

# FlatFormer: Flattened Window Attention for Efficient Point Cloud Transformer

Zhijian Liu<sup>1,#</sup>, Xinyu Yang<sup>2,#</sup>, Haotian Tang<sup>1</sup>, Shang Yang<sup>3</sup>, Song Han<sup>1</sup>

<sup>1</sup>: MIT

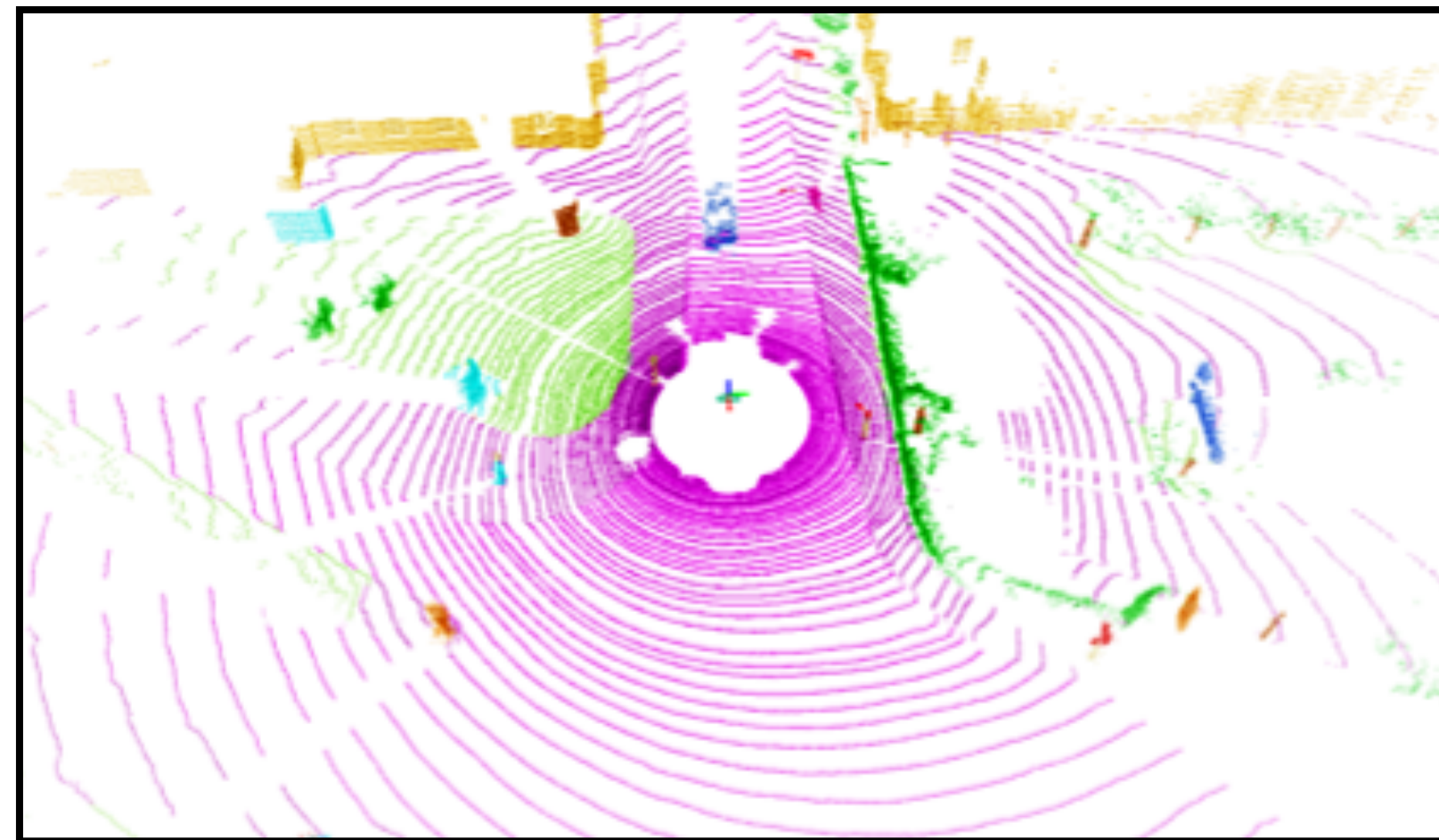
<sup>2</sup>: Shanghai Jiao Tong University

