# SparseViT: Revisiting Activation Sparsity for Efficient High-Resolution Vision Transformer

Xuanyao Chen[1,2,#],    Zhijian Liu[4,#],    Haotian Tang[4],

Li Yi[1,3],    Hang Zhao[1,3],    Song Han[4]

[1]: Shanghai Qi Zhi Institute    [2]: Fudan University    [3]: Tsinghua University    [4]: MIT

# Background

**High-resolution images have become ubiquitous!**



**Mobile Vision**

iPhone 14 Pro Max
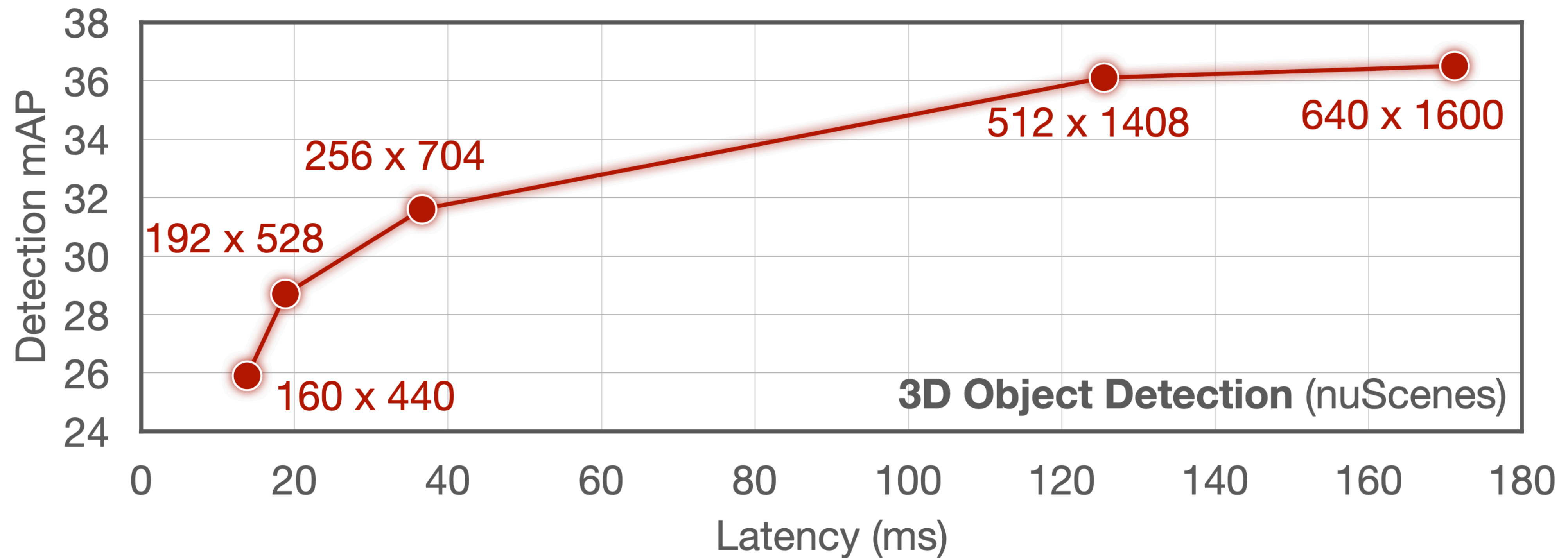
**48MP**

**Mixed Reality**

Apple Vision Pro

**23MP**

**Autonomous Driving**

Waymo Driver

**2MP** x **29**

# Background

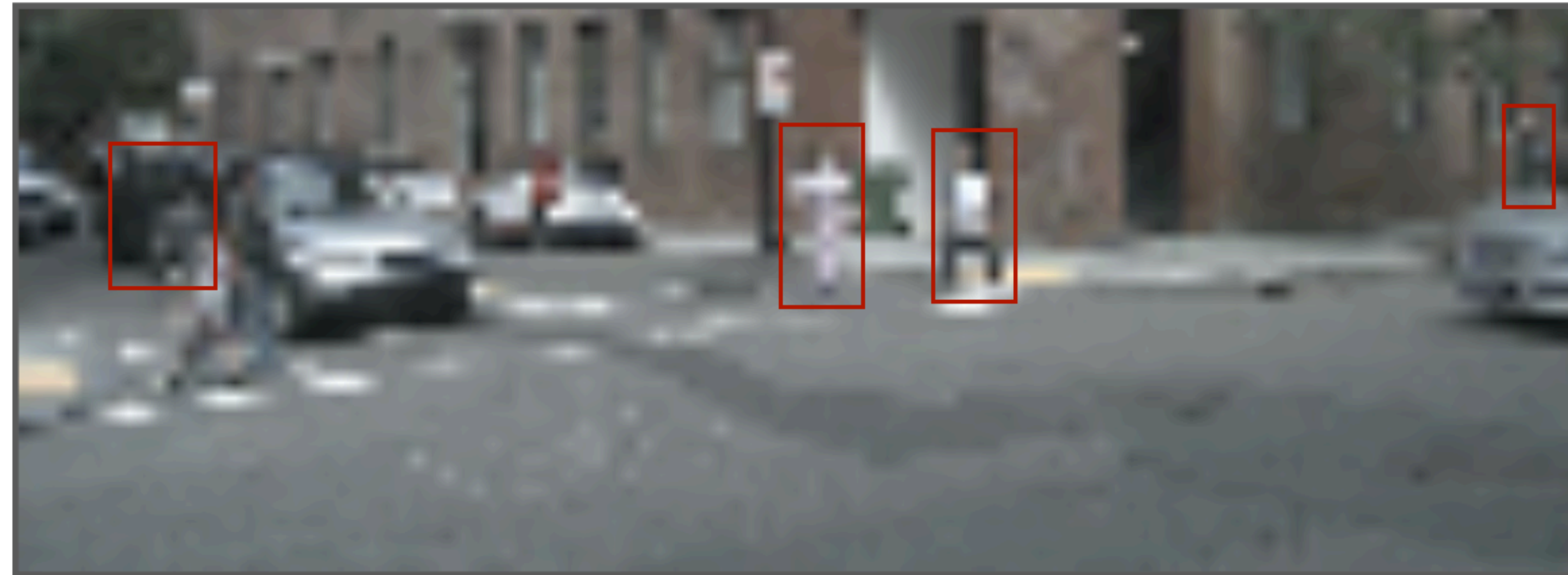**Higher resolutions deliver better accuracy but also increase computation cost!**

