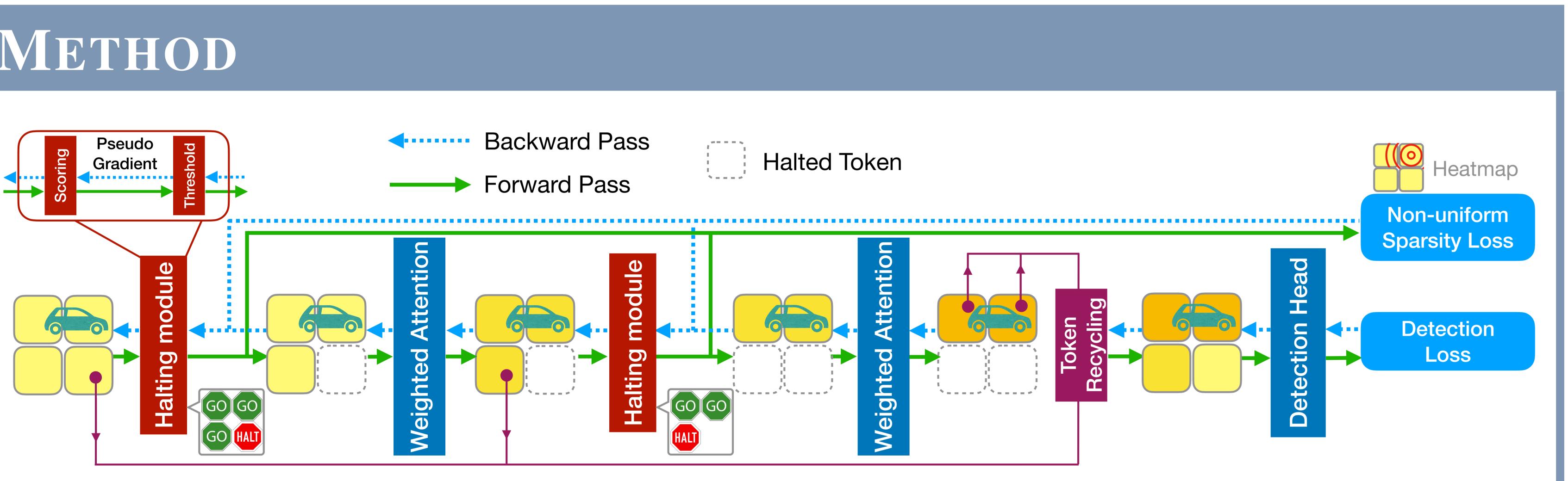


- Transformers have achieved state-of-the-art results in several computer vision tasks.
- by dynamically halting tokens based on their contribution to the detection task.

#### **METHOD**



- Our method is built upon SST, which is a transformer-based LiDAR object detector.

- During training, we define a pseudo-gradient to back-propagate through the halting module.
- To improve learning of the halting module, we leverage a weighted attention mechanism, which weights the attention given to each token based on their score.
- The whole network is trained end-to-end using a detection loss and a sparsity loss.

## **EFFICIENT TRANSFORMER-BASED 3D OBJECT DETECTION WITH DYNAMIC TOKEN HALTING** Mao Ye<sup>1,2</sup>, Gregory P. Meyer<sup>2</sup>, Yuning Chai<sup>2</sup>, and Qiang Liu<sup>1</sup> <sup>1</sup>The University of Texas at Austin <sup>2</sup>Cruise LLC

