# Efficient Transformer-based 3D Object Detection with Dynamic Token Halting

Mao Ye[*1,2], Gregory P. Meyer[†2], Yuning Chai[2], and Qiang Liu[1]

[1]The University of Texas at Austin
[2]Cruise LLC

## Abstract

*Balancing efficiency and accuracy is a long-standing problem for deploying deep learning models. The trade-off is even more important for real-time safety-critical systems like autonomous vehicles. In this paper, we propose an effective approach for accelerating transformer-based 3D object detectors by dynamically halting tokens at different layers depending on their contribution to the detection task. Although halting a token is a non-differentiable operation, our method allows for differentiable end-to-end learning by leveraging an equivalent differentiable forward-pass. Furthermore, our framework allows halted tokens to be reused to inform the model's predictions through a straightforward token recycling mechanism. Our method significantly improves the Pareto frontier of efficiency versus accuracy when compared with the existing approaches. By halting tokens and increasing model capacity, we are able to improve the baseline model's performance without increasing the model's latency on the Waymo Open Dataset.*

## 1. Introduction

Recent progress has shown the potential of applying the transformer architecture, which has been widely used by the natural language processing community [72, 25], to computer vision tasks [11, 42]. Transformers already meet or exceed the performance of convolutional neural networks. However, transformer architectures often suffer from high latency, which is crucial for real-time safety-critical or edge-computing applications.

This paper explores how to speed-up transformer-based 3D object detection. Our method is inspired by network pruning, which increases efficiency by removing less important parts of the model. However, instead of pruning neurons like in [41, 43, 84], our approach prunes or halts tokens. We want to reduce superfluous tokens because the computational complexity of the transformer's attention mechanism increases with the number of tokens. Unlike network pruning, deciding whether or not to halt a token needs to depend on the input, since the importance of a token will depend on the particular example.

There are several recent works that attempt to dynamically halt tokens throughout a vision transformer [58, 87, 54, 16, 31, 70]. However, these works consider the task of image classification, while we focus on 3D object detection. Going from image classification to 3D object detection has its challenges, but also its benefits. A challenge that arises with object detection is that any token could contain an object; therefore, we require a method that can detect objects from all tokens regardless of whether or not they were halted. This issue does not occur for image classification with vision transformers. For image classification, an artificial token is often added to classify the image, and this token is prevented from being halted. Since the halting of a token is a non-differentiable operation, a new design of the computation graph is needed to define pseudo-gradients that enable the back-propagation during training. A benefit with object detection is that the labels contain not only the object's classification but also its location, and our approach is able to leverage the localization of the objects to help the model learn which tokens are more important than others.

Our method is designed for 3D object detection for autonomous driving. For safety-critical tasks, it is often important that the model is non-stochastic. Therefore, unlike [58, 16], our halting module is designed to be deterministic.

Our contributions can be summarized as follows:

- We propose a deterministic module that progressively halts tokens throughout the transformer, and a simple but effective token recycling mechanism to reuse the information from the halted tokens.

- An equivalent differentiable forward-pass is proposed to overcome the non-differentiability of token halting, and a theoretical analysis is conducted to evaluate the accuracy of the pseudo-gradient.

---

[*]mao.ye@getcruise.com
[†]greg.meyer@getcruise.com

- A non-uniform token sparsity loss is employed to improve the learning of the halting module by utilizing the ground-truth bounding boxes.

## 2. Related Work

### 2.1. Dynamic Transformer

The idea of adapting the number of tokens within a transformer to improve performance has recently been explored. [58, 87] learn a token selection module to dynamically halt tokens at inference using the Gumbel-softmax trick [24] and ACT [21], respectively. [54, 16, 31, 70] propose different heuristics based on the attention weights to halt or aggregate tokens. [26] combines both token selection and aggregation. [80] proposes a slow-fast token update that applies token-wise transformations on the halted tokens and attention-based transformations on those that are not halted. [101] proposes to globally aggregate the tokens into a smaller set of new tokens using a reconstruction loss, which encourages the new tokens to preserve as much information as possible. Instead of adaptively removing or combining tokens, [76, 100] consider using tokens with adaptive spatial size for different inputs, and [46] simultaneously selects the tokens, attention heads, and attention windows.

The majority of the existing works are designed for image classification, while we consider the task of 3D object detection. Consequently, the prior work cannot be directly applied to object detection and significant algorithmic changes are required. Vision transformers for image classification use an artificial token, typically referred to as the `cls` token, to perform classification. However, for object detection, all tokens are used by a detection head. Therefore, our method requires a way to aggregate information from all tokens regardless of if or when a token was halted. Additionally, the aggregation of halted tokens needs to be differentiable. As a result, our proposed method differs significantly from the previous work.

Two prior works that do consider the task of object detection are [60] and [92]. [60] applies dynamic token selection within the Deformable DETR framework [99], and [92] employs non-uniform point cloud downsampling for a point-based transformer. We consider [60] orthogonal to our approach as it focuses on improving the convergence of the transformer detection head, while we focus on improving the backbone. Furthermore, Deformable DETR selects a fixed subset of keys/values for each query. We select a dynamic subset of tokens, which in turn reduces all queries, keys, and values. In our case, the token selection is more challenging because we need to determine the token importance in the early stages of the backbone, where the token features are less informative. However, Deformable DETR is applied to the detection head and has access to deep features extracted from the backbone. In terms of [92], it at-

tempts to downsample the background points, while our approach tries to halt any unnecessary tokens. Due to the different goals, the algorithms are considerably different.

### 2.2. Efficient Network Architecture

Trading off between network efficiency and accuracy has been a long-standing problem. Common approaches such as network pruning [41, 18, 43, 84, 86, 69, 86, 89], network quantization [59, 94, 93, 52, 3, 40], and neural architecture search [62, 39, 38, 68, 36, 6, 19, 29, 8, 37] have been proposed to push the limit of the Pareto front. Compared with our method, those approaches aim to search for a fixed network architecture that is not adaptive to the input. On the other hand, learning a network with a dynamic computation graph has also been explored in various directions including adaptive resolution [53, 83, 47], depth [17, 73, 75, 79], and channels [33, 90, 4]. Those directions are orthogonal to this work and could be performed in combination with our method.

### 2.3. LiDAR-based 3D Object Detection

Existing networks for 3D detection can be classified based on the representation of the 3D scene. The most popular approach is to represent the scene using a voxel grid. Vote3Deep [12] was one of the first to use a uniform voxel grid to represent the point cloud. The representation has been further improved upon by [97, 81] using a small PointNet [56] to learn a better voxel representation, and by [20] using sparse 3D convolutions to improve the efficiency. The efficiency was improved further by [82, 27] using 2D convolutions instead of 3D convolutions. Another popular approach is to directly process the point representation [55, 64, 9, 78]. These methods are usually built on PointNet [56] and require a nearest neighbor search between points which can be difficult to scale. Finally, there are a handful of methods that represent the 3D scene using a range image [49, 32, 48, 50, 15, 7].

Due to the recent progress in transformers for computer vision [11], the transformer architecture has also been applied to 3D object detection. Existing works include transformer-based backbones [45, 23, 13, 66, 14, 74] for the voxel-based representation, [51, 92] for the point-based representation, and [22] for a combination of both the point and voxel representation. Furthermore, transformers have been used to improve the detection head [98] and for sensor fusion [2, 91, 77].

Similar to our work, [67] proposes to use a foreground point selection to remove LiDAR points that do not belong to objects. This idea was later applied to a transformer-based detector [66]. This selection process is based purely on whether the points/voxels belong to foreground objects. In comparison, our dynamic token halting is learned based on a token's contribution to the detection task. That is, im-

portant background tokens can be kept while less important foreground tokens can also be removed. In addition, our framework incrementally halts tokens, and all tokens are used to inform the final predictions, while in [67, 66] voxels/points are simply removed early in the network. Lastly, the voxel selection in [66] is applied at the end of feature learning and thus has only a mild impact on latency. In addition, [10] proposes to select a subset of voxels as input to sparse convolutions. An important difference between our approach and [10] is how non-differentiable operations are handled. We propose an equivalent differentiable forward-pass to enable full differentiability with theoretical support. [10] uses a non-differentiable threshold during training; therefore, voxels that are not selected do not receive end-to-end gradient updates.

## 3. Background

### 3.1. 3D Object Detection

Given $P = \{(x, y, z, r, t)_i\}$, a point cloud of 3D coordinates $(x, y, z)$, intensity $r$, and elongation $t$ measurements, we want the model to predict a set of 3D bounding boxes for the objects of interest. The bounding box $b = (l_x, l_y, l_z, w_x, w_y, w_z, \alpha)$ parameterizes the location of an object by its center position $(l_x, l_y, l_z)$, dimensions $(w_x, w_y, w_z)$, and heading angle $\alpha$.

### 3.2. Transformer-based 3D Object Detector

Our method is built upon the recent Single-stride Sparse Transformer (SST), which is a LiDAR-based 3D object detector [13]. It uses dynamic voxelization [96] to create a set of voxel features, $f_0$, in the bird's eye-view, and each voxel is treated as a token within the transformer. For this reason, we will use the terms voxels and tokens interchangeably. SST uses regional grouping to divide the voxel grid into non-overlapping regions, and it applies sparse regional attention (SRA) to the voxels within each region. Specifically,

$$\text{SRA}(f_l) = \text{MLP}(\text{LN}(f_l')) + f_l' \quad (1)$$

where

$$f_l' = \text{MSA}(\text{LN}(f_l), \text{PE}(f_l)) + f_l, \quad (2)$$

and LN, MLP, PE, MSA denote layer normalization [1], token-wise multi-layer perceptron, position encoding, and multi-head self-attention [72], respectively. Since objects can lay within several regions, two SRA operations are performed sequentially where the regional grouping is shifted for one of the operations, similar to the Swin Transformer [42]. Thus, the $(l + 1)$-th features are generated by

$$f_{l+1} = \text{SSRA}(\text{SRA}(f_l)), \quad (3)$$

where SSRA is the shifted sparse region attention. After several rounds of attentions, the bird's eye-view feature map

is constructed from the final voxel features, and the feature map is passed to a detection head to predict the object bounding boxes.