<sup>3</sup>: Tsinghua University



# Background

Point cloud deep learning is widely used in real-world applications



**Data:** 3D Point Cloud

**Autonomous  
Driving**








**Augmented Reality  
/ Mixed Reality**



**Applications**

# Background

State-of-the-art point cloud models are based on sparse convolution

Method					Metrics									
Date	Name	Modalities	Map data	External data	mAP	mATE (m)	mASE (1-IOU)	mAOE (rad)	mAVE (m/s)	mAAE (1-acc)	NDS	PKL *	FPS (Hz)	Stats
		Lidar	All	All										
> 2023-03-10	FocalFormer3D-TT/	Lidar	no	no	0.705	0.243	0.238	0.321	0.200	0.131	0.739	0.564	n/a	
> 2023-03-06	Real-Aug	Lidar	no	no	0.689	0.237	0.227	0.320	0.184	0.131	0.734	0.597	n/a	
> 2022-10-31	LinK	Lidar	no	no	0.698	0.238	0.229	0.312	0.234	0.136	0.734	0.575	n/a	
> 2022-05-30	LargeKernel-L	Lidar	no	no	0.688	0.244	0.230	0.312	0.241	0.132	0.728	0.581	n/a	
> 2022-11-11	SphereFormer	Lidar	no	no	0.685	0.242	0.233	0.353	0.191	0.131	0.728	0.616	n/a	
> 2022-09-27	MDRNet-L	Lidar	no	no	0.684	0.242	0.229	0.319	0.227	0.129	0.728	0.581	n/a	
> 2021-07-30	DAA_AVP	Lidar	no	no	0.697	0.237	0.229	0.376	0.248	0.124	0.727	0.937	n/a	
> 2022-08-18	MGTANet	Lidar	no	no	0.675	0.250	0.233	0.308	0.193	0.124	0.727	0.538	n/a	
> 2023-03-10	FocalFormer3D	Lidar	no	no	0.687	0.254	0.242	0.340	0.218	0.126	0.726	0.586	n/a	
> 2021-04-19	CenterPoint-VID	Lidar	no	no	0.674	0.255	0.235	0.339	0.233	0.128	0.718	0.555	n/a	



# Background

## Sparse convolution requires specialized system support to run on GPUs

mit-han-lab / torchsparse Public

Code Issues 29 Pull requests 4 Discussions Actions Projects 1 Security Insights Settings

master 7 branches 4 tags

Go to file Add file Code

About [MLSys'22] TorchSparse: Efficient Point Cloud Inference Engine torchsparse.mit.edu acceleration point-cloud pytorch

Releases 4 v1.4.0 Latest on Jun 24, 2021 + 3 releases

File	Update	Time
.github	Update pre-commit to avoid installing clang-format locally ...	2 years ago
docs	Update FAQ.md (#103)	2 years ago
examples	Add optimizations introduced in MLSys paper (#150)	8 months ago
torchsparse	Fix save_dir in the tune function (#190)	3 months ago
.gitignore	Reformat codebase and add pre-commit (#81)	2 years ago
.pre-commit-config.yaml	Include paper, website information (#177)	6 months ago
LICENSE	Reformat codebase and add pre-commit (#81)	2 years ago
README.md	Include paper, website information (#177)	6 months ago
requirements.txt	Support cached_property for python version < 3.8 (#171)	7 months ago
setup.cfg	Reformat codebase and add pre-commit (#81)	2 years ago
setup.py	Reformat codebase and add pre-commit (#81)	2 years ago

```
from torch import nn
from torchsparse import nn as tsnn
from torchsparse import SparseTensor

# Model definition
model = nn.Sequential(
    tsnn.Conv3d(4, 16, 3, stride=2, padding=1),
    tsnn.BatchNorm(16),
    tsnn.ReLU(True),
).cuda()

# Tensor definition
x = SparseTensor(feats, coords).cuda()

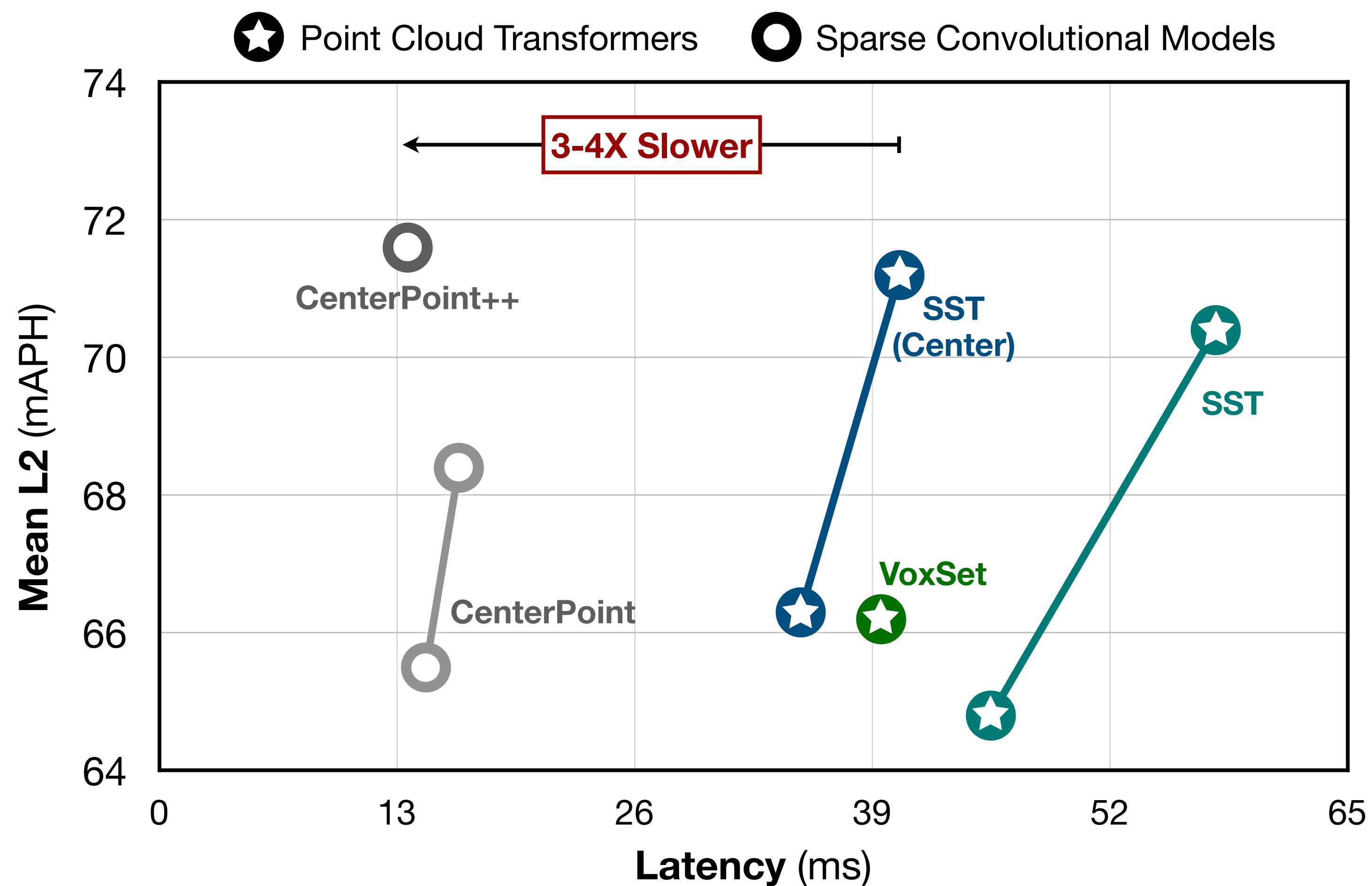
# Forward and backward pass
y = model(x)
loss = criterion(y, ...)
loss.backward()
```