# Activation Pruning

**Sparse high-resolution features are better than dense low-resolution ones.**



**Uniform Resizing**

**Low** Resolution (0.5X)
**Dense** Pixels (100%)

**Activation Pruning**

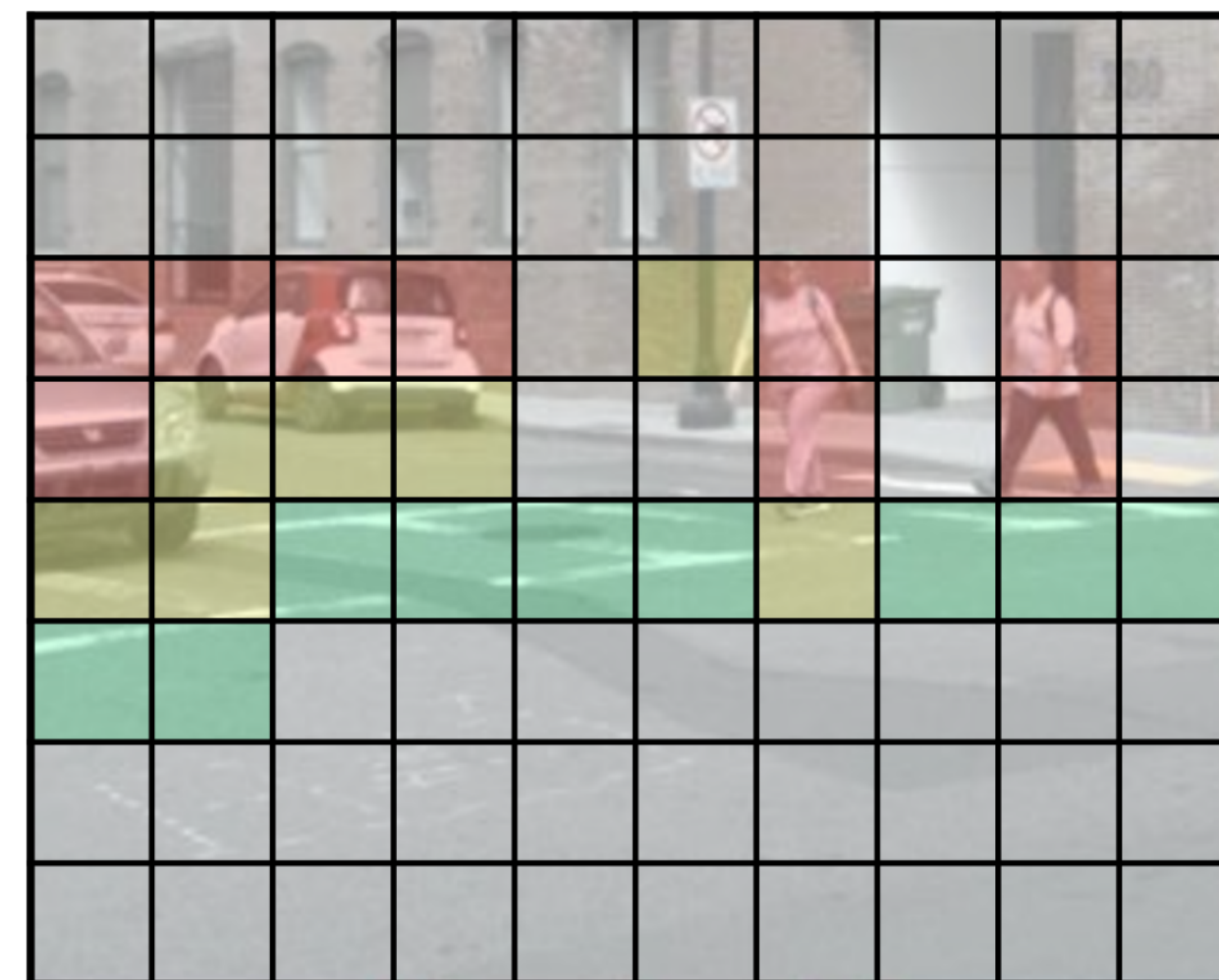**High** Resolution (1X)
**Sparse** Pixels (25%)
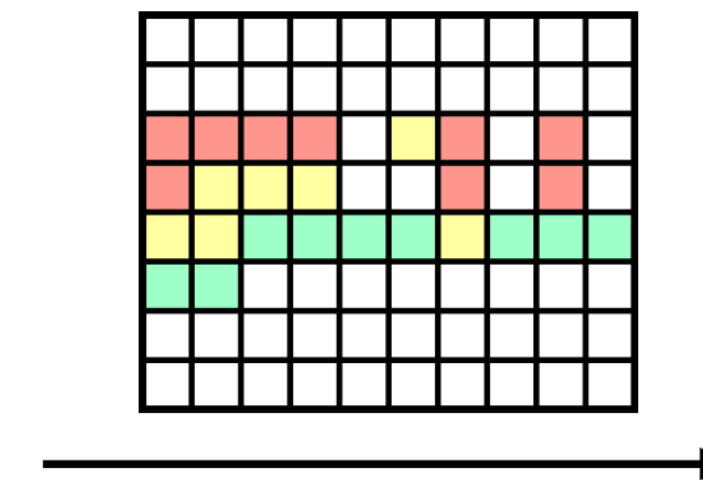
# SparseViT — Sparse Vision Transformers