## 4. Proposed Method

**Notation** Without loss of generality and to simplify the notation, we assume that each layer of our transformer backbone processes the tokens as follows:

$$f_{l+1} = \phi_2(\text{SA}(\phi_1(f_l)), f_l), \quad (4)$$

where $\phi_1$ and $\phi_2$ are non-linear token-wise operations, and SA is some general self-attention mechanism (not necessarily the SRA used by SST). Note that each basic block (*i.e.* (3)) of SST can be viewed as two layers using this notation due to the usage of region shift.

### 4.1. Overview

We are motivated by the idea that different tokens may vary in importance to the detection task. Furthermore, a token may be useful at the early stages of the network but less informative at the later stages. For example, detecting a vehicle on an empty road may require several late stage tokens from the vehicle itself, but only a few early stage tokens from the surrounding road. Alternatively, detecting a camouflaged pedestrian may require several late stage tokens from both the person and the surrounding environment. Our goal is to learn a token halting mechanism that identifies which tokens should be kept and forwarded to the subsequent layer. Given $N$ token features, $f_l = \{f_{l,i} : i \in \{1, \ldots, N\}\}$, from the $l$-th layer, the halting module outputs a binary mask, $k_l \in \{0, 1\}^N$, indicating which of the tokens are forwarded to the next layer. Specifically, the subsequent layer's tokens are computed as

$$f_{l+1} = \phi_2(\text{SA}(\phi_1(\hat{f}_l)), \hat{f}_l), \quad (5)$$

where $\hat{f}_l = \{f_{l,i} : k_{l,i} = 1\}$ are the tokens kept by the halting module. Refer to Figure 1 for illustration of our proposed method.

### 4.2. The Halting Module

The halting module outputs a binary mask $k_l$ that determines whether a token is forwarded to the next layer. To create such a mask, our halting module first computes a non-negative score $s$ for each token. The architecture used to produce the score could be anything; it could be as simple as a MLP or as complex as a transformer. Afterwards, the mask can be obtained by thresholding the score, $k_l = \psi(s_l)$, where $\psi(s_l) = \mathbb{I}\{s_l \geq u\}$ and $u$ is a threshold. The threshold, $u$, could be a fixed value or selected dynamically using the distribution of the scores.

Each layer has its own halting module as the distribution of features shifts throughout the network. In practice, we
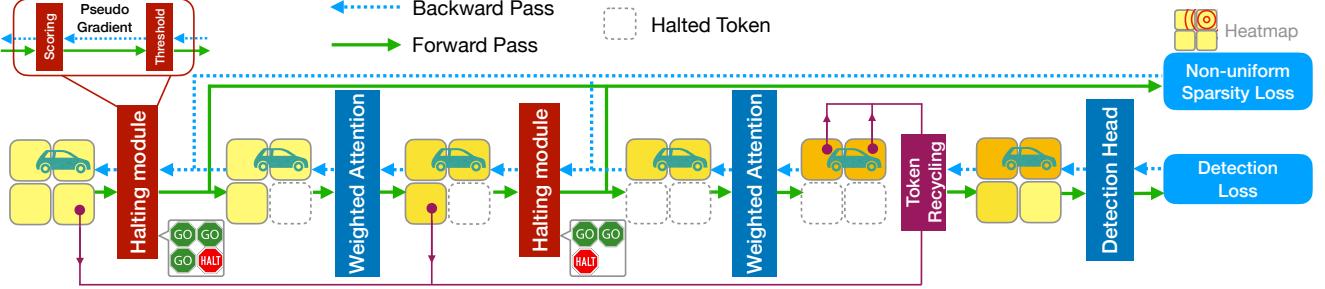
Figure 1. Given a set of tokens, the halting module produces a score for each token, and the tokens with scores below a threshold will be halted. At inference, only tokens that are not halted are forwarded to the next layer. However, during training, all the tokens are forwarded to the next layer, but halted tokens are prevented from interacting with the other tokens. We do this in order to obtain a pseudo-gradient to back-propagate through the non-differentiable threshold operation. After the last attention layer, halted tokens are recycled and combined with the non-halted ones and forwarded to the detection head. The whole network is trained end-to-end using both a detection loss and a non-uniform sparsity loss.

found that it is better to use a more complicated architecture for the halting module during the earlier layers, while a simple architecture is sufficient for for the later layers. Refer to Section 6.1, for more details on the specific architecture used by our method.

### 4.3. Weighted Self-Attention

To improve the learning of the halting module, we propose a weighted self-attention operation $\text{WSA}(f_l, s_l)$, which weights the tokens based on the score $s_l$ produced by the halting module. The output of the weight self-attention for the $i$-th token is defined as follows:

$$\text{WSA}(f_l, s_l)_i = \frac{\sum_j \exp(P_{ij}) s_{l,j} v_j}{\sum_k \exp(P_{ik}) s_{l,k}} \qquad (6)$$

where

$$P = (W_Q f_l)(W_K f_l)/\sqrt{d}, \qquad v = W_V f_l, \qquad (7)$$

and $W_Q$, $W_K$, and $W_V$ are the queue, key, and value matrices, respectively. This definition assumes the standard self-attention operation; however, it could be easily applied to other variants.

By using the token score $s_l$ as the weight, we force the attention given to a particular token to be proportional to that token's score. Intuitively, this encourages the halting module to increase the score of a token if it plays an important role within the attention mechanism. Refer to Section 5, for an in-depth analysis of how the weighted self-attention affects the learning of the halting module.

### 4.4. Token Recycling

We observe that even when a token is halted early, its features are still useful to inform the final predictions of the model. Instead of throwing away the halted tokens, we recycle them by directly forwarding them to the detection head.

Recall that SST constructs a bird's eye-view feature map, $f_{\text{BEV}}$, from the output of the transformer, and passes that feature map to the detection head. Since our method is built upon SST, we also construct a bird's eye-view feature map $\hat{f}_{\text{BEV}}$; however, in our case, the feature map is assembled from both the output of the transformer and all the halted tokens using their final features at the layer they are halted. Refer to Figure 1 for an illustration of this process.

### 4.5. Equivalent Differentiable Forward-Pass

The token halting operation is non-differentiable, which is an issue for training with gradient descent. To overcome this issue, we introduce an equivalent differentiable forward-pass (EDF) during training. The goal of EDF is to produce a forward-pass that generates an equivalent output but is differentiable.

Unlike during inference, where the halted tokens are not forwarded to the subsequent layer, EDF forwards all the tokens but prevents "halted" tokens from interacting with other tokens by using the mask as a multiplier on their score. Specifically,

$$f_{l+1} = \phi_2(\text{WSA}(\phi_1(f_l), s_l \circ k_{0:l}), f_l), \qquad (8)$$

where

$$k_{0:l} = k_0 \circ \ldots \circ k_l \qquad (9)$$

and $\circ$ denotes an element-wise multiply. Note that, we set $k_0 := \mathbf{1}$. For each token, we use its feature at the layer it was halted to construct $f_{\text{BEV}}$. Similarly, EDF defines the feature map as

$$f_{\text{BEV}} = \sum_{l=1}^{L+1} (k_{0:l-1} - k_{0:l}) \circ f_l. \qquad (10)$$

For convenience, we set $k_{0:L+1} := \mathbf{0}$ where $L$ is the total number of layers. Although inference and training have different forward-passes due to EDF, it is trivial to verify that $f_{\text{BEV}} = \hat{f}_{\text{BEV}}$.

The binary mask $k_l$ is produced by thresholding the token's score $s_l$, and this thresholding function $\psi(s_l)$ has zero gradients almost everywhere. To enable back-propagation, we use the straight-through estimator (STE) [5]. STE defines a pseudo-gradient during the back-propagation, by replacing the derivative of the threshold function with the derivative of some other activation function

$$\psi'(s_l) := \sigma'(s_l). \tag{11}$$

A common choice for $\sigma$ is the identity function.

By combining the EDF with the STE, our computation graph at training is fully differentiable. In Section 5, we provide a detailed analysis of the (pseudo)-gradients used to update the halting module.

### 4.6. Non-Uniform Token Sparsity Loss

We want to encourage the halting module to output a sparse binary mask, and a common approach [28, 85] is to penalize the scores $s_l$ with a $\ell_1$ penalty (*i.e.* LASSO [71]). Such a penalty is applied uniformly to all the tokens. On the other hand, we find that the performance can be significantly improved by applying a non-uniform penalty to the tokens. Our intuition is that tokens belonging to a foreground object are usually more important than ones belong to the background. Therefore, on average, we want foreground tokens to have a larger score. To accomplish this, we leverage the grouth-truth bounding boxes and create a heatmap in a similar fashion as [95]. Tokens that are within a bounding box have a positive value in the heatmap between $[0, 1]$, and a token's value increases as it gets closer to the center of object. Tokens that do not fall inside any bounding box have a value of zero in the heatmap. Any difference between $s_l$ and the heatmap is penalized using focal loss [35]. Such a loss applies a uniform sparse penalty to the background token and a non-uniform penalty to the foreground tokens based on their distance to the object's center. See Appendix A.1.1 for more details.

