 1) *Significant sampling* and 2) *Explicit geometric transformation learning* to better answer two questions raised above: I) **Which** features to attend and II) **How** to aggregate. Based on objects predicted by past frames, a significant sampling technique only saves features of objects information-rich in the past. It promotes the medium between multiple frames which enables only a modicum of features will be maintained as the multi-frame medium in memory but still reserving enough information from past frames. On the other hand, the explicit geometric transformation learning technique learns the object-centric transformation of these sparse features among frames with explicit object transformation supervision. The joint efforts of both components drastically whittle away the computation burden under long-time dependence but also facilitate the learning of geometric transformation.

To prove the feasibility of SUIT, we conduct extensive multi-frame experiments on a large-scale automatic driving dataset nuScenes [16] and Waymo [23]. Figure 1 shows the accuracy v.s. time sequence length and the cost v.s. times sequence length. Our SUIT achieves a larger boost as the time sequence extends whilst keeping a lower memory consumption increment.

In brief, our contributions can be summarized as follows:

- We propose a multi-frame detector framework SUIT based on Significant sampling and Explicit geometric transformation learning.
- The proposed SUIT better answers the two basis: adequate medium to preserve cross-frame information and better geometric transformation understanding.
- Significant sampling sets significant features as medium-saving cross-frame information manifests its advantage of memory consumption and makes it possible for object-centric geometric transformation learning. Explicit geometric transformation learning modules realize the second basis thus promoting multi-frame estimation accuracy.
- Extensive experiments on nuScenes [16] and Waymo [23] prove the effectiveness and efficiency of our SUIT.

## II. RELATED WORK

**Single-Frame 3D Object Detection.** Existing works on 3D object detection with point clouds alone can be mainly divided into two categories: point-based and grid-based. Point-based methods [4], [3], [6], [24], [25], [26], [27] directly learn the 3D representation of the scene by applying PointNet [28], [29] to process the irregular point cloud and generate the final proposal for each point. These methods can preserve accurate geometry information from the raw points. However, due to a large amount of time spent on organizing the irregular points [30], operation on large-scale points will lead to non-trivial costs on computation and memory.

The grid-based schemes [1], [2], [31], [5], [32], [7], [11] first convert sparse point clouds into grid forms using 3D voxels or 2D pillars and then utilize the 3D/2D convolution networks to process the fixed size input. VoxelNet [1] is the first end-to-end voxel-based 3D object detector, which uses PointNet [28], [29] to encode the points within a voxel. SECOND [2] further introduces 3D sparse convolution to tackle a large number of empty voxels in the outdoor scene for acceleration. PointPillars [5] uses pillars rather than voxels to divide the point cloud. Compared to point-based methods, grid-based methods are more efficient in computation and easier to deploy on autonomous driving vehicles. In this work, we choose the popular grid-based CenterPoint [7] as our baseline model.

**Multi-Frame 3D Object Detection.** While handling point cloud sequence, a straightforward method is to concatenate points from the sweep [16], [17], [2], which densifies the point cloud but ignores the temporal correlation between the neighbor sweeps. Although it's simple, the rapid growth of computation and memory usage makes it unable to handle long time series. Instead, some recent approaches [18], [20], [12], [13], [14], [33], try to model the temporal interaction at the feature level. The point cloud from the different frames will be voxelization and encoded independently and merged at the Bird's-eye-view features. However, the point cloud data is read and processed repeatedly, and the sequence length is still bounded by its memory burden. 3DVID [19] introduces a message-passing network and a GRU-based solution to aggregate the spatiotemporal information. TCTR [20] uses a transformer to fuse the middle-level feature, which further explores the spatial, temporal, and channel correlations among different frames. 3D-MAN [12] stores a series of proposals and their corresponding feature from the historical frame in a memory bank. A transformer is adopted to align and encode these features from different perspectives together with the current frame. MPPNet [33] is a two-stage point-based method, which uses a series of proxy points for multi-frame feature encoding and integrating. A series of self and global attention propagate the previous features in the whole trajectory. INT [14] extends the storage of a fixed number of frames into an infinite by an on-stream training and prediction framework. Compared to early works, our method explores the necessity of dense features, as well as the interaction between the neighbor frames.

## III. PRELIMINARIES

**Single-Frame 3D Detector.** A common voxel-based 3D object detection pipeline can be summarized as voxelization, 3D backbone, 2D backbone, and 3D detection heads. Point cloud will be voxelized after voxelization and be represented as Bird's-eye-view features  $F$  after the 3D&2D backbone. Our single-frame 3D detector baseline is implemented by the first stage of a commonly used 3D object detection head [7]. Before heads, feature map  $F_t(I_t) \in \mathbb{R}^{H \times W}$  are produced by 2D backbone with input point cloud  $I_t$ . For the popular adopted center head, the center locations of any probable objects are shown at the peak of a heatmap, as well as the