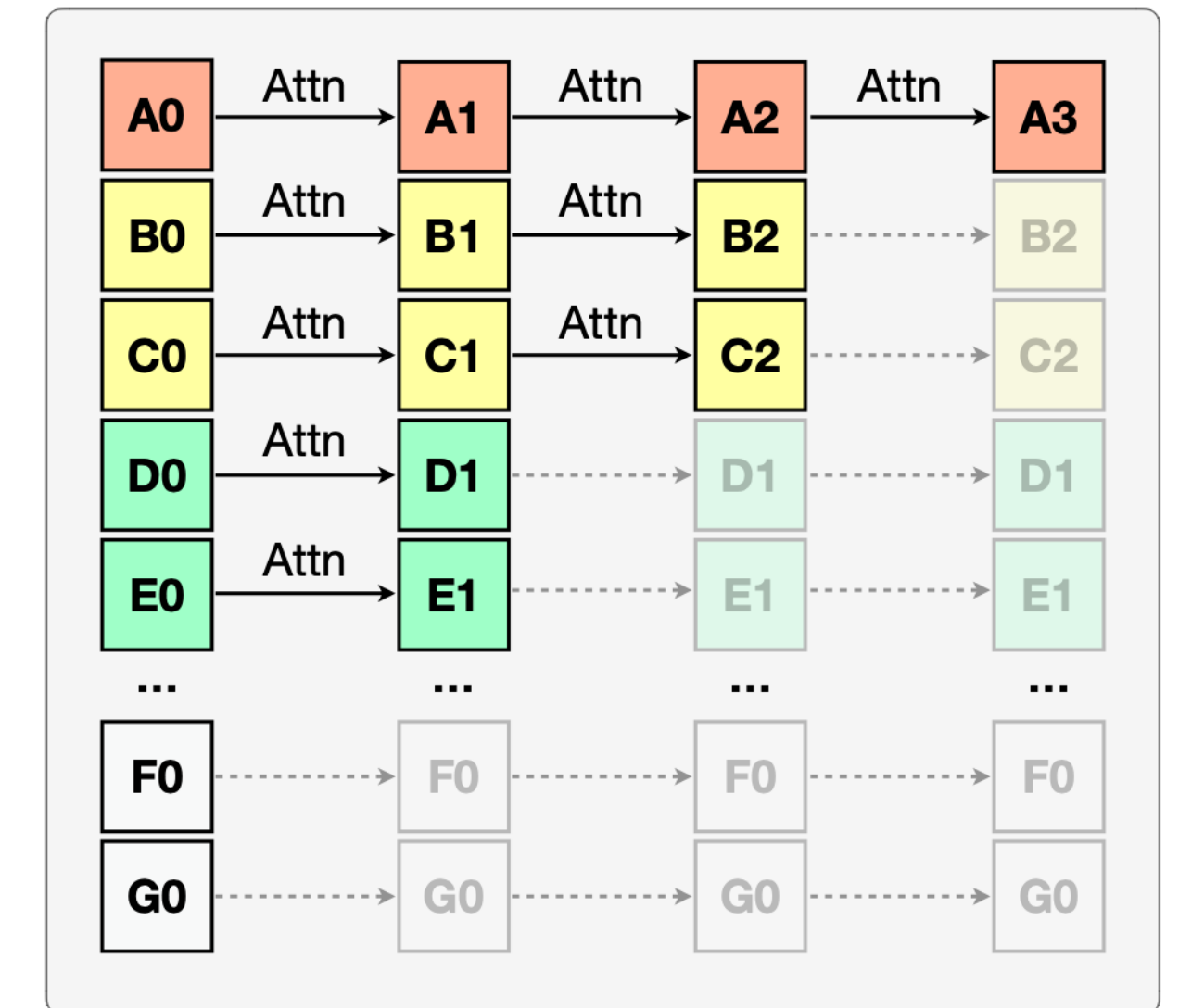**Step I. Window activation pruning (with non-uniform sparsity)**