### 4.7. Losses

The feature map $f_{\text{BEV}}$ is forwarded to the Center-Point [88] detection head for predictions, and we train the model end-to-end with a total loss defined as

$$\mathcal{L} = \lambda_b \mathcal{L}_b + \lambda_h \mathcal{L}_h + \lambda_s \mathcal{L}_s, \tag{12}$$

where $\mathcal{L}_b$ and $\mathcal{L}_h$ are the box and heatmap regression losses defined by [88], and $\mathcal{L}_s$ is our non-uniform sparsity loss.

## 5. Analyzing the Pseudo-Gradient

We demonstrate that the pseudo-gradient provided by our proposed EDF and the STE gives a high-quality update direction for the halting module. To simplify the analysis,

we consider a reduced model with only one attention layer. In this case, $f_{\text{BEV}} = \{q_1, q_2\}$ where we define

$$q_1 = (\mathbf{1} - k) \circ f, \; q_2 = k \circ \phi_2(\text{WSA}(\phi_1(f), s \circ k), f). \tag{13}$$

Note that we drop the subscripts to simplify the notation. Furthermore, we use a STE of $\psi'(s) = 1$. Consider the situation that the $i$-th feature is halted in the first layer, *i.e.* $k_i = 0$, while it is actually better to forward to the next layer. In such case, a feature map $\tilde{f}_{\text{BEV}} = \{\tilde{q}_1, \tilde{q}_2\}$, where

$$\tilde{q}_1 = (\mathbf{1} - k - \mathbf{1}_i) \circ f \tag{14}$$
$$\tilde{q}_2 = (k + \mathbf{1}_i) \circ \phi_2(\text{WSA}(\phi_1(f), s \circ (k + \mathbf{1}_i)), f), \tag{15}$$

and $\mathbf{1}_i = [0, 0, \ldots, 1, 0, 0, \ldots]$ (only the $i$-th index being 1), is expected to have a smaller loss. Specifically,

$$\Delta_i := \mathcal{L}(\tilde{q}_1, \tilde{q}_2) - \mathcal{L}(q_1, q_2) < 0, \tag{16}$$

where $\mathcal{L}$ is a detection-related loss. We are able to show that

$$\Delta_i \approx \frac{\partial \mathcal{L}(q_1, q_2)}{\partial s_i} + O(u), \tag{17}$$

where the approximate equality has precision up to the second-order Taylor approximation. That is to say, our pseudo-gradient is almost identical to the change of loss, and the additional error term has a magnitude proportional to the threshold $u$. Since the threshold is, in general, very small (a practical choice being $u \sim 0.01$), the pseudo-gradient gives an accurate update direction, pushing the halting module towards a direction that gives a better halting decision. The analysis of when the token is forwarded but it is better to halt is similar and thus omitted. See Appendix A.1.2 for more details and a derivation.

## 6. Experiment

### 6.1. Setup and Implementation Details

**Dataset** We evaluate our method using the Waymo Open Dataset (WOD) [65] which contains 1,150 sequences where 798 sequences are for training, 202 for validation, and 150 for testing. The entire dataset contains more than 200k frames, and each frame contains a LiDAR point cloud covering a $150m \times 150m$ area.

**Backbone** For our experiments, we use the default backbone configuration of SST [13]. It uses a dynamic voxelization technique similar to PointPillars [27] with the LiDAR point cloud as the input. Furthermore, SST uses four consecutive SRA blocks followed by four dense convolutional layers. Each self-attention operation uses 8 heads, 128 input channels, and 256 hidden channels. Each spatial region contains $14 \times 14 \times 1$ voxels, and each voxel has a size of $0.32m \times 0.32m \times 6m$. We use the CenterPoint [88] one-stage detection head and single frame of input.

**Halting Module** We employ two dynamic halting modules before the first and second SRA blocks. Since we obtain a very high-level of token sparsity after the second halting module with negligible performance drop (see below), adding more modules to further increase sparsity gives limited speed-up. Like mentioned in Section 4.2, we use a complex architecture for the first halting module and simple architecture for the second one. For the first halting, we use a lightweight U-Net [61] architecture with MobileNetV2 blocks [62], and a linear layer with sigmoid activation to obtain the halting score. For the second halting module, we simply apply a one-layer MLP on the input feature to produce the score. For both halting modules, the first 32 features of a token are used as input. To re-use the latent features extracted by the halting module, we fuse the penultimate latent features back into the token features.

The halting threshold $u$ is adjusted to obtain different levels of token sparsity. The threshold is set such that only a certain quantile of tokens are halted. We found that setting the threshold based on the score quantiles instead of a fixed threshold helps stabilize training. We refer readers to Appendix A.2.1 for more details.

**Training** We train the network for 24 epochs using AdamW optimizer [44] and a one-cycle learning rate. The learning rate starts at $4 \times 10^{-5}$, and increases to $1 \times 10^{-3}$ during the first $10\%$ of iterations. Afterwards, the learning rate is annealed with consine decay for the rest of the iterations. We apply standard data augmentation (random flop, rotation, and scaling) during the training. The loss weights are set to $\lambda_b = 2.0$, $\lambda_h = 1.0$, and $\lambda_s = 0.5$.

## 6.2. Efficiency and Accuracy Trade-off

The goal of this section is to study the efficiency and accuracy trade-off achieved by varying the sparsity of the tokens. We compare our dynamic token halting approach with other model scaling approaches, including changing the latent dimension of the attention mechanism (*i.e.* width scaling) and the number of attention heads (*i.e.* # head scaling). Furthermore, we adapted AViT [87], a dynamic token halting approach for image classification, to the 3D object detection task. We evaluate the latency of the backbone on a high-end NVIDIA A100 GPU and report the relative speed-up compared to the original architecture. Figure 2 plots the Pareto frontier of detection performance versus speed-up achieved by the various approaches. It is clear from the figure, that our method provides the best efficiency and accuracy trade-off. In fact, our method can achieve over a $50\%$ speed-up with only a slight impact to performance. Moreover, we find that AViT does not perform well, which indicates that a straightforward adaptation of a dynamic halting approach for image classification does not perform well when applied to the 3D object detection. We refer readers
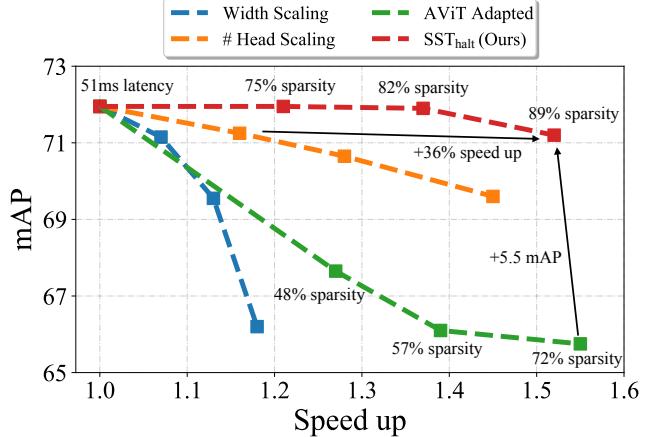


Figure 2. The accuracy-efficiency trade-off. The mAP is the average of L1 and L2 vehicle AP. The sparsity is measured by averaging the percentage of halted tokens across all layers.

to Appendix A.2.2 for more details on the experiment and a completed list of the numerical results.

## 6.3. Comparisons with the State-of-the-Art

In the previous section, we demonstrated that our proposed method can significantly reduce the latency of SST without dramatically impacting its performance. The goal of this section is to leverage the latency savings provided by our dynamic token halting to improve the performance of SST while maintaining its runtime. We improve the performance of SST by simply increasing the capacity of the detection head. Specifically, we add an additional convolutional block and a feature pyramid network [34]. The token sparsity is adjusted to ensure this improved SST has the same latency as the original model. Refer to Appendix A.2.3 for a detailed description of the architecture.

Table 1 compares the performance of our improved model, referred to as $\text{SST}_{\text{halt}}^{++}$, to the baseline model, SST [13], on the Waymo Open Dataset validation split. We observe that the performance of our proposed method outperforms the baseline (SST). Compared to the original SST, we see a 1% to 3.5% AP/APH improvement for all classes. Recall that the improvement in detection performance is accomplished without increasing the latency of the model. Furthermore, our method meets or exceeds the performance of all other state-of-the-art models. Table 2 shows the performance on the Waymo Open Data test split, and we observe a similar improvement.