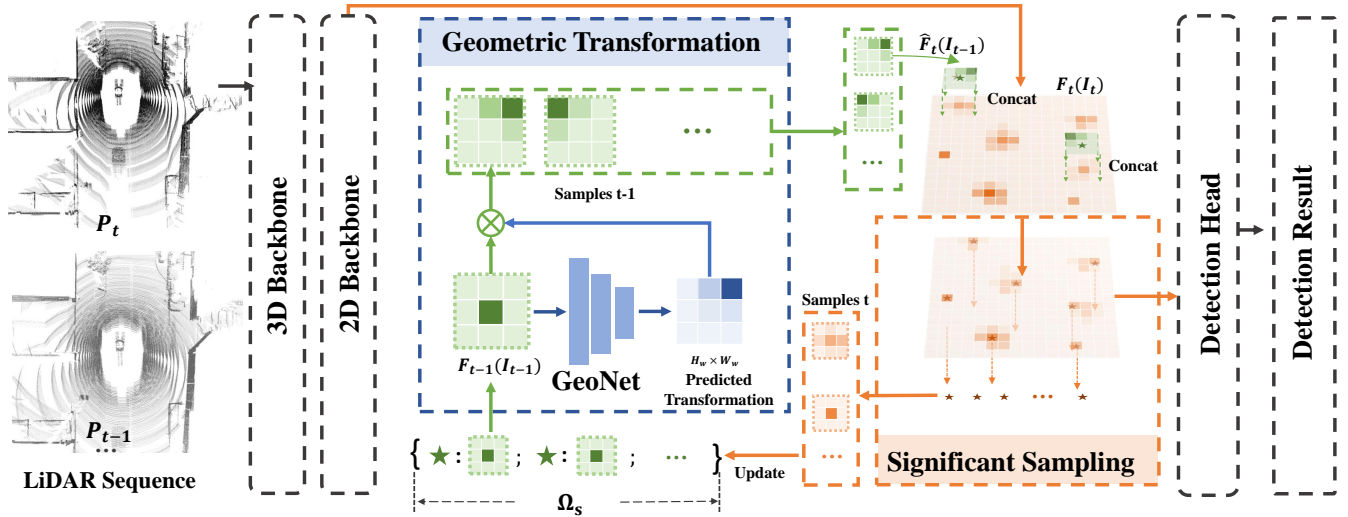


Fig. 2. Architecture overview. Our multi-frame pipeline is positioned between the 2D backbone and detection head, which comprises two key components: geometric transformation and significant sampling. Specifically, geometric transformation aligns features across different time frames, while significant sampling identifies potential targets in the feature map and reduces the impact of irrelevant objects on the detection results.

categories of objects. Fundamentally, we use Bird’s-eye-view coordinate  $\Omega = \{r|r = (x, y)\}$  to represent the location where  $x, y$  are restricted on the BEV feature map  $F$ . Then, we define the random variable  $R_t \in \mathbb{R}^{H \times W}$  to represent the probability of an object being at a given location at current time  $t$ . By knowledge, center heatmap heads produce the following heatmap

$$p(R_t|I_t). \quad (1)$$

**Multi-Frame 3D Detector.** The difference between single-frame and feature-level multi-frame detectors happens after 2D Backbone. Our multi-frame 3D detector shares the same architecture single-frame 3D detector including voxelization, 3D backbone, and detection head. But, our feature-level multi-frame 3D detector aggregates a temporal unified feature from the past features  $F_t(I_t), F_{t-1}(I_{t-1}), \dots$  after 2D backbone and before detection head. Considering the nonalignment among these past features, the aggregated temporal unified feature can be modeled as follows.

Our target is to learn the temporal unified feature  $\hat{F}_t(I_t, I_{t-1}, \dots)$ . The hat notation means the feature map is predicted by temporal transformation. In convenience, we first introduce a 2-frame aggregation  $f_2$ , *i.e.*

$$\hat{F}_t(I_t, I_{t-1}) = f_2 [F_t(I_t), F_{t-1}(I_{t-1})]. \quad (2)$$

Noted that, to align features  $F_t(I_t), F_{t-1}(I_{t-1})$ , we need to first transform geometric information of  $F_{t-1}(I_{t-1})$  to the current frame. To process  $k$  frames simultaneously, we apply a geometric transformation on past features, *i.e.* learning  $\hat{F}_t(I_{t-1})$ . After alignment, we can simply merge two frames

$$\hat{F}_t(I_t, I_{t-1}) = [F_t(I_t) \oplus \hat{F}_t(I_{t-1})]. \quad (3)$$

The problem is then transformed to how to realize the geometric transformation. We rephrase the transformation

with a perspective of Bayesian estimation. We inferred that the transformed past features can be described by a conditional probability-based transformation.

$$\hat{F}_t(I_{t-1}) = \sum_{R_t \in \Omega} p(R_t|I_{t-1})F_{t-1}(I_{t-1}), \quad (4)$$

We expand  $p(R_t|I_{t-1})$  with empirical Bayesian

$$\begin{aligned} p(R_t|I_{t-1}) &= \int_{\Omega} p(R_t|R_{t-1}, I_{t-1})p(R_{t-1}|I_{t-1}) dR_{t-1} \\ &= \sum_{R_{t-1} \in \Omega} p(R_t|R_{t-1}, I_{t-1})p(R_{t-1}|I_{t-1}). \end{aligned} \quad (5)$$

From the above equation, we have the following observation: (i)  $p(R_{t-1}|I_{t-1})$  is the product of the heatmap head of the single-frame detector according to Equation 1. (ii) Summation overall space  $\Omega$  is memory consuming. Therefore, a small subspace that can also approximate the equation will be preferred.

## IV. METHOD

Based on these observations, we propose our multi-frame detector SUIT. Section IV-A illustrates the significant sampling approach for the subspace  $\Omega_s$ , which produce and store a series of sparse BEV features with rich information. Section IV-B shows the geometric transformation learning that aggregates features from different times. Finally, we illustrate the overall multi-frame detector in Section IV-C.

### A. Significant sampling

The output of the heatmap head in our single-frame detector, denoted by  $p(R_{t-1}|I_{t-1})$ , is a sparse distribution. This is because the heatmap only indicates a few locations where objects may be present. Thus, a subspace  $\Omega_s \subset \Omega$  with a significant distribution can be identified, where  $|\Omega_s| \ll$

$|\Omega|$ . To save computation resources while approximating the original target, we propose a significant sampling process consisting of three steps: coarse sampling, refined sampling, and relaxation.

**Coarse Sampling.** We use a threshold  $\alpha$  to select  $R_{t-1}$  which  $p(R_{t-1}|I_{t-1}) < \alpha$  in Equation 5, where  $\alpha$  is a small real number called *significant threshold*.