**Input Image**
(or Input Feature Map)

**Window Importance**
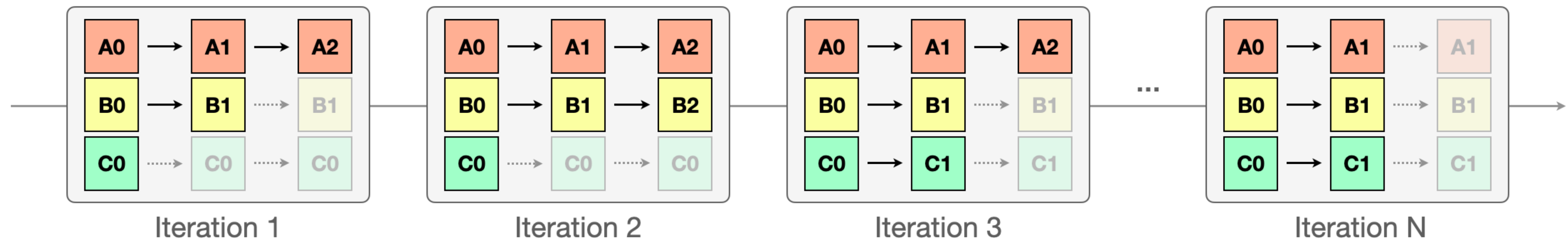(L2 Activation Magnitude)

**Sorting & Gathering**

**Sparse Window Attention**

# SparseViT — Sparse Vision Transformers

## Step II. Sparsity-aware adaptation

**Goal**: Assess the model's accuracy under different activation sparsity settings both **efficiently** and **accurately**.
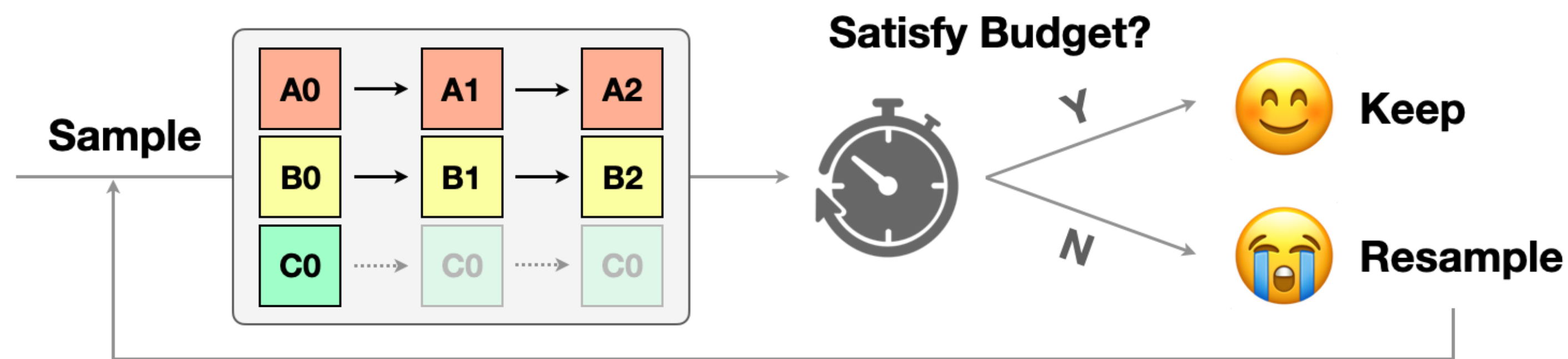