## 6.4. Additional Analysis

We find that dynamic token halting is not only an effective approach to obtain a better efficiency-accuracy trade-off, but also it has the additional effect of aiding in the detection of long-range and difficult objects. In Table 3, we show the breakdown of detection performance at different

| Method | Publication Date | TS | Vehicle | | Pedestrian | | Cyclist | |
|---|---|---|---|---|---|---|---|---|
| | | | AP/APH L1 | AP/APH L2 | AP/APH L1 | AP/APH L2 | AP/APH L1 | AP/APH L2 |
| PointPillar [27] | CVPR 2019 | × | 72.1/71.5 | 63.6/63.1 | 70.6/56.7 | 62.8/50.3 | 64.4/62.3 | 61.9/59.9 |
| PV-RCNN [63] | CVPR 2020 | ✓ | 77.5/76.9 | 69.0/68.4 | 75.0/65.6 | 66.0/57.6 | 67.8/66.4 | 65.4/64.0 |
| RangeDet [15] | ICCV 2021 | × | 72.9/72.3 | 64.0/63.6 | 75.9/71.9 | 67.6/63.9 | 65.7/64.4 | 63.3/62.1 |
| MVP++ [57] | CVPR 2021 | ✓ | 74.6/— | —/— | 78.0/— | —/— | —/— | —/— |
| LiDAR R-CNN [30] | CVPR 2021 | ✓ | 76.0/75.5 | 68.3/67.9 | 71.2/58.7 | 63.1/51.7 | 68.6/66.9 | 66.1/64.4 |
| CenterPoint [88] | CVPR 2021 | ✓ | 76.1/75.5 | 68.0/67.5 | 76.1/65.1 | 68.1/57.9 | —/— | —/— |
| RSN [67] | CVPR 2021 | × | 75.1/74.6 | 66.0/65.5 | 77.8/72.7 | 68.3/63.7 | —/— | —/— |
| IA-SSD [92] | CVPR 2022 | ✓ | 70.5/69.7 | 61.6/61.0 | 69.4/58.5 | 60.3/50.7 | 67.7/65.3 | 65.0/62.7 |
| LidarNAS [37] | ECCV 2022 | × | 75.6/— | —/— | 77.4/— | —/— | —/— | —/— |
| VoTr-TSD [45] | ICCV 2021 | ✓ | 74.9/74.3 | 65.9/65.3 | —/— | —/— | —/— | —/— |
| M3DETR [22] | WACV 2022 | × | 75.7/75.1 | 66.6/66.0 | 65.0/56.4 | 56.0/48.4 | 65.4/64.2 | 62.7/61.5 |
| SWFormer [66] | ECCV 2022 | × | **77.8/77.3** | 69.2/68.8 | 80.9/72.7 | 72.5/64.9 | —/— | —/— |
| SST [13] | CVPR 2022 | × | 74.2/73.8 | 65.5/65.1 | 78.7/69.6 | 70.0/61.7 | —/— | —/— |
| SST*$_{center}$ [13] | CVPR 2022 | × | 76.2/75.7 | 67.7/67.2 | 79.9/71.4 | 72.7/64.8 | 67.7/66.3 | 65.2/63.8 |
| SST$_{halt}^{++}$ (Ours) | - | × | 77.7/77.1 (+1.5) | **69.5/69.0**(+1.8) | **80.9/73.0** (+1.3) | **74.0/66.5**(+1.5) | **70.0/68.6**(+2.3) | **67.3/66.0**(+2.2) |

Table 1. Performance comparison on the Waymo Open Dataset validation split. All methods take a single frame of LiDAR data as input. TS denotes whether or not a two-stage detection head is used. Methods below the first middle separator are transformer-based detectors. Note that, our proposed method, SST$_{halt}^{++}$, does not apply any test-time augmentations, use a model ensemble, or use a two-stage detection head. *The performance of SST$_{center}$ is based on our own implementation of SST with a CenterPoint detection head. The bold/underscored numbers correspond to the best and second best approach. The blue numbers in the parentheses are the average AP/APH improvement we obtain by applying our token halting approach to SST.
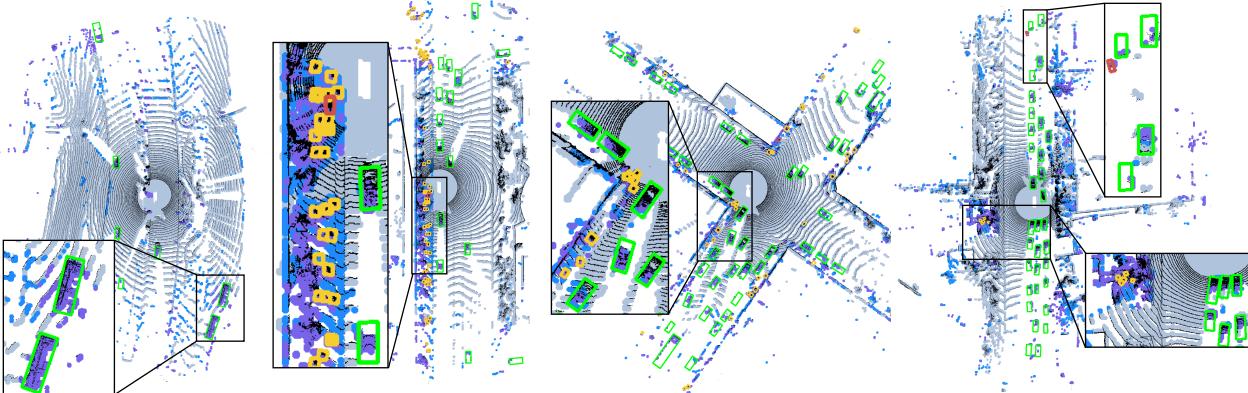


Figure 3. Visualization of the token halting process. The grey/blue voxels are halted at the first/second layer, and the purple voxels are not halted. The green, yellow, and red boxes are the vehicles, pedestrians, and cyclists detections, respectively.

| Method | Vehicle | | Pedestrian | |
|---|---|---|---|---|
| | AP/APH L1 | AP/APH L2 | AP/APH L1 | AP/APH L2 |
| PointPillar | 68.6/68.1 | 60.5/60.1 | 68.0/55.5 | 61.4/50.1 |
| CenterPoint_2f | 80.2/79.7 | 72.2/71.8 | 78.3/72.1 | 72.2/66.4 |
| SST$_{center}$ | 80.0/79.6 | 72.1/71.7 | 79.3/71.3 | 73.4/65.9 |
| SST$_{halt}^{++}$ (Ours) | **81.2/80.7** (+1.2) | **73.7/73.2** (+1.6) | **80.4/72.8** (+1.3) | **74.7/67.5** (+1.5) |

Table 2. Performance on the Waymo Open Dataset test split. The setting is the same as Table 1.

distances and difficulty levels for the original SST, the SST with dynamic token halting (denoted as SST$_{halt}$), and our improved SST$_{halt}^{++}$. Compared to the baseline, SST$_{halt}$ has better accuracy for the difficult L2 objects while the same performance for L1 objects. Moreover, for both SST$_{halt}$ and SST$_{halt}^{++}$, the improvement becomes more significant for long-range objects. We believe this result can be explained

intuitively. For difficult and long-range objects, their points are more sparse, which may make it challenging for the model to distinguish the foreground objects from the background. Our dynamic token halting method removes the majority of the background tokens, which perhaps increases the signal-to-noise ratio for the detector.

## 6.5. Ablation Study

We conducted an ablation study to show the efficacy of the various aspects of our proposed method. For each experiment, we keep the sparsity/latency the same and study how the performance changes when we remove different components. Specifically, we are interested in how the performance changes when the token recycling is removed and when the non-uniform sparsity loss is replaced with an uniform sparsity loss. Table 4 summarizes the result.

| Config | Vehicle BEV AP | | | | | |
| | [0, 30) L1 | [0, 30) L2 | [30, 50) L1 | [30, 50) L2 | [50, ∞) L1 | [50, ∞) L2 |
|---|---|---|---|---|---|---|
| $SST_{center}$ | 92.2 | 91.0 | 74.6 | 68.0 | 52.7 | 40.8 |
| $SST_{halt}$ | 92.2 (+0.0) | 91.1 (+0.1) | 74.6 (+0.0) | 68.5 (+0.5) | 52.7 (+0.0) | 41.2 (+0.4) |
| $SST_{halt}^{++}$ | 92.6 (+0.4) | 91.4 (+0.3) | 76.2 (+1.6) | 70.0 (+2.0) | 55.4 (+2.7) | 43.4 (+2.6) |

Table 3. Breakdown of the BEV vehicle detection performance. $SST_{center}$ is our implementation of SST with a CenterPoint detection head. $SST_{halt}$ is the SST model with our dynamic halting (the model explored in section 6.2). $SST_{halt}^{++}$ is the model we proposed in section 6.3. The blue numbers in the parentheses are the improvement we obtained compared with $SST_{center}$.

| Config | Vehicle | | Pedestrian | | Cyclist | |
| | AP/APH L1 | AP/APH L2 | AP/APH L1 | AP/APH L2 | AP/APH L1 | AP/APH L2 |
|---|---|---|---|---|---|---|
| Full | 75.7/75.2 | 67.7/67.2 | 78.7/70.4 | 72.3/64.5 | 70.0/68.5 | 67.8/66.2 |
| w/o voxel recycle | 61.9/61.0 (-14.0) | 54.9/54.0 (-13.0) | 76.8/68.8 (-1.8) | 69.7/62.2 (-2.5) | 63.7/62.3 (-6.2) | 61.3/60.0 (-6.4) |
| w/o non-uniform sparsity | 74.5/74.0 (-1.2) | 66.1/65.6 (-1.6) | 77.7/68.1 (-1.6) | 70.2/61.3 (-2.6) | 64.2/62.7 (-5.8) | 61.8/60.3 (-6.0) |

Table 4. Ablation study where "w/o non-uniform sparsity" denotes that we use a uniform sparsity penalty instead. The red numbers in the parentheses are the average AP/APH decrease when a certain component is removed.

We observe a considerable drop in performance, especially for vehicles and cyclists, when we do not recycle the halted tokens. We do not believe this is because the halting module is not keeping important tokens. If it was not capturing important tokens, we would have seen a significant drop in performance even with token recycling because when a token is halted, it will only be processed by one or two attention blocks which is only a quarter or half of the backbone. However, in Section 6.2, there is very little performance drop even with very high sparsity. We believe the performance degradation is due to the halted tokens still having useful semantic information for the detection and thus should be used by the detection head.

The non-uniform sparsity loss is also shown to be useful. To demonstrate this, we replaced the non-uniform sparsity loss with a uniform sparsity loss, *i.e.* $\ell_1$ loss. Since our framework is fully differentiable, the model is still able to achieve good results with a uniform sparsity loss. However, the non-uniform loss provides a better signal to the model. To understand how, we plot the change in sparsity of background and foreground tokens when training with and without our non-uniform sparsity loss in Figure 4.