**Refined Sampling.** Based on the coarse sampling, we have sampled all the points of interest but many of them actually belong to the same Gaussian peak. To mitigate this issue, we propose a refinement stage based on Non-maximum Suppression (NMS). Specifically, we leverage the location and scale information obtained from the Regression heads of the single frame detection to filter out overlapping boxes. By applying NMS, we ensure that only the most promising and representative boxes are retained while discarding redundant ones."

**Relaxation.** After refined sampling, we find that simply using the feature  $F_{t-1}(I_{t-1})$  extracted at the sampled points loses too much information. Therefore, we make a neighborhood relaxation  $\mathcal{B}(r_{t-1})$  around the sample points  $r_{t-1}$ . We try three relaxation methods – rectangle window, circular window, and square window. *Rectangle window* takes points inside the box predicted by Regression heads of single frame detector. *Circular window* takes points from the circle with radius  $\rho$  around  $r_{t-1}$ . *Square window* takes points inside the  $H_s \times W_s$  square centered at  $r_{t-1}$ . our significant sampling sample a minority of locations  $\Omega_s \subset \Omega$ . Noted that the scale of  $\Omega_s$  ( $10^1$ ) is greatly smaller than that of  $\Omega$  ( $10^4$ ). Now, we can simplify the Equation 5 by sampling only locations with significant influence, *i.e.*

$$\begin{aligned} p(R_t|I_{t-1}) \\ \approx \sum_{R_{t-1} \in \Omega_s} p(R_t|R_{t-1}, I_{t-1})p(R_{t-1}|I_{t-1}). \end{aligned} \quad (6)$$

Consequently, combining Equation 4 and 6, we get

$$\begin{aligned} \hat{F}_t(I_{t-1}) \\ \approx \sum_{R_t \in \Omega} \left[ \sum_{R_{t-1} \in \Omega_s} p(R_t|R_{t-1}, I_{t-1})p(R_{t-1}|I_{t-1}) \right] F_{t-1}(I_{t-1}) \\ = \sum_{R_t \in \Omega} p(R_t|R_{t-1}, I_{t-1}) \left[ \sum_{R_{t-1} \in \Omega_s} p(R_{t-1}|I_{t-1})F_{t-1}(I_{t-1}) \right]. \end{aligned} \quad (7)$$

Therefore, only Feature  $F_{t-1}(I_{t-1})$  sampled in  $\Omega_s$  need to be stored for next frame inference which shows our memory efficiency.

### B. Geometric transformation learning

Geometric transformation can be divided into two guises: ego-car transformation and object-centric transformation. For ego-car transformation, we directly transform the feature map of the last frame to the current frame BEV coordinate so that  $p(R_{t-1}|I_{t-1})$  will be directly predicted under the current BEV coordinate. This transformation can be carried

out at no additional cost. On the other hand, the object-centric transformation is more complex and requires explicit consideration in our method. We introduce our object-centric transformation learning in two parts: theory and implementation.

**Theory** The object-centric transformation is based on the velocities of objects. To acquire ground-truth velocities of each sample from  $\Omega_s$  of past frames, we first assign each sample with the specific ground-truth object. The velocity of the successfully assigned sample is equal to the velocity of the assigned object while the velocity of the unassigned sample is equal to 0. We denote the ground-truth velocity for  $r_{t-1}$  as  $v_{t-1}^{(r_{t-1})} \in \mathbb{R}^2$ , then the ground-truth distribution for  $p$  can be derived as

$$p(\hat{R}_t|r_{t-1}, I_{t-1}) = \begin{cases} 1, & \hat{R}_t = r_{t-1} + v_{t-1}^{(r_{t-1})} \cdot \tau \\ 0, & \text{otherwise.} \end{cases} \quad (8)$$

where  $\tau$  is the time gap between frames. However, learning the probability across  $\hat{R}_t$  on all the point cloud ranges is tedious and difficult. To tackle the issue, we make the following local approximation. We limit the space of  $R_t$  inside the neighbor of  $r_{t-1}$ , *i.e.* a  $H_w \times W_w$  window  $\mathcal{B}(r_{t-1})$ . We define a probability  $q(\hat{R}_t|r_{t-1}, I_{t-1})$  to represent the ground-truth geometric transformation probability, where  $\hat{R}_t$  is limited inside the window  $\mathcal{B}(r_{t-1})$ . Assume that, the choice of  $H_w, W_w$  can adequately take  $v_{t-1}^{(r_{t-1})} \cdot \tau \in \mathcal{B}(r_{t-1})$  for all  $r_{t-1} \in \Omega_s$ , we have

$$q(\hat{R}_t|r_{t-1}, I_{t-1}) = \begin{cases} 1, & \hat{R}_t = v_{t-1}^{(r_{t-1})} \cdot \tau \\ 0, & \text{otherwise.} \end{cases} \quad (9)$$

Then we define the  $q_\phi(R_t|r_{t-1}, I_{t-1})$  to represent the predicted geometric transformation probability, where  $R_t$  is also limited inside the window  $\mathcal{B}(r_{t-1})$ . During optimization, we minimize the divergence between probability, *i.e.*

$$\min_{\phi} \mathcal{D}(q(\hat{R}_t|r_{t-1}, I_{t-1}), q_\phi(R_t|r_{t-1}, I_{t-1})) \quad (10)$$

Finally, the global predicted geometric transformation probability can be approximated with local probability by

$$p(R_t|r_{t-1}, I_{t-1}) = \begin{cases} q(R_t|r_{t-1}, I_{t-1}), & R_t \in \mathcal{B}(r_{t-1}) \\ 0, & \text{otherwise.} \end{cases} \quad (11)$$

**Implementation** We predict geometric transformation probability of position  $r_{t-1} \in \Omega_s$  from the feature  $F_{t-1}(I_{t-1})$  sampled around  $r_{t-1}$ . Given relaxation form Section IV-A, we can acquire a  $H_s \times W_s \times C$  local feature map. Then we use a simple three-layer CNN called **GeoNet** on the local feature map to produce the geometric transformation probability  $1 \times 1 \times (H_w W_w)$ . During the training phase, we use Cross Entropy Loss to optimize the distribution.