Installation: `pip install --upgrade git+https://github.com/mit-han-lab/torchsparse.git`



# Point Cloud Transformers

Achieve comparable accuracy but lag far in latency (3-4X slower)



# Point Cloud Transformers

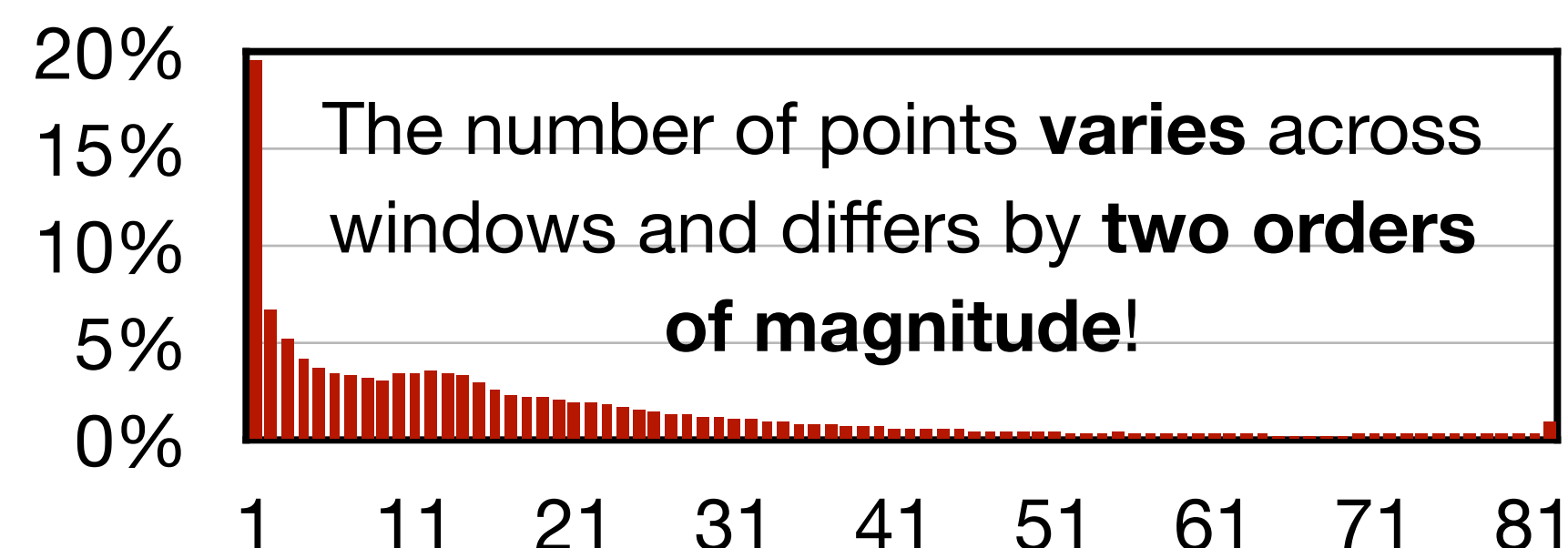
Achieve comparable accuracy but lag far in latency (3-4X slower)

**Global PCTs:** Apply MHSA globally across the entire point cloud.

The runtime of global PCTs grows **quadratically** as the number of points grows. The model takes almost **1 second** to run with **32k** input points.