With the uniform sparsity loss, as training progresses, the ratio of kept foreground tokens increases while the ratio of background tokens decreases. We believe this occurs because the model gradually becomes better at the detection task and in turn becomes better at selecting useful tokens. However, during the early stages of training, the model fails to make a good token halting decision, which negatively impacts the learning of the model. As a result, the model ends up in a sub-optimal state. When the non-uniform sparsity loss is applied, we observe two interesting phases during learning. In the first phase (within the first epoch), the model is learning to select tokens primarily based on whether they belong to a foreground object. In this phase, the sparsity of foreground tokens rapidly decreases
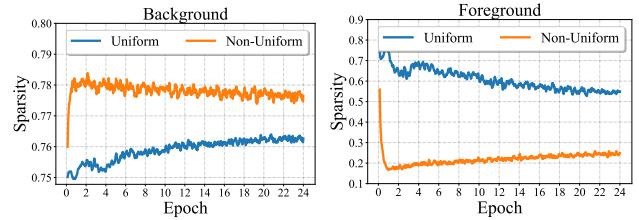


Figure 4. The sparsity of background and foreground tokens when training with the uniform and non-uniform sparsity loss.

while the sparsity of background tokens rapidly increases. Afterwards, the model goes into a second learning phase, and it starts to adjust the selection of tokens based on their contribution to the detection task. In this phase, the model learns to halt useless foreground tokens while keeping useful background tokens.

## 6.6. Visualization

To better understand the token halting process, we visualize the tokens halted at different layers with different colors in Figure 3. From the figure, we observe that most of the tokens belong to roads, sides of buildings, or foliage, which do not contribute much to the detection. Appropriately, these tokens are halted at the very beginning by the model. In addition, our model tends to keep tokens that are near the objects we want to detect, since those tokens are usually semantically informative. For example, in the leftmost portion of the second scene, many of the background voxels near pedestrians are kept while the voxels in the rightmost portion are almost all halted as there are no pedestrians. Lastly, we observed that tokens belonging to objects remained unhalted with a significantly higher chance as those tokens are in general critical for the detection.

# 7. Conclusion

In this work, we proposed an approach to dynamically halt tokens in order to speed-up transformer-based 3D object detectors. Our method significantly improves the Pareto frontier of the accuracy and efficiency trade-off. By leveraging our improved model efficiency, we are able to dramatically increase the performance of the baseline model without sacrificing latency. There are several interesting directions for future research. We believe there are additional architectural changes that could be made to further exploit the high-level of token sparsity and to further reduce model latency. It would also be interesting to explore how to generalize this idea to multi-modal and multi-frame object detectors.

# References

[1] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016. 3

[2] Xuyang Bai, Zeyu Hu, Xinge Zhu, Qingqiu Huang, Yilun Chen, Hongbo Fu, and Chiew-Lan Tai. Transfusion: Robust lidar-camera fusion for 3d object detection with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1090–1099, 2022. 2

[3] Ron Banner, Yury Nahshan, and Daniel Soudry. Post training 4-bit quantization of convolutional networks for rapid-deployment. *Advances in Neural Information Processing Systems*, 32, 2019. 2

[4] Babak Ehteshami Bejnordi, Tijmen Blankevoort, and Max Welling. Batch-shaping for learning conditional channel gated networks. In *International Conference on Learning Representations*, 2020. 2

[5] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013. 5

[6] Han Cai, Chuang Gan, Tianzhe Wang, Zhekai Zhang, and Song Han. Once-for-all: Train one network and specialize it for efficient deployment. *arXiv preprint arXiv:1908.09791*, 2019. 2

[7] Yuning Chai, Pei Sun, Jiquan Ngiam, Weiyue Wang, Benjamin Caine, Vijay Vasudevan, Xiao Zhang, and Dragomir Anguelov. To the point: Efficient 3d object detection in the range image with graph convolution kernels. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16000–16009, 2021. 2

[8] Arnav Chavan, Zhiqiang Shen, Zhuang Liu, Zechun Liu, Kwang-Ting Cheng, and Eric P Xing. Vision transformer slimming: Multi-dimension searching in continuous optimization space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4931–4941, 2022. 2

[9] Qi Chen, Lin Sun, Zhixin Wang, Kui Jia, and Alan Yuille. Object as hotspots: An anchor-free 3d object detection approach via firing of hotspots. In *European conference on computer vision*, pages 68–84. Springer, 2020. 2

[10] Yukang Chen, Yanwei Li, Xiangyu Zhang, Jian Sun, and Jiaya Jia. Focal sparse convolutional networks for 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5428–5437, 2022. 3

[11] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. 1, 2

[12] Martin Engelcke, Dushyant Rao, Dominic Zeng Wang, Chi Hay Tong, and Ingmar Posner. Vote3deep: Fast object detection in 3d point clouds using efficient convolutional neural networks. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1355–1361. IEEE, 2017. 2

[13] Lue Fan, Ziqi Pang, Tianyuan Zhang, Yu-Xiong Wang, Hang Zhao, Feng Wang, Naiyan Wang, and Zhaoxiang Zhang. Embracing single stride 3d object detector with sparse transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8458–8468, 2022. 2, 3, 5, 6, 7

[14] Lue Fan, Feng Wang, Naiyan Wang, and Zhaoxiang Zhang. Fully sparse 3d object detection. *arXiv preprint arXiv:2207.10035*, 2022. 2

[15] Lue Fan, Xuan Xiong, Feng Wang, Naiyan Wang, and Zhaoxiang Zhang. Rangedet: In defense of range view for lidar-based 3d object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2918–2927, 2021. 2, 7

[16] Mohsen Fayyaz, Soroush Abbasi Kouhpayegani, Farnoush Rezaei Jafari, Eric Sommerlade, Hamid Reza Vaezi Joze, Hamed Pirsiavash, and Juergen Gall. Ats: Adaptive token sampling for efficient vision transformers. *arXiv preprint arXiv:2111.15667*, 2021. 1, 2

[17] Michael Figurnov, Maxwell D Collins, Yukun Zhu, Li Zhang, Jonathan Huang, Dmitry Vetrov, and Ruslan Salakhutdinov. Spatially adaptive computation time for residual networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1039–1048, 2017. 2

[18] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. *arXiv preprint arXiv:1803.03635*, 2018. 2

[19] Chengyue Gong, Dilin Wang, Meng Li, Xinlei Chen, Zhicheng Yan, Yuandong Tian, qiang liu, and Vikas Chandra. NASVit: Neural architecture search for efficient vision transformers with gradient conflict aware supernet training. In *International Conference on Learning Representations*, 2022. 2

[20] Benjamin Graham, Martin Engelcke, and Laurens Van Der Maaten. 3d semantic segmentation with submanifold sparse convolutional networks. In *Proceedings of the*

*IEEE conference on computer vision and pattern recognition*, pages 9224–9232, 2018. 2

[21] Alex Graves. Adaptive computation time for recurrent neural networks. *arXiv preprint arXiv:1603.08983*, 2016. 2

[22] Tianrui Guan, Jun Wang, Shiyi Lan, Rohan Chandra, Zuxuan Wu, Larry Davis, and Dinesh Manocha. M3detr: Multi-representation, multi-scale, mutual-relation 3d object detection with transformers. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 772–782, 2022. 2, 7

[23] Chenhang He, Ruihuang Li, Shuai Li, and Lei Zhang. Voxel set transformer: A set-to-set approach to 3d object detection from point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8417–8427, 2022. 2

[24] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*, 2016. 2

[25] Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT*, pages 4171–4186, 2019. 1

[26] Zhenglun Kong, Peiyan Dong, Xiaolong Ma, Xin Meng, Wei Niu, Mengshu Sun, Bin Ren, Minghai Qin, Hao Tang, and Yanzhi Wang. Spvit: Enabling faster vision transformers via soft token pruning. *arXiv preprint arXiv:2112.13890*, 2021. 2

[27] Alex H Lang, Sourabh Vora, Holger Caesar, Lubing Zhou, Jiong Yang, and Oscar Beijbom. Pointpillars: Fast encoders for object detection from point clouds. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12697–12705, 2019. 2, 5, 7

[28] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient convnets. *arXiv preprint arXiv:1608.08710*, 2016. 5

[29] Yanyu Li, Geng Yuan, Yang Wen, Eric Hu, Georgios Evangelidis, Sergey Tulyakov, Yanzhi Wang, and Jian Ren. Efficientformer: Vision transformers at mobilenet speed. *arXiv preprint arXiv:2206.01191*, 2022. 2

[30] Zhichao Li, Feng Wang, and Naiyan Wang. Lidar r-cnn: An efficient and universal 3d object detector. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7546–7555, 2021. 7

[31] Youwei Liang, Chongjian GE, Zhan Tong, Yibing Song, Jue Wang, and Pengtao Xie. EVit: Expediting vision transformers via token reorganizations. In *International Conference on Learning Representations*, 2022. 1, 2

[32] Zhidong Liang, Ming Zhang, Zehan Zhang, Xian Zhao, and Shiliang Pu. Rangercnn: Towards fast and accurate 3d object detection with range image representation. *arXiv preprint arXiv:2009.00206*, 2020. 2

[33] Ji Lin, Yongming Rao, Jiwen Lu, and Jie Zhou. Runtime neural pruning. *Advances in neural information processing systems*, 30, 2017. 2