### C. Detection Head

Finally, Equation 3 gives the aligned aggregation. For a longer sequence, we use induction. Taking  $t+1$  for example,



Method	Backbone	Information	mAP	NDS	Car	Truck	C.V.	Bus	Trailer	Barrier	Motor.	Bicycle	Ped.	T.C.
PointPillars [5]	Pillar	L	30.5	45.3	68.4	23.0	4.1	28.2	23.4	38.9	27.4	1.1	59.7	30.8
WYSIWYG [32]	Pillar	L	35.0	41.9	79.1	30.4	7.1	46.6	40.1	34.7	18.2	0.1	65.0	28.8
3DVID [19]	Pillar	L+T	45.4	53.1	79.7	33.6	18.1	47.1	43.0	48.8	40.7	7.9	76.5	58.8
TCTR [20]	Pillar	L+T	50.5	-	83.2	51.5	15.6	<b>63.7</b>	33.0	53.8	54.0	22.6	74.9	52.5
TransPillars [13]	Pillar	L+T	52.3	-	84.0	<b>52.4</b>	18.9	62.0	34.3	55.1	55.2	27.6	77.9	55.4
<b>SUIT(Ours)</b>	Pillar	L+T	<b>57.1</b>	<b>65.5</b>	<b>84.6</b>	48.8	<b>23.9</b>	57.5	<b>53.8</b>	<b>63.2</b>	<b>58.6</b>	<b>30.5</b>	<b>80.3</b>	<b>69.6</b>
CBGS [31]	Voxel	L	52.8	63.3	81.1	48.5	10.5	54.9	42.9	65.7	51.5	22.3	80.1	70.9
CenterPoint [7]	Voxel	L	58.0	65.5	84.6	51.0	17.5	<b>60.2</b>	53.2	70.9	53.7	28.7	83.4	76.7
<b>SUIT(Ours)</b>	Voxel	L+T	<b>62.8</b>	<b>68.9</b>	<b>85.9</b>	<b>53.7</b>	<b>26.4</b>	59.0	<b>54.9</b>	<b>71.6</b>	<b>68.9</b>	<b>42.9</b>	<b>85.6</b>	<b>79.7</b>

TABLE I

PERFORMANCE COMPARISONS ON THE nuSCENES TEST SET. WE REPORT THE OVERAL MAP, NDS, AND MAP FOR EACH DETECTION CATEGORY, WHERE “L” DENOTE LiDAR MODALITY, AND T DENOTES MULTI-FRAME TEMPORAL INPUT. “-” REPRESENTS THE UNKNOWN INFORMATION. WE HIGHLIGHT THE BEST RESULTS WITH PILLAR-BASED BACKBONE IN **BLUE** AND WITH VOXEL-BASED BACKBONE IN **RED**.

Method	Frames	mAP	NDS	Latency(ms)	Memory(GB)
Centerpoint [7]	1	59.6	66.8	<b>78.6</b>	<b>5.9</b>
INT [14]	1	58.5	65.5	84.1	12.1
INT [14]	2	60.9	66.9	84.1	12.1
INT [14]	10	61.8	67.3	84.1	12.1
<b>SUIT(Ours)</b>	2	<b>63.1</b>	<b>68.4</b>	107.5	6.2
<b>SUIT(Ours)</b>	10	<b>63.7</b>	<b>68.6</b>	107.5	6.2

TABLE II

PERFORMANCE COMPARISONS ON THE nuSCENES VALIDATION SET. WE REPORT THE OVERAL MAP, NDS. ALL MODELS ARE BUILT ON THE VOXELNET BACKBONE. WE HIGHLIGHT THE BEST RESULTS IN **BLOD**.

we can replace  $F_t(I_t)$  with  $\hat{F}_t(I_{t-1})$ , *i.e.*

$$\hat{F}_{t+1}(I_{t+1}, I_t) = f_2 \left[ F_{t+1}(I_{t+1}), \hat{F}_t(I_{t-1}) \right]. \quad (12)$$

After the acquirement of the temporal unified feature, the following detection heads – heatmap head and regression head are just the same as the single-frame 3D detector.

## V. EXPERIMENTS

We evaluate our method on two widely used large-scale datasets including nuScenes [16] and Waymo [23]. We conduct extensive ablation studies to validate our design choices.

### A. Experimental Setup

**Dataset.** The nuScenes [16] is a large-scale dataset for autonomous driving. It contains 1,000 driving sequences, which are divided into 700, 150, and 150 for training, validation, and testing, respectively. Each sequence lasts about 20 seconds, with a LiDAR sampling frequency of 20HZ. nuScenes provides annotations every 0.5 seconds including 10 classes with a long-tail distribution. The Waymo Open Dataset [23] is consisting of 798 training scenes and 202 validation scenes. Data collection and 3D annotation are both at 10 Hz frequency. Notably, nuScenes annotates 10 classes for 3D detection task, while Waymo provides annotations of 3 classes for its 3D detection benchmark.