Guo et al., “PCT: Point Cloud Transformer”, CVM 2021

**Window PCTs:** Apply MHSA to a set of non-overlapping windows.



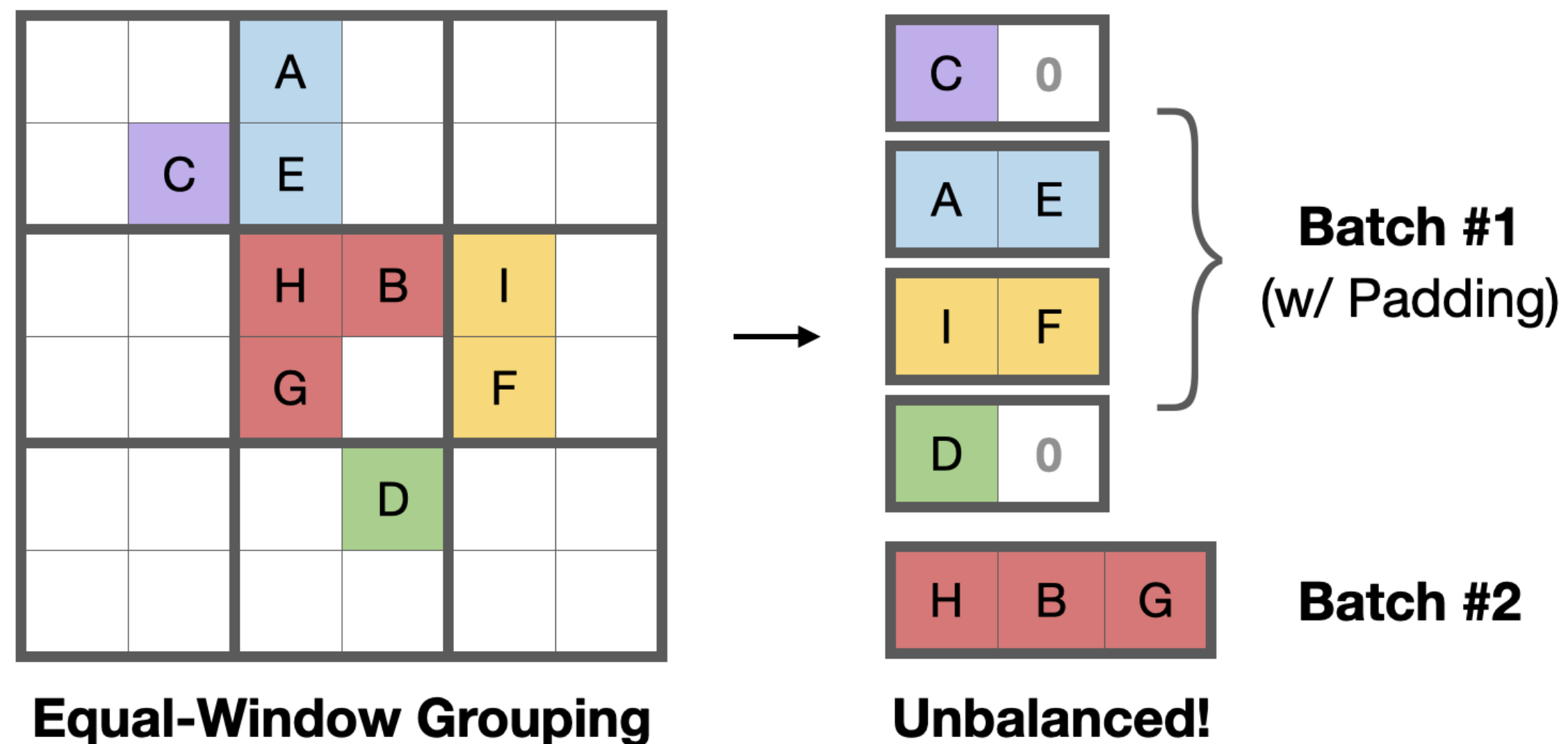
Window PCTs suffer from the **padding and partitioning overhead** due to the unbalanced workload across windows.