[34] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the*

[35] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017. 5

[36] Chenxi Liu, Liang-Chieh Chen, Florian Schroff, Hartwig Adam, Wei Hua, Alan L Yuille, and Li Fei-Fei. Auto-deeplab: Hierarchical neural architecture search for semantic image segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 82–92, 2019. 2

[37] Chenxi Liu, Zhaoqi Leng, Pei Sun, Shuyang Cheng, Charles R Qi, Yin Zhou, Mingxing Tan, and Dragomir Anguelov. Lidarnas: Unifying and searching neural architectures for 3d point clouds. In *European Conference on Computer Vision*, pages 158–175. Springer, 2022. 2, 7

[38] Chenxi Liu, Barret Zoph, Maxim Neumann, Jonathon Shlens, Wei Hua, Li-Jia Li, Li Fei-Fei, Alan Yuille, Jonathan Huang, and Kevin Murphy. Progressive neural architecture search. In *Proceedings of the European conference on computer vision (ECCV)*, pages 19–34, 2018. 2

[39] Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search. *arXiv preprint arXiv:1806.09055*, 2018. 2

[40] Xingchao Liu, Mao Ye, Dengyong Zhou, and Qiang Liu. Post-training quantization with multiple points: Mixed precision without mixed precision. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 8697–8705, 2021. 2

[41] Zhuang Liu, Jianguo Li, Zhiqiang Shen, Gao Huang, Shoumeng Yan, and Changshui Zhang. Learning efficient convolutional networks through network slimming. In *Proceedings of the IEEE international conference on computer vision*, pages 2736–2744, 2017. 1, 2

[42] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10012–10022, 2021. 1, 3

[43] Zhuang Liu, Mingjie Sun, Tinghui Zhou, Gao Huang, and Trevor Darrell. Rethinking the value of network pruning. *arXiv preprint arXiv:1810.05270*, 2018. 1, 2

[44] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019. 6

[45] Jiageng Mao, Yujing Xue, Minzhe Niu, Haoyue Bai, Jiashi Feng, Xiaodan Liang, Hang Xu, and Chunjing Xu. Voxel transformer for 3d object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3164–3173, 2021. 2, 7

[46] Lingchen Meng, Hengduo Li, Bor-Chun Chen, Shiyi Lan, Zuxuan Wu, Yu-Gang Jiang, and Ser-Nam Lim. Adavit: Adaptive vision transformers for efficient image recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12309–12318, 2022. 2

[47] Yue Meng, Chung-Ching Lin, Rameswar Panda, Prasanna Sattigeri, Leonid Karlinsky, Aude Oliva, Kate Saenko, and Rogerio Feris. Ar-net: Adaptive frame resolution for efficient action recognition. In *European Conference on Computer Vision*, pages 86–104. Springer, 2020. 2

[48] Gregory P Meyer, Jake Charland, Shreyash Pandey, Ankit Laddha, Shivam Gautam, Carlos Vallespi-Gonzalez, and Carl K Wellington. Laserflow: Efficient and probabilistic object detection and motion forecasting. *IEEE Robotics and Automation Letters*, 6(2):526–533, 2020. 2

[49] Gregory P Meyer, Ankit Laddha, Eric Kee, Carlos Vallespi-Gonzalez, and Carl K Wellington. Lasernet: An efficient probabilistic 3d object detector for autonomous driving. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12677–12686, 2019. 2

[50] Gregory P Meyer and Niranjan Thakurdesai. Learning an uncertainty-aware object detector for autonomous driving. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 10521–10527. IEEE, 2020. 2

[51] Ishan Misra, Rohit Girdhar, and Armand Joulin. An end-to-end transformer model for 3d object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2906–2917, 2021. 2

[52] Markus Nagel, Mart van Baalen, Tijmen Blankevoort, and Max Welling. Data-free quantization through weight equalization and bias correction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1325–1334, 2019. 2

[53] Mahyar Najibi, Bharat Singh, and Larry S Davis. Autofocus: Efficient multi-scale inference. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9745–9755, 2019. 2

[54] Bowen Pan, Rameswar Panda, Yifan Jiang, Zhangyang Wang, Rogerio Feris, and Aude Oliva. Ia-red^2: Interpretability-aware redundancy reduction for vision transformers. *Advances in Neural Information Processing Systems*, 34:24898–24911, 2021. 1, 2

[55] Charles R Qi, Or Litany, Kaiming He, and Leonidas J Guibas. Deep hough voting for 3d object detection in point clouds. In *proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9277–9286, 2019. 2

[56] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 652–660, 2017. 2

[57] Charles R Qi, Yin Zhou, Mahyar Najibi, Pei Sun, Khoa Vo, Boyang Deng, and Dragomir Anguelov. Offboard 3d object detection from point cloud sequences. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6134–6144, 2021. 7

[58] Yongming Rao, Wenliang Zhao, Benlin Liu, Jiwen Lu, Jie Zhou, and Cho-Jui Hsieh. Dynamicvit: Efficient vision transformers with dynamic token sparsification. *Advances in neural information processing systems*, 34:13937–13949, 2021. 1, 2

[59] Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi. Xnor-net: Imagenet classification using binary convolutional neural networks. In *European conference on computer vision*, pages 525–542. Springer, 2016. 2

[60] Byungseok Roh, JaeWoong Shin, Wuhyun Shin, and Saehoon Kim. Sparse DETR: Efficient end-to-end object detection with learnable sparsity. In *International Conference on Learning Representations*, 2022. 2

[61] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 6, 16

[62] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4510–4520, 2018. 2, 6

[63] Shaoshuai Shi, Chaoxu Guo, Li Jiang, Zhe Wang, Jianping Shi, Xiaogang Wang, and Hongsheng Li Pv-rcnn. Point-voxel feature set abstraction for 3d object detection. 2020 ieee. In *CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10526–10535, 2020. 7

[64] Shaoshuai Shi, Xiaogang Wang, and Hongsheng Li. Pointr-cnn: 3d object proposal generation and detection from point cloud. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 770–779, 2019. 2

[65] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, Vijay Vasudevan, Wei Han, Jiquan Ngiam, Hang Zhao, Aleksei Timofeev, Scott Ettinger, Maxim Krivokon, Amy Gao, Aditya Joshi, Yu Zhang, Jonathon Shlens, Zhifeng Chen, and Dragomir Anguelov. Scalability in perception for autonomous driving: Waymo open dataset. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 5

[66] Pei Sun, Mingxing Tan, Weiyue Wang, Chenxi Liu, Fei Xia, Zhaoqi Leng, and Dragomir Anguelov. Swformer: Sparse window transformer for 3d object detection in point clouds. *arXiv preprint arXiv:2210.07372*, 2022. 2, 3, 7

[67] Pei Sun, Weiyue Wang, Yuning Chai, Gamaleldin Elsayed, Alex Bewley, Xiao Zhang, Cristian Sminchisescu, and Dragomir Anguelov. Rsn: Range sparse net for efficient, accurate lidar 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5725–5734, 2021. 2, 3, 7

[68] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International conference on machine learning*, pages 6105–6114. PMLR, 2019. 2

[69] Hidenori Tanaka, Daniel Kunin, Daniel L Yamins, and Surya Ganguli. Pruning neural networks without any data by iteratively conserving synaptic flow. *Advances in Neural Information Processing Systems*, 33:6377–6389, 2020. 2

[70] Yehui Tang, Kai Han, Yunhe Wang, Chang Xu, Jianyuan Guo, Chao Xu, and Dacheng Tao. Patch slimming for efficient vision transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12165–12174, 2022. 1, 2

[71] Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996. 5

[72] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 1, 3

[73] Andreas Veit and Serge Belongie. Convolutional networks with adaptive inference graphs. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 3–18, 2018. 2

[74] Haiyang Wang, Chen Shi, Shaoshuai Shi, Meng Lei, Sen Wang, Di He, Bernt Schiele, and Liwei Wang. Dsvt: Dynamic sparse voxel transformer with rotated sets. *arXiv preprint arXiv:2301.06051*, 2023. 2

[75] Xin Wang, Fisher Yu, Zi-Yi Dou, Trevor Darrell, and Joseph E Gonzalez. Skipnet: Learning dynamic routing in convolutional networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 409–424, 2018. 2

[76] Yulin Wang, Rui Huang, Shiji Song, Zeyi Huang, and Gao Huang. Not all images are worth 16x16 words: Dynamic transformers for efficient image recognition. *Advances in Neural Information Processing Systems*, 34:11960–11973, 2021. 2

[77] Yikai Wang, TengQi Ye, Lele Cao, Wenbing Huang, Fuchun Sun, Fengxiang He, and Dacheng Tao. Bridged transformer for vision and point cloud 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12114–12123, 2022. 2

[78] Kelvin Wong, Shenlong Wang, Mengye Ren, Ming Liang, and Raquel Urtasun. Identifying unknown instances for autonomous driving. In *Conference on Robot Learning*, pages 384–393. PMLR, 2020. 2

[79] Zuxuan Wu, Tushar Nagarajan, Abhishek Kumar, Steven Rennie, Larry S Davis, Kristen Grauman, and Rogerio Feris. Blockdrop: Dynamic inference paths in residual networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8817–8826, 2018. 2

[80] Yifan Xu, Zhijie Zhang, Mengdan Zhang, Kekai Sheng, Ke Li, Weiming Dong, Liqing Zhang, Changsheng Xu, and Xing Sun. Evo-vit: Slow-fast token evolution for dynamic vision transformer. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 2964–2972, 2022. 2

[81] Yan Yan, Yuxing Mao, and Bo Li. Second: Sparsely embedded convolutional detection. *Sensors*, 18(10):3337, 2018. 2, 16