**Experimental Details.** We follow CenterPoint [7] to set our backbone network. For nuScenes data, we set the range of the point cloud space as  $[-54.0m, 54.0m]$  for the  $X$  and  $Y$  axes, and  $[-5m, 3m]$  for the  $Z$  axis. The size of each voxel

Methods	Frames	VEH	PED	CYC	mAPH	Latency(ms)
CenterPoint [7]	1	66.2	62.6	67.6	65.5	71.7
CenterPoint [7]	2	67.3	67.5	69.9	68.2	90.9
3D-MAN [12]	15	67.1	59.0	-	-	-
RSN [22]	3	69.1	-	-	-	<b>67.5</b>
INT [34]	2	69.4	69.1	<b>72.6</b>	70.3	74.0
<b>SUIT</b>	2	<b>70.0</b>	<b>70.3</b>	72.4	<b>70.9</b>	79.3

TABLE III

PERFORMANCE COMPARISONS ON THE WAYMO VALIDATION SET. WE REPORT THE LEVEL\_2 MAPH(%) FOR EACH CATEGORY. ALL MODELS ARE BUILT ON THE VOXELNET BACKBONE. WE HIGHLIGHT THE BEST RESULTS IN **BLOD**.

is set to  $(0.075m, 0.075m, 0.2m)$ . For Waymo data, we set the range of the point cloud space as  $[-75.2m, 75.2m]$  for the  $X$  and  $Y$  axes, and  $[-2m, 4m]$  for the  $Z$  axis. The size of each voxel is set to  $(0.1m, 0.1m, 0.15m)$ . We trained our model using Adam optimizer with one cycle learning rate policy, weight decay 0.01, and momentum 0.85 to 0.95 on 8 Tesla V100 GPUs. The total batch size of our model is 32. We use the official CenterPoint [7] checkpoints as our pretrained model and fine-tune 6 epochs. During training, we use random flipping, global scaling, global rotation, and global translation as our data augmentation.

**Metrics.** We use mean Average Precision (mAP) [35] and NuScenes Detection Score (NDS) [16] to evaluate our performance. The mAP stands for the average of precision on 0.5m, 1m, 2m, 4m four different distance thresholds. NDS is a particular metric on nuScenes [16] dataset. It is the average weighted summation of mAP and True Positive (TP) metrics. TP metrics average five individual metrics: translation (ATE), velocity (AVE), scale (ASE), orientation (AOE), and attribute (AAE) errors. The weighted average is calculated by:  $NDS = \frac{1}{10} [5mAP + \sum_{mTP \in TP} (1 - \min(1, mTP))]$ . We follow the official evaluation metrics mAP and mAPH (mAP weighted by heading) on Waymo Open Dataset. The metrics are further split into two difficulty levels according to the point numbers in GT boxes: LEVEL\_1 ( $>5$ ) and LEVEL\_2 ( $\geq 1$ ).