Fan et al., “Embracing Single Stride 3D Object Detector with Sparse Transformer”, CVPR 2022

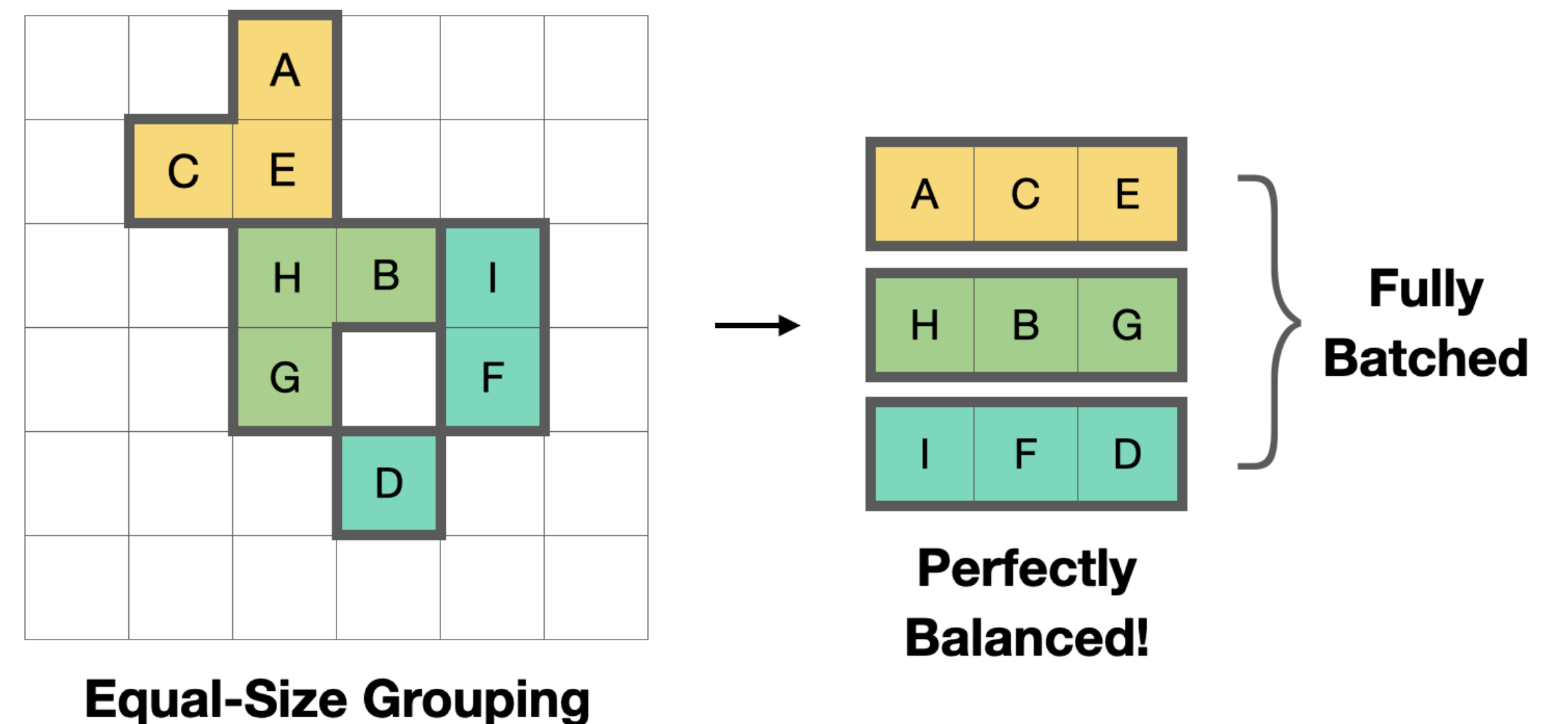
# FlatFormer: Flattened Window Attention

Equal-size grouping trades spatial proximity for computational regularity

Equal-window grouping maintains perfect **spatial proximity** but breaks the **computational regularity**.



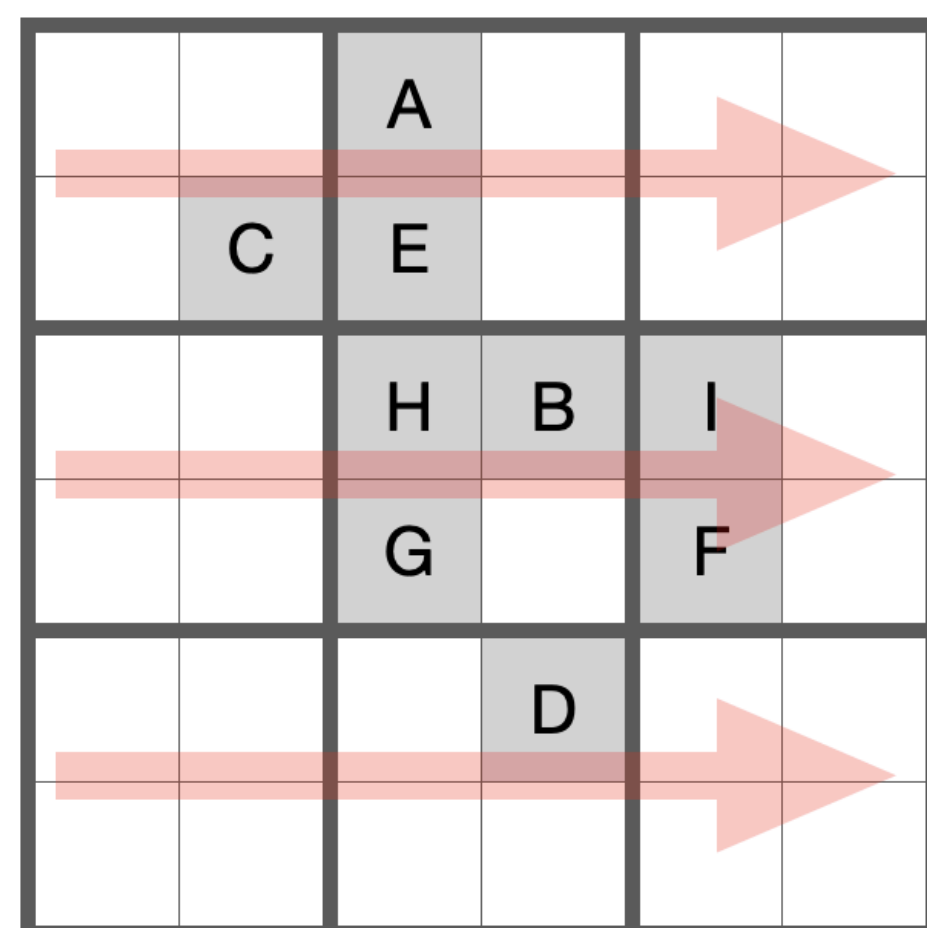
Equal-size grouping ensures balanced **computation workload** but cannot guarantee the **geometric locality**.