As the original model is trained with **only dense activations**, we **improve its sparsity awareness** by finetuning it with **randomly sampled** layerwise activation sparsity configurations at each training iteration.
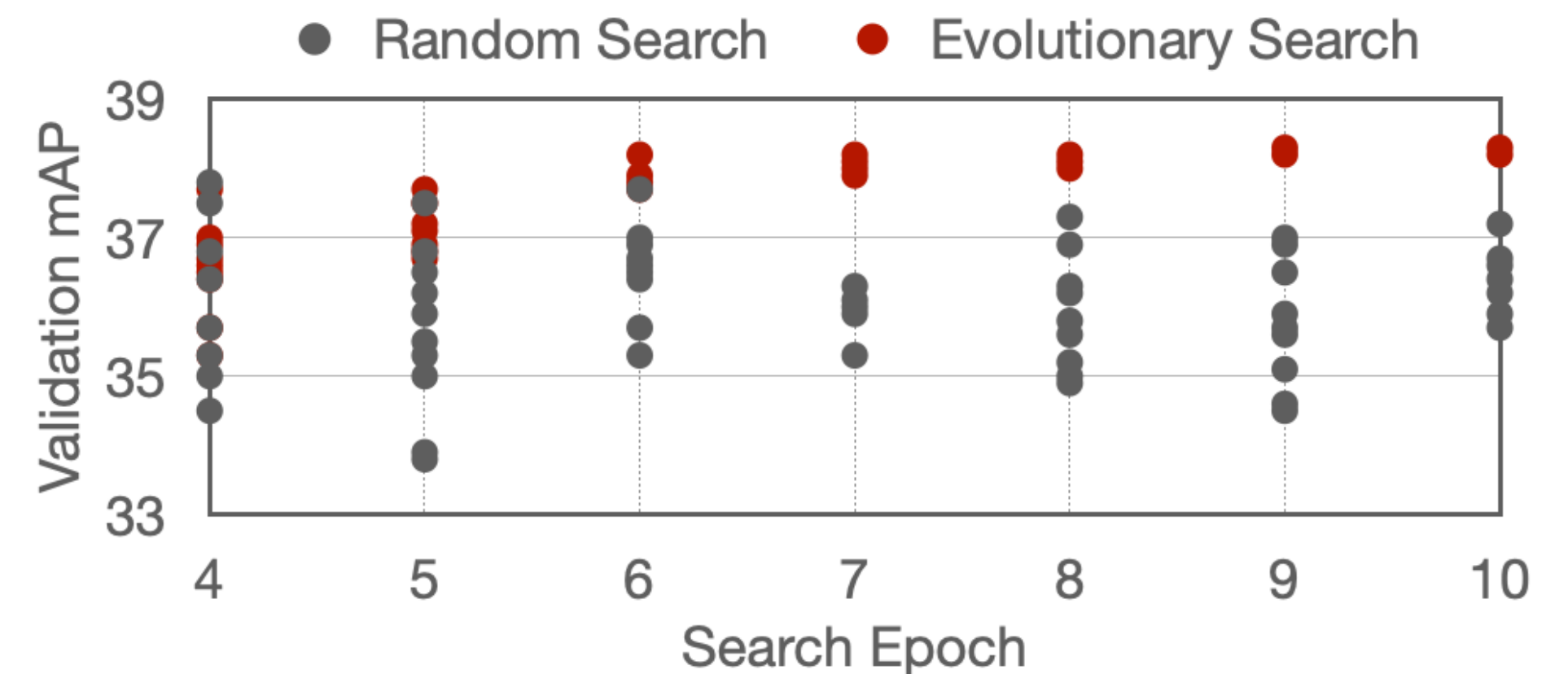
# SparseViT — Sparse Vision Transformers