### B. Main Results

**nuScenes Results.** Table I shows our comparison with state-of-the-art single-frame and multi-frame methods on

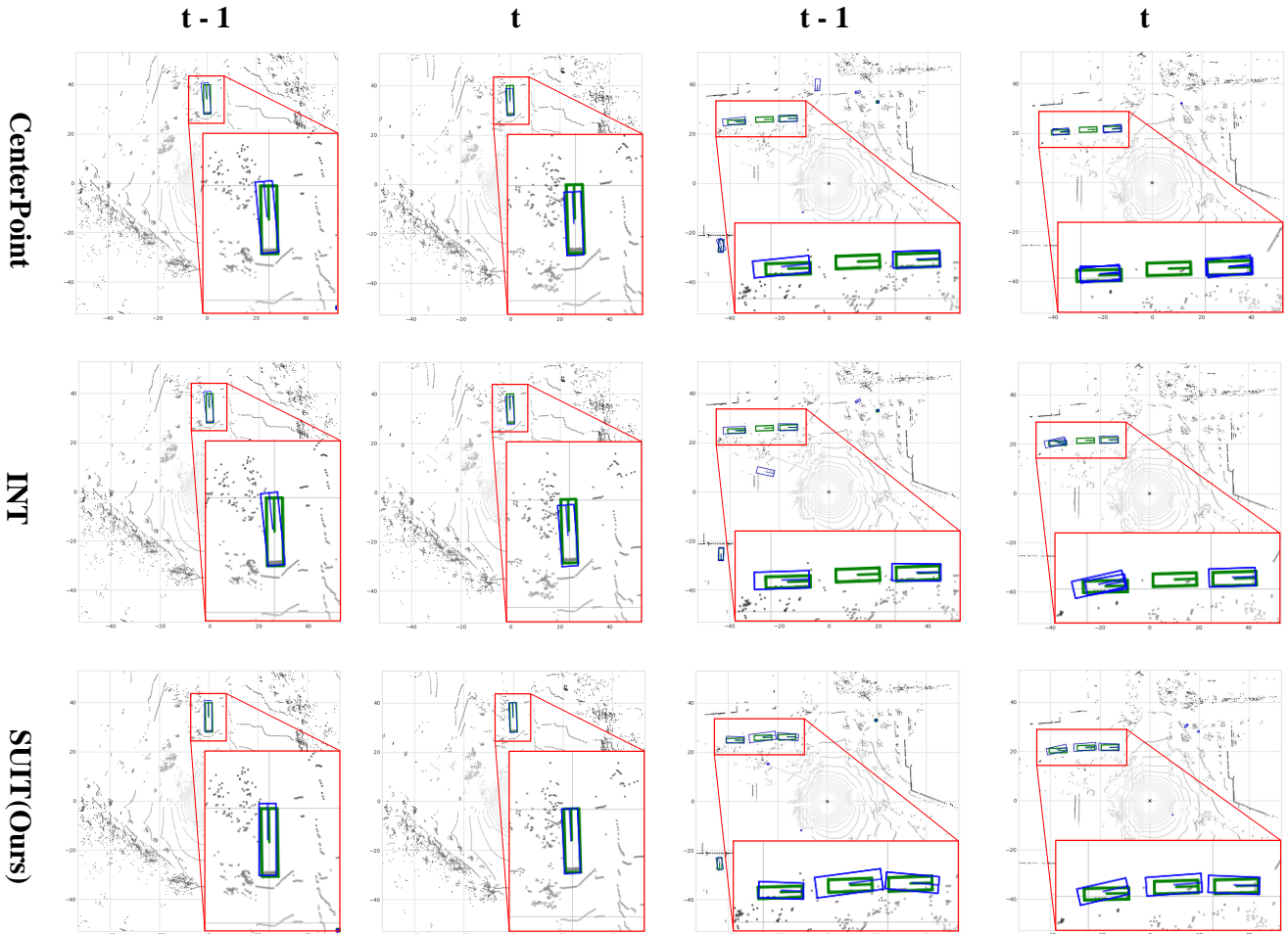


Fig. 3. Qualitative comparison of temporal 3D detector between INT [14] (*middle*) and ours (*bottom*) on nuScenes dataset. We also provide the single-frame baseline CenterPoint [7](*top*) to show its detection result. Boxes in **blue** are the predicted results while the **green** boxes are the ground truth.

nuScenes [16] test split. Since nuScenes provides 10 samples between annotated frames, we follow the common technique in INT [14] which aggregates the point cloud of those interpolated samples to represent one frame. We first compare our multi-frame improvement with the single-frame baseline Centerpoint [7]. According to the 8<sup>th</sup> and 9<sup>th</sup> row, our multi-frame mechanism bring 4.8 point improvement on mAP and 3.4 improvement on NDS. To compare with other state-of-the-art methods, we mainly select pillar-based methods for the lack of voxel-based multi-frame results on nuScenes [16] dataset. Comparing the 5<sup>th</sup> and 6<sup>th</sup> row, we surpass the former multi-frame state-of-the-art TransPillars [13] by 4.8 point on mAP. 3DVID [19], TCTR [20] and TransPillars [13] utilize implicit GRU or Transformer on dense features to learn a geometric transformation. Our higher performance proves the importance of explicit supervision of geometric transformation.

Table II shows our comparison with state-of-the-art multi-frame methods on nuScenes [16] validation split. On the 2-frame paradigm, we outperform the former state-of-the-art INT [14] with 2.2 point on mAP, 1.5 point on NDS under the

voxel-based backbone. On a longer sequence, our 10-frame detector surpasses 10-frame detector of INT [14] by 2.8 point on mAP under the voxel-based backbone. Our SUIT achieves the best mAP and the best memory cost among multi-frame models. As the frame length increases, SUIT still runs with low latency and memory cost, benefiting from the sparse information. INT [14] use a loose uniform distribution assumption when handling object-centric geometric transformation. The loose assumption, according to the comparison, harms temporal information learning. Therefore, our explicit geometric transformation learning drastically surpasses the INT [14].

**Waymo Results.** We also make comparisons on the Waymo dataset as in Table III. Note that nuScenes and Waymo differ a lot in terms of sensor configurations and evaluation classes, the detection model is hard to be reused among datasets, which is a consensus in previous work. Thus we train a detector on the Waymo dataset with the same model architecture on the nuScenes dataset. The results in Table III show that our model consistently outperforms previous methods and boosts the detector by on the three challenging classes

Top K	$\alpha$	Refined	NDS	mAP
500	-	-	67.7	62.6
200	-	-	68.1	62.6
200	0.1	-	68.2	62.7
200	0.2	-	68.2	62.7
200	0.1	✓	<b>68.4</b>	<b>63.1</b>

TABLE IV

ABLATIONS ON THE SAMPLING SCHEMES. WE HIGHLIGHT THE BEST RESULTS IN **BLOD**.

Windows	Relaxation	NDS	mAP
Rectangle	-	65.9	61.3
Circular	$\rho=2$	68.0	62.9
Square	$H_s = W_s=3$	68.4	<b>63.1</b>
Square	$H_s = W_s=5$	67.5	62.9
Square	$H_s = W_s=7$	<b>68.5</b>	62.9

TABLE V

ABLATIONS ON THE SAMPLING RELAXATION. WE HIGHLIGHT THE BEST RESULTS IN **BLOD**.

for the Waymo evaluation benchmark. Experimental results demonstrate that our method generalizes and scales well on both datasets with state-of-the-art performance.

### C. Qualitative Results.

**Multi-frame detection results.** We present visualizations of multi-frame detection results in Figure 3. We compare single-frame detection baseline Centerpoint [7] and the previous state-of-the-art multi-frame detector INT [14]. Compared with the single-frame detector, we can refine the current frame detection results by the last frame. Box circled in red notation is refined based on the last frames by our method, but the single-frame cannot build the temporal refinement. Compared with the previous, we predict a more accurate result under a lower cost on memory.

### D. Ablation Studies

**Analysis on sample stages.** Table IV shows the investigation of two stages of significant sampling – coarse sampling and refined sampling. First, we do not use coarse sampling and refined sampling but only leave samples with top-K probability. Comparing the 1<sup>st</sup> and 2<sup>nd</sup> rows, we observe that leaving the top 500 samples performs worse than the top 200 although more samples reserve more information during sampling. It shows that introducing more noise object-centric samples disturb the following object-centric geometric transformation learning. On the other hand, when we only use coarse sampling, we tune the significant threshold  $\alpha$ . Comparing the 3<sup>rd</sup> and 4<sup>th</sup> rows, we observe that a larger threshold, filtering more samples, however, leads to poorer performance. It illustrates that, on the other hand, too less samples can omit too many significant samples thus weakening our model. Finally, we compare the importance of refined sampling between 3<sup>rd</sup> and 5<sup>th</sup> rows. With refined NMS operation, we can whittle down noise samples thus helping the geometric transformation learning.