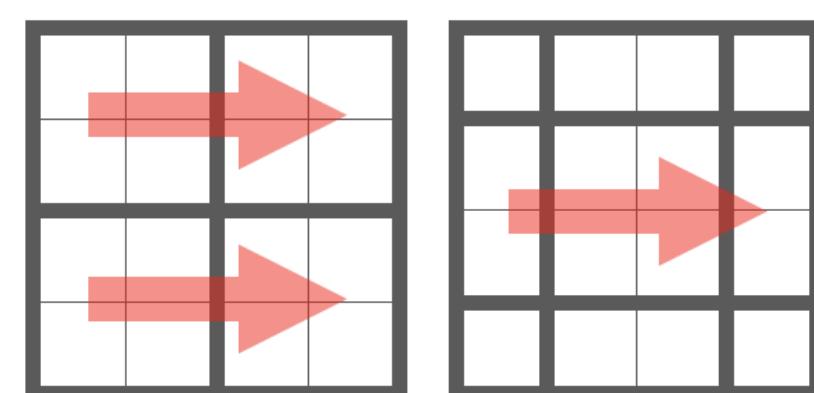


# FlatFormer: Flattened Window Attention

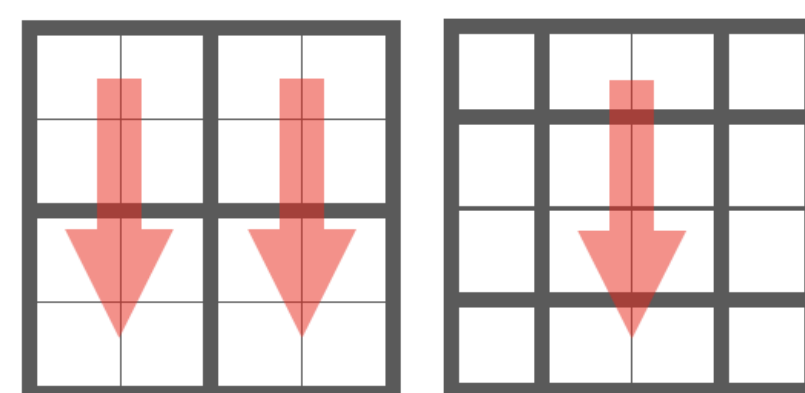
Window-based sorting preserves geometric locality after flattening



Sort all points first by **window coordinates** and then by **local coordinates** within the window.