[82] Bin Yang, Wenjie Luo, and Raquel Urtasun. Pixor: Real-time 3d object detection from point clouds. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 7652–7660, 2018. 2

[83] Le Yang, Yizeng Han, Xi Chen, Shiji Song, Jifeng Dai, and Gao Huang. Resolution adaptive networks for efficient inference. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2369–2378, 2020. 2

[84] Mao Ye, Chengyue Gong, Lizhen Nie, Denny Zhou, Adam Klivans, and Qiang Liu. Good subnetworks provably exist: Pruning via greedy forward selection. In *International Conference on Machine Learning*, pages 10820–10830. PMLR, 2020. 1, 2

[85] Mao Ye and Yan Sun. Variable selection via penalized neural network: a drop-out-one loss approach. In *International Conference on Machine Learning*, pages 5620–5629. PMLR, 2018. 5

[86] Mao Ye, Lemeng Wu, and Qiang Liu. Greedy optimization provably wins the lottery: Logarithmic number of winning tickets is enough. *Advances in Neural Information Processing Systems*, 33:16409–16420, 2020. 2

[87] Hongxu Yin, Arash Vahdat, Jose M Alvarez, Arun Mallya, Jan Kautz, and Pavlo Molchanov. A-vit: Adaptive tokens for efficient vision transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10809–10818, 2022. 1, 2, 6, 16

[88] Tianwei Yin, Xingyi Zhou, and Philipp Krahenbuhl. Center-based 3d object detection and tracking. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11784–11793, 2021. 5, 7

[89] Hao Yu and Jianxin Wu. A unified pruning framework for vision transformers. *arXiv preprint arXiv:2111.15127*, 2021. 2

[90] Jiahui Yu and Thomas S Huang. Universally slimmable networks and improved training techniques. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1803–1811, 2019. 2

[91] Yanan Zhang, Jiaxin Chen, and Di Huang. Cat-det: Contrastively augmented transformer for multi-modal 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 908–917, 2022. 2

[92] Yifan Zhang, Qingyong Hu, Guoquan Xu, Yanxin Ma, Jianwei Wan, and Yulan Guo. Not all points are equal: Learning highly efficient point-based detectors for 3d lidar point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18953–18962, 2022. 2, 7

[93] Ritchie Zhao, Yuwei Hu, Jordan Dotzel, Chris De Sa, and Zhiru Zhang. Improving neural network quantization without retraining using outlier channel splitting. In *International conference on machine learning*, pages 7543–7552. PMLR, 2019. 2

[94] Shuchang Zhou, Yuxin Wu, Zekun Ni, Xinyu Zhou, He Wen, and Yuheng Zou. Dorefa-net: Training low bitwidth convolutional neural networks with low bitwidth gradients. *arXiv preprint arXiv:1606.06160*, 2016. 2

[95] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points. *arXiv preprint arXiv:1904.07850*, 2019. 5, 14

[96] Yin Zhou, Pei Sun, Yu Zhang, Dragomir Anguelov, Jiyang Gao, Tom Ouyang, James Guo, Jiquan Ngiam, and Vijay Vasudevan. End-to-end multi-view fusion for 3d object detection in lidar point clouds. In *Conference on Robot Learning*, pages 923–932. PMLR, 2020. 3

[97] Yin Zhou and Oncel Tuzel. Voxelnet: End-to-end learning for point cloud based 3d object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4490–4499, 2018. 2

[98] Zixiang Zhou, Xiangchen Zhao, Yu Wang, Panqu Wang, and Hassan Foroosh. Centerformer: Center-based transformer for 3d object detection. *arXiv preprint arXiv:2209.05588*, 2022. 2

[99] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable DETR: Deformable transformers for end-to-end object detection. In *International Conference on Learning Representations*, 2021. 2

[100] Yichen Zhu, Yuqin Zhu, Jie Du, Yi Wang, Zhicai Ou, Feifei Feng, and Jian Tang. Make a long image short: Adaptive token length for vision transformers. *arXiv preprint arXiv:2112.01686*, 2021. 2

[101] Zhuofan Zong, Kunchang Li, Guanglu Song, Yali Wang, Yu Qiao, Biao Leng, and Yu Liu. Self-slimming vision transformer, 2022. 2

# A. Appendix

## A.1. Additional Details

### A.1.1 Non-Uniform Sparsity Loss

For each token, we determine whether or not it lies within a ground-truth bounding box. If the $i$-th token falls inside a bounding box, the corresponding heatmap's value is defined as

$$m_i = \exp\left(-\frac{(x_i - l_x)^2 + (y_i - l_y)^2}{2\sigma^2}\right), \tag{18}$$

where $x_i$ and $y_i$ are the $xy$-coordindates of the token, $l_x$ and $l_y$ are the $xy$-coordinates of the box's center, and $\sigma$ is a hyper-parameter that controls the smoothness of the heatmap. For tokens that do not lie within a bounding box, the heatmap's value is set to zero. We create the heatmap for each object class, and we take the maximize over all heatmaps to obtain the final class-agnostic heatmap used by our non-uniform sparsity loss. The non-uniform sparsity loss is similar to focal loss [95], and it is defined as follows:

$$\mathcal{L}_s = -\sum_{l=1}^{L} \sum_{i \in \mathcal{K}} \frac{1}{|\mathcal{K}|} \left[ (1 - s_{l,i})^\alpha \log(s_{l,i}) \mathbb{I}_{m_i \geq 1-\epsilon} + (1 - m_i)^\gamma s_i^\alpha \log(1 - s_{l,i}) \mathbb{I}_{m_i < 1-\epsilon} \right], \tag{19}$$

where $\mathcal{K} = \{i : k_{0:l-1,i} = 1\}$ is the set of tokens that have not been halted before the $l$-th layer, $s_{l,i}$ is the score for the $i$-th token at the $l$-th layer, $\alpha = 2$ and $\gamma = 4$ are hyper-parameters, and $\epsilon = 10^{-4}$ is used to improve numerical stability.

### A.1.2 Analyzing the Pseudo-Gradient

In Section 5, we claim the following:

$$\Delta_i \approx \frac{\partial \mathcal{L}(q_1, q_2)}{\partial s_i} + O(u), \tag{20}$$

where

$$\Delta_i := \mathcal{L}(\tilde{q}_1, \tilde{q}_2) - \mathcal{L}(q_1, q_2) \tag{21}$$

is the difference in the detection loss when the $i$-th token is halted instead of being forwarded in a single layer network. In other words, we claim that the (pseudo)-gradient of $\mathcal{L}(\tilde{q}_1, \tilde{q}_2)$ with respect to $s_i$ provided by our proposed EDF and the STE is a reasonable proxy of $\Delta_i$.

To prove this claim, we begin by computing

$$\frac{\partial \mathcal{L}(q_1, q_2)}{\partial s_i} = \left\langle \frac{\partial \mathcal{L}(q_1, q_2)}{\partial q_1}, \frac{\partial q_1}{\partial s_i} \right\rangle + \left\langle \frac{\partial \mathcal{L}(q_1, q_2)}{\partial q_2}, \frac{\partial q_2}{\partial s_i} \right\rangle. \tag{22}$$

Recall that $q_1 = (\mathbf{1} - k) \circ f$ and $q_2 = k \circ \phi_2(\text{WSA}(\phi_1(f), s \circ k), f)$. Using the definition of the STE,

$$\frac{\partial q_1}{\partial s_i} = -\frac{\partial(k \circ f)}{\partial s_i} = -\mathbf{1}_i \circ f \tag{23}$$

and

$$\frac{\partial q_2}{\partial s_i} = k \circ \frac{\partial \phi_2(\text{WSA}(\phi_1(f), s \circ k), f)}{\partial s_i} + \frac{\partial k}{\partial s_i} \circ \phi_2(\text{WSA}(\phi_1(f), s \circ k), f) \tag{24}$$

$$= k \circ \frac{\partial \phi_2(\text{WSA}(\phi_1(f), s \circ k), f)}{\partial s_i} + \mathbf{1}_i \circ \phi_2(\text{WSA}(\phi_1(f), s \circ k), f). \tag{25}$$

Next, let us compute $\Delta_i$. Using Taylor series approximation,

$$\Delta_i \approx \left\langle \frac{\partial \mathcal{L}(q_1, q_2)}{\partial q_1}, \tilde{q}_1 - q_1 \right\rangle + \left\langle \frac{\partial \mathcal{L}(q_1, q_2)}{\partial q_2}, \tilde{q}_2 - q_2 \right\rangle. \tag{26}$$

Recall that $\tilde{q}_1 = (\mathbf{1} - k - \mathbf{1}_i) \circ f$ and $\tilde{q}_2 = (k + \mathbf{1}_i) \circ \phi_2(\text{WSA}(\phi_1(f), s \circ (k + \mathbf{1}_i)), f)$; therefore,

$$\tilde{q}_1 - q_1 = -\mathbf{1}_i \circ f \tag{27}$$

and

$$\tilde{q}_2 - q_2 = (k + \mathbf{1}_i) \circ \phi_2(\text{WSA}(\phi_1(f), s \circ (k + \mathbf{1}_i)), f) - k \circ \phi_2(\text{WSA}(\phi_1(f), s \circ k), f) \tag{28}$$
$$= k \circ \phi_2(\text{WSA}(\phi_1(f), s \circ (k + \mathbf{1}_i)), f) - k \circ \phi_2(\text{WSA}(\phi_1(f), s \circ k), f) \tag{29}$$
$$+ \mathbf{1}_i \circ \phi_2(\text{WSA}(\phi_1(f), s \circ (k + \mathbf{1}_i)), f). \tag{30}$$