$F_{t-j}(I_{t-j})$	$p_\theta(R_{t-j} I_{t-j})$	$p_\phi(\tilde{R}_t R_{t-j}, I_{t-j})$	NDS	mAP
✗	✗	✗	66.5	60.9
✓	✗	✗	68.1	62.6
✓	✓	✗	68.1	62.6
✓	✓	✓	<b>68.4</b>	<b>63.1</b>

TABLE VI

ABLATIONS ON THE GEOGRAPHIC TRANSFORMATION COMPONENTS. WE HIGHLIGHT THE BEST RESULTS IN **BLOD**.

	Components	NDS	mAP
0	Baseline	66.8	59.6
1	+dense BEV	67.9	60.8
	+GeoNet	68.2	62.9
2	+sparse BEV	68.1	62.6
	+speed	67.8	62.2
3	+sparse BEV	68.1	62.6
	+GeoNet	<b>68.4</b>	<b>63.1</b>

TABLE VII

ABLATION ON THE GEO NET. WE HIGHLIGHT THE BEST RESULTS IN **BLOD**.

**Analysis on relaxation.** Table V studies three relaxation methods. We have the following observations. (i) Rectangle window in the 1<sup>st</sup> achieves the worst performance. The potential reason is the error brought by Regression heads since it uses the predicted box as a neighbor. (ii) Circular windows on 2<sup>nd</sup> and 3<sup>rd</sup> rows achieve better than rectangle windows and are similar to square windows. The circular windows have an irregular boundary, thus less efficient than square windows (iii) Square windows perform the best on the mAP metric. Also, the regular boundary is more efficient when accessing memory-storing features. The different window sizes can achieve similar performance but considering the efficiency of small window sizes, we finally choose square window relaxation with the window size equal to 3.

**Analysis on components.** Table VI shows the ablation studies of components in Equation 7. Based on the 1<sup>st</sup> row, when feature  $F_{t-j}(I_{t-j})$  attending, we find the importance of saving features from the last frames. Based on the 2<sup>nd</sup> row, when the prior  $p(R_{t-j}|I_{t-j})$  attending, we observe a tiny promotion. It shows that a uniform prior distribution is strong enough. Comparing the 3<sup>rd</sup> and 4<sup>th</sup> row, a huge difference is made by geometric transformation. By considering  $p_\phi(\tilde{R}_t|R_{t-j}, I_{t-j})$ , we find a 0.6% on mAP, which proves the significance of geometric transformation.

**Analysis on GeoNet** Table VII presents the ablation study results that demonstrate the effectiveness of the proposed GeoNet. In the absence of the proposed significant sampling process, our GeoNet aligns all possible targets on the entire dense BEV feature map, resulting in a modest 2.1 point increase in mAP. A comparison of the last two rows indicates that directly utilizing predicted velocities to align objects fails to enhance the detector’s performance. This observation suggests that inaccurate velocity predictions can adversely affect the performance of multi-frame detectors.

Sampling relaxation	Local window	NDS	mAP
$H_s = W_s = 3$	$H_w = W_w = 3$	68.2	62.59
$H_s = W_s = 3$	$H_w = W_w = 7$	<b>68.4</b>	<b>63.1</b>
$H_s = W_s = 3$	$H_w = W_w = 15$	68.4	63.0
$H_s = W_s = 7$	$H_w = W_w = 7$	68.3	62.8
$H_s = W_s = 7$	$H_w = W_w = 15$	68.3	62.9

TABLE VIII

ABLATIONS ON THE GEOGRAPHIC LEARNING LOCAL WINDOWS. ALSO, WE CONSIDER THE INFLUENCE OF SAMPLING RELAXATION. WE HIGHLIGHT THE BEST RESULTS IN **BLOD**.

Table VIII shows additional analyses on GeoNet. We first study the influence of local window size. Comparing the first three rows, we observe that both too-large and too-small local windows downgrade the geometric learning performance. The small local window can miss the geometric transformation of long-distance (e.g. fast vehicles) but is easy to train. Contrariwise, the large local window can cover larger geometric transformations but is harder to train.  $H_w = W_w = 7$  achieves an equipose between large and small local windows. The study of sampling relaxation receives a similar observation as that in Section V-D. There is no significant influence on relaxation size, but smaller sizes are more efficient.

## VI. CONCLUSION

In this paper, we present a single-stage 3D detection framework SUIT, which aims to exploit the potential of temporal information. Particularly, we extract the slight sparse features from redundant temporal information, which are then fused across frames to boost detection accuracy under low memory cost. Furthermore, we introduce an explicit geometric transformation that learns the relationships throughout the trajectory of an object to enhance the alignment. With the proposed end-to-end detector, SUIT achieves significant improvements over strong baselines on the challenging nuScenes and Waymo benchmark datasets. It is noteworthy that our model is a general multi-frame approach that can be extended to other related tasks such as segmentation and tracking.

## REFERENCES

- [1] Y. Zhou and O. Tuzel, "Voxelnet: End-to-end learning for point cloud based 3d object detection," in *CVPR*, 2018.
- [2] Y. Yan, Y. Mao, and B. Li, "Second: Sparsely embedded convolutional detection," *Sensors*, 2018.
- [3] Z. Yang, Y. Sun, S. Liu, X. Shen, and J. Jia, "Std: Sparse-to-dense 3D object detector for point cloud," in *ICCV*, 2019.
- [4] S. Shi, X. Wang, and H. Li, "Pointrenn: 3D object proposal generation and detection from point cloud," in *CVPR*, 2019.
- [5] A. H. Lang, S. Vora, H. Caesar, L. Zhou, J. Yang, and O. Beijbom, "Pointpillars: Fast encoders for object detection from point clouds," in *CVPR*, 2019.
- [6] Z. Yang, Y. Sun, S. Liu, and J. Jia, "3DSSD: Point-based 3D single stage object detector," in *CVPR*, 2020.
- [7] T. Yin, X. Zhou, and P. Krähenbühl, "Center-based 3d object detection and tracking," in *CVPR*, 2021.
- [8] S. Shi, C. Guo, L. Jiang, Z. Wang, J. Shi, X. Wang, and H. Li, "PV-RCNN: Point-voxel feature set abstraction for 3d object detection," in *CVPR*, 2020.