## Step III. Resource-constrained search

**Goal**: Discover the **optimal** layerwise activation sparsity configuration **under a given latency budget.**
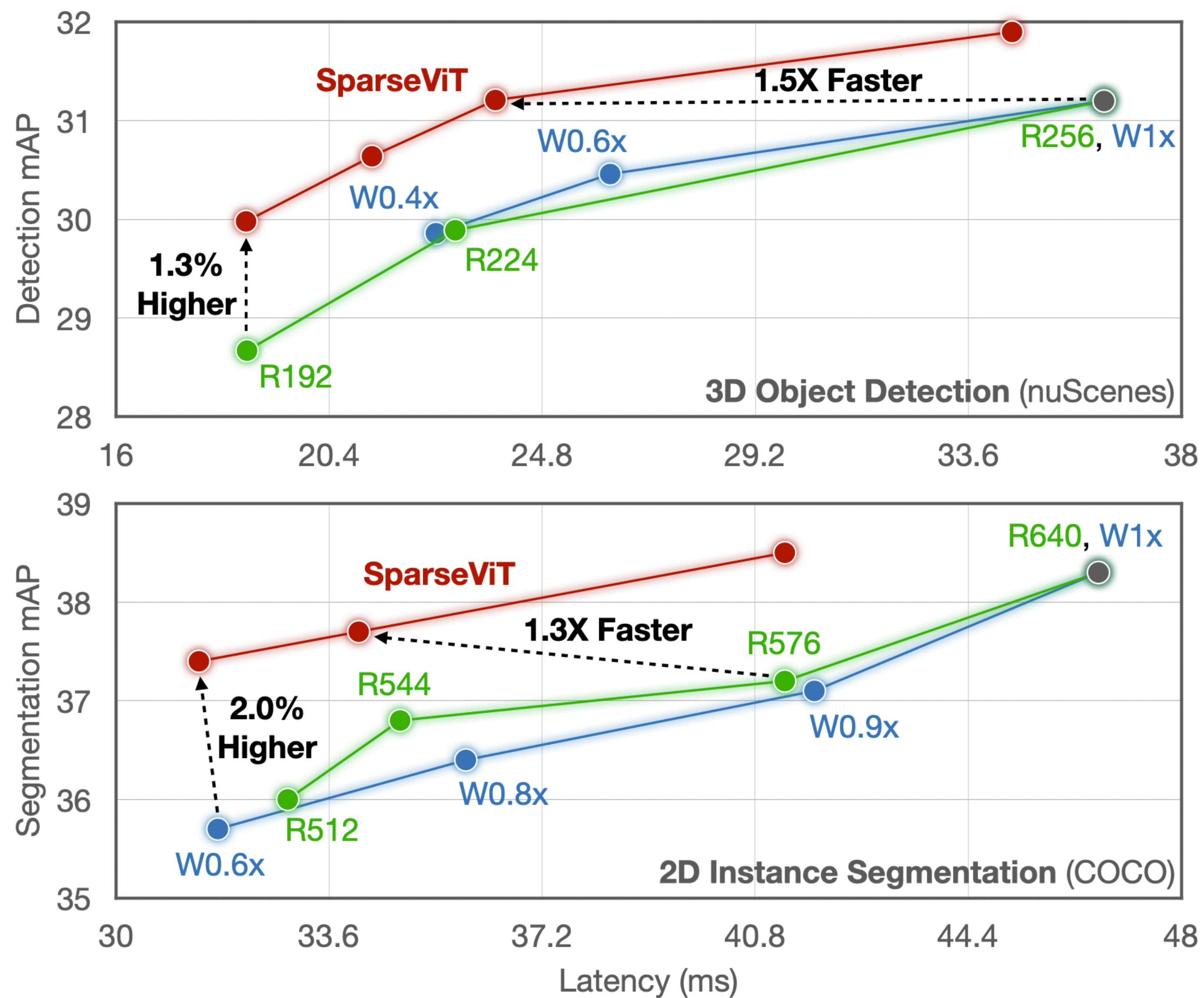


We enforce the latency constraint using **rejection sampling** (repeated resampling until satisfaction).

Evolutionary search is **sample-efficient!**

# Results



3D Object Detection (nuScenes) — Detection mAP vs Latency (ms)
- SparseViT (red), with points at R192, R224, R256/W1x, W0.4x, W0.6x
- 1.5X Faster, 1.3% Higher

2D Instance Segmentation (COCO) — Segmentation mAP vs Latency (ms)
- SparseViT (red), with points at R512, R544, R576, R640/W1x, W0.6x, W0.8x, W0.9x
- 1.3X Faster, 2.0% Higher

# Visualizations



SparseViT learns to prune **irrelevant background** windows while retaining **informative foreground** ones!