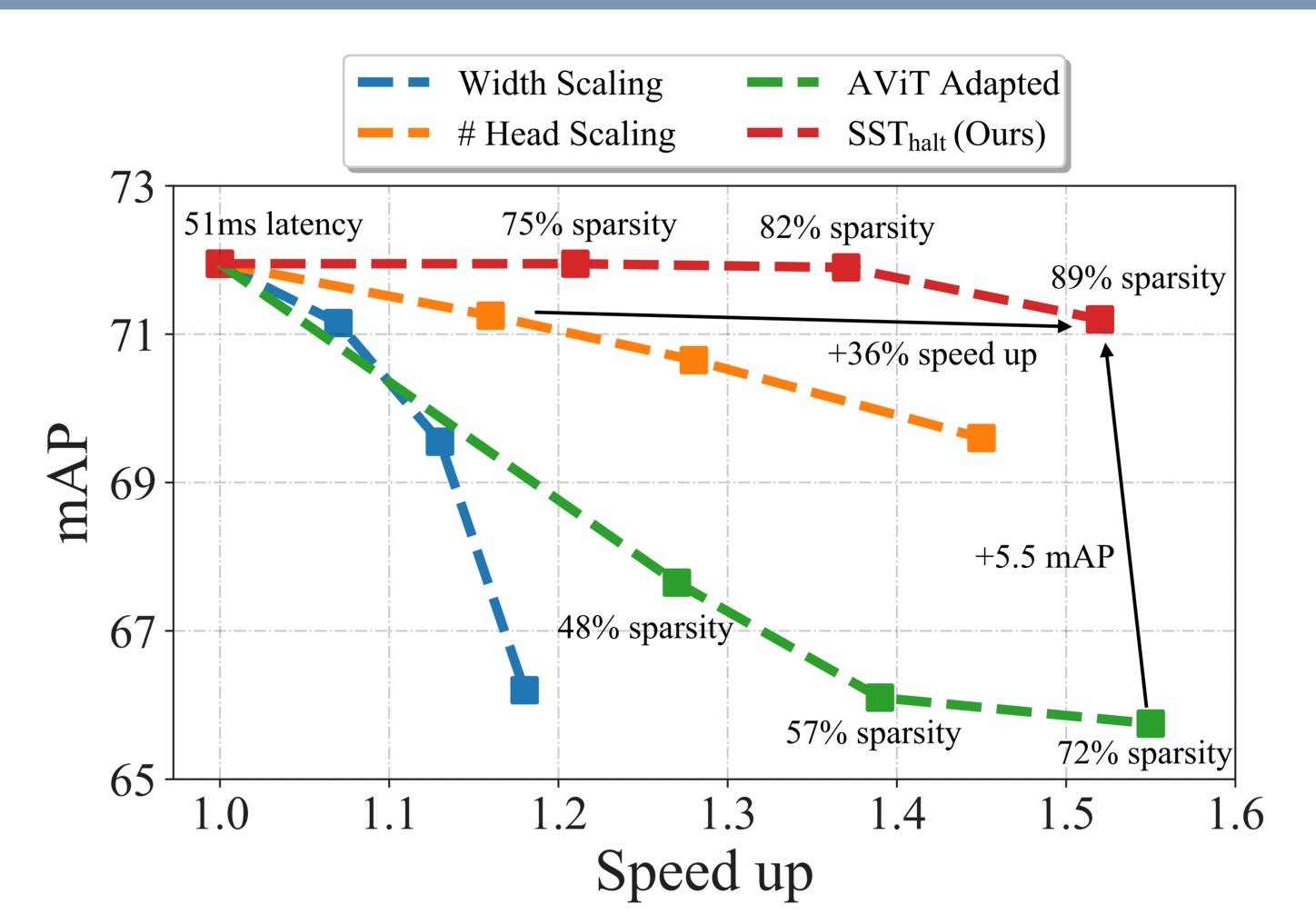
• However, transformers suffer from high-latency due the quadratic complexity of attention. • In this work, we propose a method for accelerating transformer-based 3D object detectors

• As a result, our method significantly improves the Pareto frontier of latency versus accuracy.

• First, the input LiDAR point cloud is voxelized, and each voxel is treated as a token. • Before each layer, a halting module scores the tokens and halts those with a low score.