Again, using Taylor series approximation,

$$k \circ \phi_2(\text{WSA}(\phi_1(f), s \circ (k + \mathbf{1}_i)), f) - k \circ \phi_2(\text{WSA}(\phi_1(f), s \circ k), f) \approx k \circ \frac{\partial \phi_2(\text{WSA}(\phi_1(f), s \circ k), f)}{\partial s_i} \tag{31}$$

and

$$\mathbf{1}_i \circ \phi_2(\text{WSA}(\phi_1(f), s \circ (k + \mathbf{1}_i)), f) \approx \mathbf{1}_i \circ \phi_2(\text{WSA}(\phi_1(f), s \circ k), f) + \mathbf{1}_i \circ \frac{\partial \phi_2(\text{WSA}(\phi_1(f), s \circ k), f)}{\partial s_i} \tag{32}$$

As a result,

$$\tilde{q}_2 - q_2 \approx k \circ \frac{\partial \phi_2(\text{WSA}(\phi_1(f), s \circ k), f)}{\partial s_i} + \mathbf{1}_i \circ \phi_2(\text{WSA}(\phi_1(f), s \circ k), f) + \mathbf{1}_i \circ \frac{\partial \phi_2(\text{WSA}(\phi_1(f), s \circ k), f)}{\partial s_i}. \tag{33}$$

Comparing Eq. (23) to Eq. (27) and Eq. (25) to Eq. (33), we see that

$$\Delta_i - \frac{\partial \mathcal{L}(q_1, q_2)}{\partial s_i} \approx \left\langle \frac{\partial \mathcal{L}(q_1, q_2)}{\partial q_2}, \mathbf{1}_i \circ \frac{\partial \phi_2(\text{WSA}(\phi_1(f), s \circ k), f)}{\partial s_i} \right\rangle \tag{34}$$

$$= \frac{\partial \mathcal{L}(q_1, q_2)}{\partial q_{2,i}} \frac{\partial \phi_2(\text{WSA}(\phi_1(f), s \circ k), f)_i}{\partial s_i} \tag{35}$$

$$= \frac{\partial \mathcal{L}(q_1, q_2)}{\partial q_{2,i}} \frac{\partial \phi_2(\text{WSA}(\phi_1(f), s \circ k), f)_i}{\partial \text{WSA}(\phi_1(f), s \circ k)_i} \frac{\partial \text{WSA}(\phi_1(f), s \circ k)_i}{\partial (s \circ k)_i} s_i, \tag{36}$$

where

$$\frac{\partial \phi_2(\text{WSA}(\phi_1(f), s \circ k), f)_i}{\partial s_i} = \frac{\partial \phi_2(\text{WSA}(\phi_1(f), s \circ k), f)_i}{\partial \text{WSA}(\phi_1(f), s \circ k)_i} \frac{\partial \text{WSA}(\phi_1(f), s \circ k)_i}{\partial (s \circ k)_i} \frac{\partial (s \circ k)_i}{\partial s_i} \tag{37}$$

and

$$\frac{\partial (s \circ k)_i}{\partial s_i} = \frac{\partial s_i k_i}{\partial s_i} = k_i + s_i = s_i \tag{38}$$

using the definition of the STE and the fact that $k_i = 0$. In general, the derivative of $\mathcal{L}$ and $\phi_2$ is bounded since the parameter space of the network is bounded (due to the weight decay) and the operators inside the network are Lipschitz continuous. However, $\partial \text{WSA}(\phi_1(f), s \circ k)_i / \partial (s \circ k)_i$ can be singular when all the element of $s \circ k$ are zero. We argue that in our analysis, we can still treat this term as bounded for two reasons. Firstly, in practice, we add an $\epsilon$ to the denominator of the WSA to prevent numeric instability, which makes the gradient bounded even if all the elements of $s \circ k$ are zero. Secondly, we employ gradient clipping to enforce a bound on the gradient. All of this combine, we have

$$\Delta_i \approx \frac{\partial \mathcal{L}(q_1, q_2)}{\partial s_i} + O(u). \tag{39}$$

That is, the approximation error of the pseudo-gradient is proportional to $s_i$ thanks to the usage of weighted attention. Furthermore, since $k_i = 0$, we have $|s_i| < u$ where the threshold $u$ is in general a very small value. This demonstrates that our pseudo-gradient provides useful information for updating the halting module.

| Speed Up/Sparsity | Method | Vehicle | | Pedestrian | | Cyclist | |
|---|---|---|---|---|---|---|---|
| | | AP/APH L1 | AP/APH L2 | AP/APH L1 | AP/APH L2 | AP/APH L1 | AP/APH L2 |
| 1.00/0.00 | Original | 76.2/75.7 | 67.7/67.2 | 79.9/71.4 | 72.7/64.8 | 67.7/66.3 | 65.2/63.8 |
| 1.07/0.00 | Width scale | 75.4/74.9 | 66.9/66.5 | 79.5/70.7 | 72.2/64.0 | 65.6/64.1 | 63.0/61.6 |
| 1.16/0.00 | Num head scale | 75.5/75.0 | 67.0/66.6 | 79.4/70.8 | 72.0/64.1 | 65.5/63.9 | 63.0/61.5 |
| 1.28/0.48 | AViT$_{adapted}$ [87] | 71.9/71.3 | 63.4/62.9 | 76.8/67.9 | 69.1/60.9 | 63.1/61.6 | 60.7/59.3 |
| 1.21/0.75 | Ours | 76.1/75.6 | 67.8/67.3 | 79.4/70.7 | 72.1/64.0 | 67.0/65.6 | 64.4/63.1 |
| 1.13/0.00 | Width scale | 73.8/73.3 | 65.3/64.8 | 78.2/69.0 | 70.6/62.1 | 61.7/60.2 | 59.3/57.9 |
| 1.45/0.00 | Num head scale | 73.8/73.3 | 65.4/64.9 | 78.2/68.8 | 70.7/62.1 | 62.0/60.3 | 59.6/58.0 |
| 1.39/0.57 | AViT$_{adapted}$ [87] | 70.3/69.7 | 61.9/61.4 | 76.2/67.2 | 68.4/60.1 | 60.8/59.3 | 58.5/57.0 |
| 1.37/0.82 | Ours | 76.1/75.6 | 67.7/67.2 | 79.9/71.5 | 72.6/64.7 | 67.4/66.1 | 64.8/63.6 |
| 1.18/0.00 | Width scale | 70.4/69.8 | 62.0/61.5 | 74.5/64.3 | 66.7/57.4 | 54.9/52.9 | 52.8/50.8 |
| 1.45/0.00 | Num head scale | 73.8/73.3 | 65.4/64.9 | 78.2/68.8 | 70.7/62.1 | 62.0/60.3 | 59.6/58.0 |
| 1.55/0.72 | AViT$_{adapted}$ [87] | 70.1/69.6 | 61.7/61.3 | 76.3/67.5 | 68.5/60.4 | 61.6/60.1 | 59.2/57.8 |
| 1.52/0.89 | Ours | 75.4/74.9 | 67.0/66.5 | 79.7/71.5 | 72.4/64.7 | 67.1/65.7 | 64.5/63.2 |

Table 5. Efficiency and accuracy trade-off. We report the relative backbone speed-up and the average sparsity across all the attention layers.

## A.2. Additional Experiment Details and Results

### A.2.1 Setup and Implementation Details

We use a mixed strategy to decide the halting threshold. We specific an upper and lower token score quantile (denoted as $\alpha_u$ and $\alpha_l$) and enforce that the sparsity of each layer varies within $[\alpha_l, \alpha_u]$. To achieve this, the final threshold is given by clamping the pre-specific threshold $u$ within $[Q(\alpha_l), Q(\alpha_u)]$ where $Q(\alpha_l)$ and $Q(\alpha_u)$ denote the score corresponding to the $\alpha_l$ and $\alpha_u$ quantile, respectively. We enforce such a constrain because we observe that the distribution of scores can vary considerably for different scenes during the early stages of training and selecting the threshold in this way helps to stabilize training. In Table 1, the sparsity is bounded between 80% and 90% for the first halting module, 90% and 99% for the second halting module, and the default value of $u$ is 0.01. The following technique can be used to identify $u$ for a new dataset/model: train a model for a short period, then select $u$ such that it is higher than the score of most foreground voxels and less than the score of most background voxels.

### A.2.2 Efficiency and Accuracy Trade-off

For the baselines, we vary the latent dimension of the attention mechanism by $\{16, 12, 8, 4\}$, and we vary the number of attention head by $\{8, 6, 5, 4\}$. We adapt AViT from [87], but we apply our token recycling to improve the performance. Also, we adjust the number of input token features for the halting module from 1 to 32 as this improves performance while having a negligible impact on latency. Table 5 summarizes the results. Overall, we observe that our method significantly improves over other model scaling approaches as well as AViT.

### A.2.3 The SST$^{++}_{halt}$ Architecture

For our SST$^{++}_{halt}$ architecture, we use a U-Net [61] and a single layer MLP as the first and second halting module. We find that the latent features of the U-Net contain useful semantic information. To reuse those features, we fuse the token features with the U-Net's features by applying a linear transformation and sums the features. Furthermore, we add the U-Net's feature map to the BEV feature map. To leverage the latency savings provided by halting tokens, SST$^{++}_{halt}$ uses an extra convolutional block in the detection head for a total of two convolutional blocks. The first convolutional block contains four convolutional layers. The second convolutional block contains four convolutional layer where the first layer has a stride of 2. All the layers use a kernel size of 3. Afterwards, we use the Feature Pyramid Network [34] employed by SECOND [81] to fuse the two scales of the BEV feature map and make predictions.