- [9] J. Deng, S. Shi, P. Li, W. Zhou, Y. Zhang, and H. Li, "Voxel r-cnn: Towards high performance voxel-based 3d object detection," in *AAAI*, 2021.
- [10] L. Fan, X. Xiong, F. Wang, N. Wang, and Z. Zhang, "Rangedet: In defense of range view for lidar-based 3d object detection," in *ICCV*, 2021.
- [11] Y. Hu, Z. Ding, R. Ge, W. Shao, L. Huang, K. Li, and Q. Liu, "Afdetv2: Rethinking the necessity of the second stage for object detection from point clouds," in *AAAI*, 2022.
- [12] Z. Yang, Y. Zhou, Z. Chen, and J. Ngiam, "3d-man: 3d multi-frame attention network for object detection," in *CVPR*, 2021.
- [13] Z. Luo, G. Zhang, C. Zhou, T. Liu, S. Lu, and L. Pan, "Transpillars: Coarse-to-fine aggregation for multi-frame 3d object detection," *arXiv preprint*, 2022.
- [14] J. Xu, Z. Miao, D. Zhang, H. Pan, K. Liu, P. Hao, J. Zhu, Z. Sun, H. Li, and X. Zhan, "Int: Towards infinite-frames 3d detection with an efficient framework," in *ECCV*, 2022.
- [15] Z. Zhou, X. Zhao, Y. Wang, P. Wang, and H. Foroosh, "Centerformer: Center-based transformer for 3d object detection," in *ECCV*, 2022.
- [16] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom, "nusenes: A multimodal dataset for autonomous driving," in *CVPR*, 2020.
- [17] B. Yang, M. Bai, M. Liang, W. Zeng, and R. Urtasun, "Auto4d: Learning to label 4d objects from sequential point clouds," *arXiv preprint*, 2021.
- [18] R. Huang, W. Zhang, A. Kundu, C. Pantofaru, D. A. Ross, T. Funkhouser, and A. Fathi, "An lstm approach to temporal 3d object detection in lidar point clouds," in *ECCV*, 2020.
- [19] J. Yin, J. Shen, C. Guan, D. Zhou, and R. Yang, "Lidar-based online 3d video object detection with graph-based message passing and spatiotemporal transformer attention," in *CVPR*, 2020.
- [20] Z. Yuan, X. Song, L. Bai, Z. Wang, and W. Ouyang, "Temporal-channel transformer for 3d lidar-based video object detection for autonomous driving," *TCSVT*, 2021.
- [21] R. Huang, W. Zhang, A. Kundu, C. Pantofaru, D. A. Ross, T. Funkhouser, and A. Fathi, "An lstm approach to temporal 3d object detection in lidar point clouds," in *ECCV*, 2020.
- [22] P. Sun, W. Wang, Y. Chai, G. Elsayed, A. Bewley, X. Zhang, C. Sminchisescu, and D. Anguelov, "Rsn: Range sparse net for efficient, accurate lidar 3d object detection," in *CVPR*, 2021.
- [23] P. Sun, H. Kretzschmar, X. Dotiwalla, A. Chouard, V. Patnaik, P. Tsui, J. Guo, Y. Zhou, Y. Chai, B. Caine, *et al.*, "Scalability in perception for autonomous driving: Waymo open dataset," in *CVPR*, 2020.
- [24] W. Shi and R. Rajkumar, "Point-gnn: Graph neural network for 3d object detection in a point cloud," in *CVPR*, 2020.
- [25] D. Zhou, J. Fang, X. Song, L. Liu, J. Yin, Y. Dai, H. Li, and R. Yang, "Joint 3d instance segmentation and object detection for autonomous driving," in *CVPR*, 2020.
- [26] W. Zheng, W. Tang, L. Jiang, and C.-W. Fu, "Se-ssd: Self-ensembling single-stage object detector from point cloud," in *CVPR*, 2021.
- [27] H.-X. Yu, J. Wu, and L. Yi, "Rotationally equivariant 3d object detection," in *CVPR*, 2022.
- [28] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "Pointnet: Deep learning on point sets for 3D classification and segmentation," in *CVPR*, 2017.
- [29] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "Pointnet++: Deep hierarchical feature learning on point sets in a metric space," in *NeurIPS*, 2017.
- [30] Z. Liu, H. Tang, Y. Lin, and S. Han, "Point-voxel cnn for efficient 3d deep learning," in *NeurIPS*, 2019.
- [31] B. Zhu, Z. Jiang, X. Zhou, Z. Li, and G. Yu, "Class-balanced grouping and sampling for point cloud 3d object detection," *arXiv preprint*, 2019.
- [32] P. Hu, J. Ziglar, D. Held, and D. Ramanan, "What you see is what you get: Exploiting visibility for 3d object detection," in *CVPR*, 2020.
- [33] X. Chen, S. Shi, B. Zhu, K. C. Cheung, H. Xu, and H. Li, "Mppnet: Multi-frame feature intertwining with proxy points for 3d temporal object detection," in *ECCV*, 2022.
- [34] C.-Y. Wu and P. Krahenbuhl, "Towards long-form video understanding," in *CVPR*, 2021.
- [35] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (voc) challenge," *IJCV*, 2010.