Window Shift

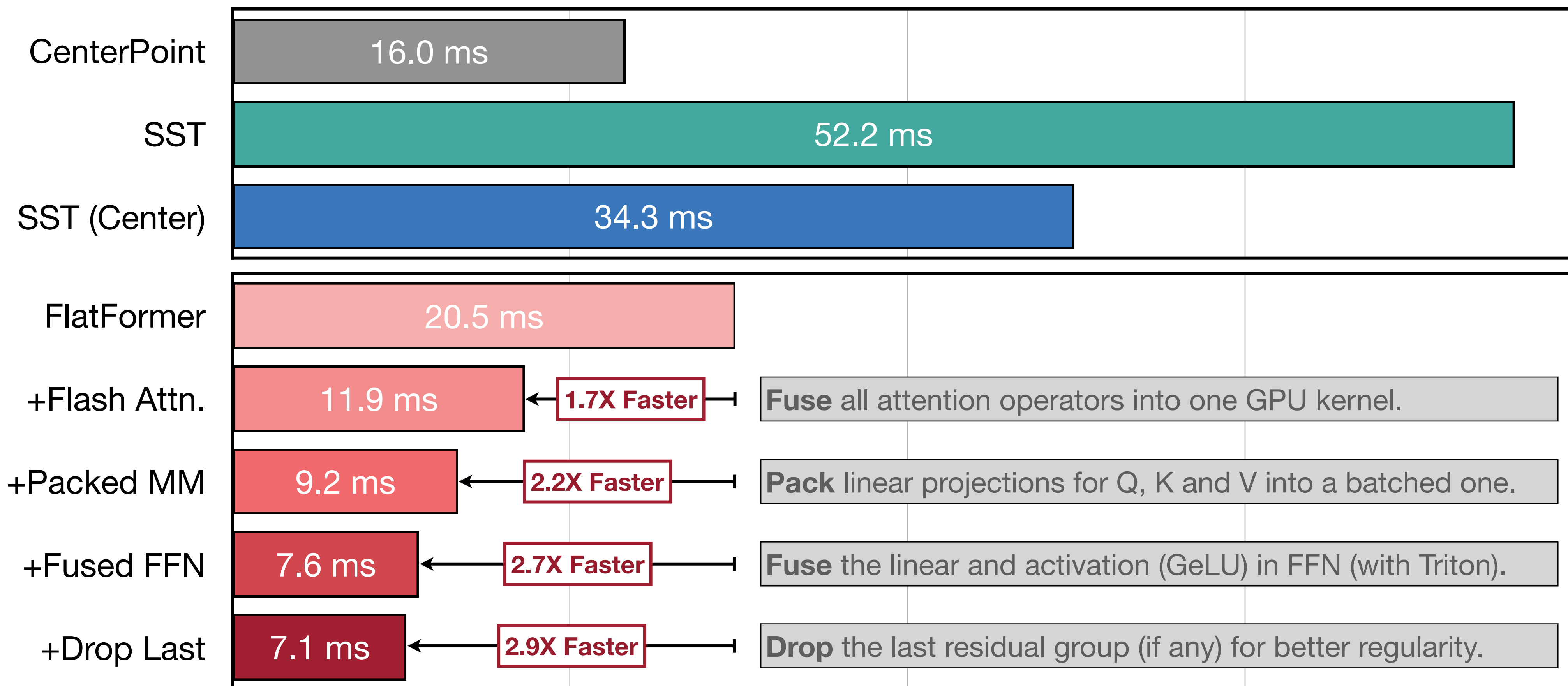


Alternate Sorting Axis

Sorting Criterion	Window Shifting	Axis Alternation	Mean L2 mAPH
<b>Random</b>	—	—	57.8
<b>Point</b>	—	Y	60.4
Window	Y	<b>N</b>	61.1
Window	<b>N</b>	Y	61.2
<b>Window</b>	<b>Y</b>	<b>Y</b>	<b>61.7</b>