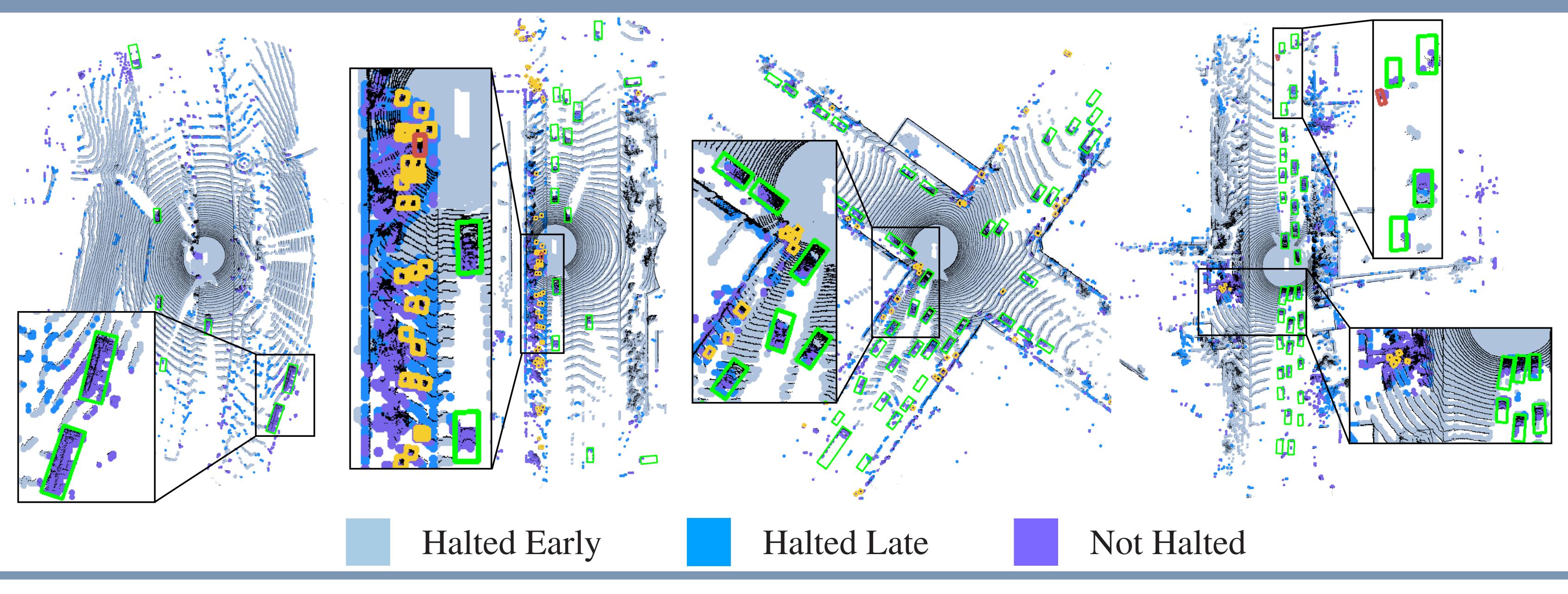
• After the last attention layer, halted tokens are recycled and forwarded to the detection head.

## **EFFICIENCY VS. ACCURACY DETECTION PERFORMANCE**



When compared to other model scaling approaches, our method provides the best efficiency and accuracy trade-off.

### VISUALIZATIONS



# CIUISE

WOD Validation Set	
--------------------	--

Method	Vehicle APH L2	Pedestrian APH L2	Cyclist APH L2
PointPillar	63.1	50.3	59.9
<b>PV-RCNN</b>	68.4	57.6	64.0
RangeDet	63.6	63.9	62.1
Lidăr R-CNN	67.9	51.7	64.4
CenterPoint	67.5	57.9	
RSN	65.5	63.7	
SST	65.1	61.7	
SWFormer	68.8	64.9	
SST <sup>++</sup> <sub>halt</sub> (Ours)	69.0	66.5	66.0

We leverage the latency savings provided by our method to improve the performance of SST while maintaining its runtime.