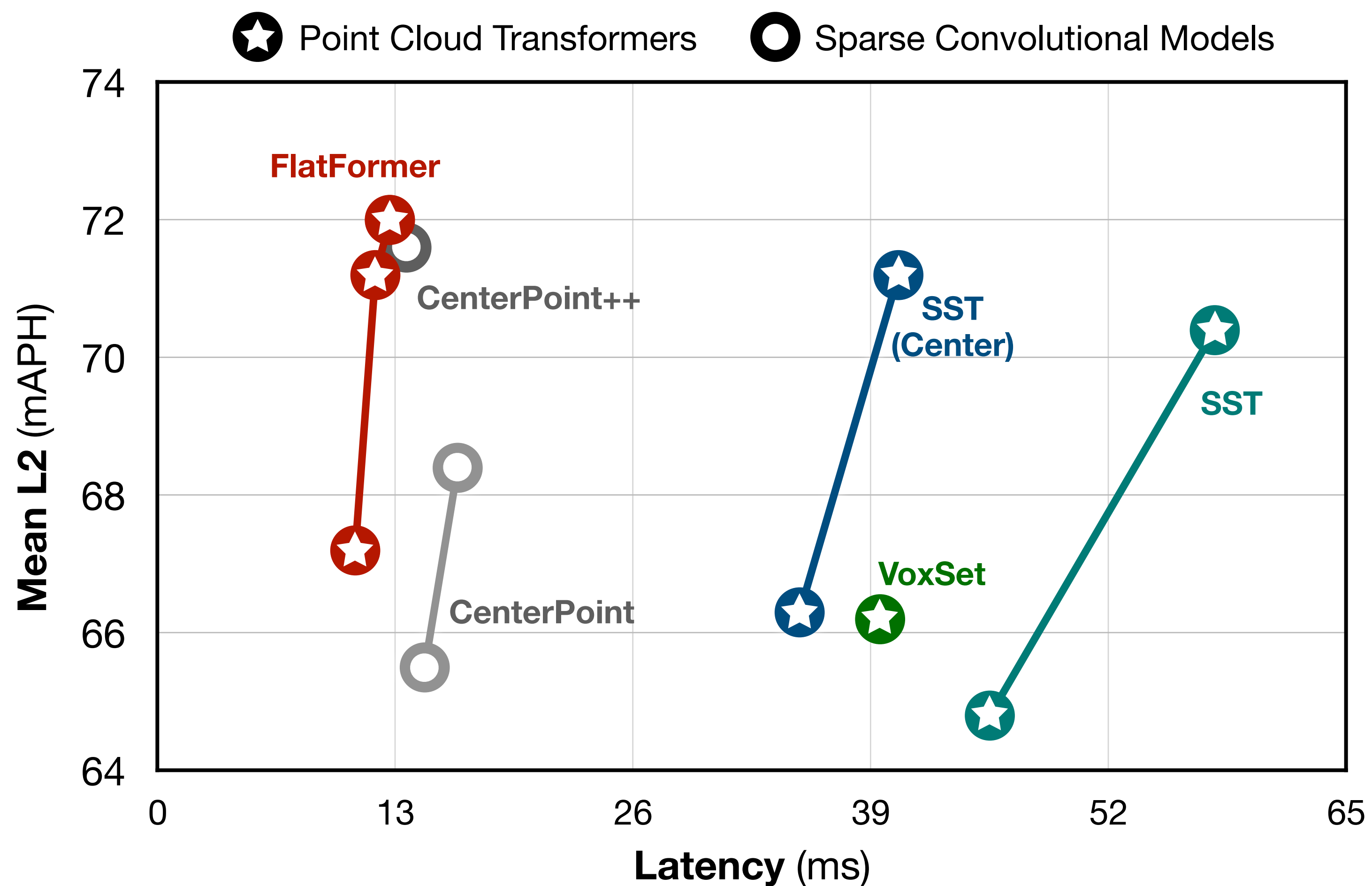
# Efficient Implementation

FlatFormer can benefit from existing system optimizations for transformers



# Results: 3D Object Detection on Waymo

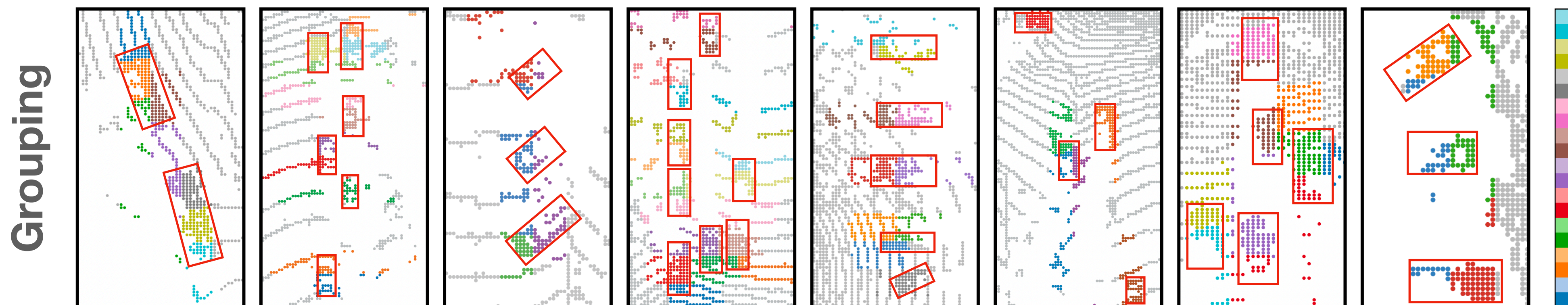
FlatFormer closes the latency gap between PCTs and SpConv-based models



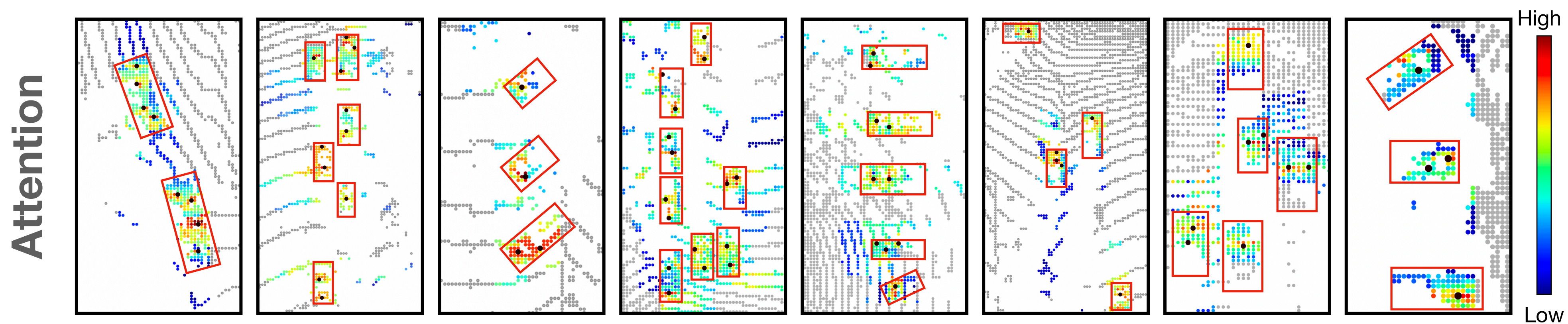


# Analysis: Grouping & Attention

Equal-size grouping is **mostly spatially local**.



Attention learns to **suppress outlier points**.





# Deployment on Edge GPUs

First point cloud transformer that achieves real-time performance on edge GPUs!

