### 30th Asia and South Pacific Design Automation Conference (ASPDAC '25) Tokyo, Japan

### APTO: Accelerating Serialization-Based Point Cloud Transformers with Position-Aware Pruning

**Qichu Sun**, Rui Meng, Haishuang Fan, Fangqiang Ding, Linxi Lu, Jingya Wu, Xiaowei Li, Guihai Yan 1. State Key Laboratory of Processors, Institute of Computing Technology, Chinese Academy of Sciences 2. University of Chinese Academy of Sciences 3. University of Edinburgh 4. YUSUR Technology Co., Ltd. sunqichu22z@ict.ac.cn

Jan. 21, 2025













# **Point Cloud Processing**

Point cloud is a representation of 3D data, containing valuable geometry & color information

**Autonomous Driving** 

### **Robotic Perception**

AR / VR



### Accurate, real-time and energy-efficient point cloud processing is crucial

# **Point Cloud Transformers**

#### **Point-Based Models**



- Farthest Point Sampling (FPS) & k-Nearest
  Neighbors (kNN) for down-sampling
- Attention mechanisms in local windows for feature computation

#### **Serialization-Based Models**



- Organizes points onto a directed curve (e.g. z-curve)
- 3D sparse convolution for down-sampling
- Larger attention windows based on the curve

Repeated point access leads to redundant computation & memory use, while window size restricts accuracy

Regular memory access, less redundant computation, and better accuracy

# **Motivation: Performance Bottleneck Analysis**

#### Inference time breakdown:

- Octformer: 46% SpConv & 27% attention
- PTv3: 32% SpConv & 50% attention
- Failed to meet real-time requirements



### SpConv & attention are two key bottlenecks of serialization-based models

Peng-Shuai Wang. 2023. Octformer: Octree-based transformers for 3d point clouds. ACM Transactions on Graphics (TOG) 42, 4 (2023), 1–11.
 Xiaoyang Wu, Li Jiang, Peng-Shuai Wang, Zhijian Liu, Xihui Liu, Yu Qiao, Wanli Ouyang, Tong He, and Hengshuang Zhao. 2024. Point Transformer V3: Simpler Faster Stronger. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 4840–4851.

### **Motivation**

# Challenges

- Chlg.1: Inefficient Neighbor Search in SpConv
- Chlg.2: Useless Attention Computation
- Chlg.3: Data Dependencies in Softmax

# > Optimizations

- Opt.1: Serialization-Based Parallel Neighbor Search
- Opt.2: Position-Aware Attention Pruning
- Opt.3: Fine-Grained Attention Dataflow

# Chlg.1: Inefficient Neighbor Search in SpConv

### > Kernel mapping in SpConv

- Stride > 1 for down-sampling
- Stride = 1 for positional encoding



Neighbor search is the most complex step

# Chlg.1: Inefficient Neighbor Search in SpConv

### > Prior works: voxel merging

• Merge and compare all voxels (>10<sup>5</sup>) serially



### Results in enormous operations and lacks parallelism

[1] Yujun Lin, Zhekai Zhang, Haotian Tang, Hanrui Wang, and Song Han. 2021. Pointacc: Efficient point cloud accelerator. In *MICRO-54: 54th Annual IEEE/ACM* International Symposium on Microarchitecture. 449–461.

# **Opt.1: Serialization-Based Parallel Neighbor Search**

### > Z-curve transform unstructured point clouds into ordered sequences





#### > Z-curve enables parallel search for voxels

- Organize voxels into 8 sets based on the last 3 bits  $\phi_0$  of their z-curve indices
- Each set performs merging and comparison in parallel



# Chlg.2: Useless Attention Computation

### > Attention in local windows

- Each voxel engages in attention calculations with every other voxel
- Unimportant voxels have minimal impact on the results

### > The increased window size

 A fixed point cloud size N and a window size w, the overall complexity is O(Nw)

ı 80

78

72

70

Fest mIoU (%)

• Grows rapidly as *w* increases



### Attention suffers from significant useless computation

# Chlg.2: Useless Attention Computation

### > Dynamic pruning method

Static pruning method

Adopt fixed pruning patterns

Significant accuracy loss

- Pre-calculate a low-precision score matrix
- Additional computation & storage overhead



### Both methods cannot suit serialization-based models

# **Opt.2: Position-Aware Attention Pruning**

- > Boundaries are more crucial than interiors
  - Once an object's boundary is identified, it becomes feasible to be segmented from the background
  - Pruning tokens from interior points minimally affects the results
- > Get neighbor counts from neighbor search in SpConv
  - Boundary points: voxels with fewer neighbors
  - Interior points: voxels with more neighbors



boundary points

interior points

# Chlg.3: Data Dependencies in Softmax

#### > Avoid numerical overflow in exponentials

- Subtract  $m_i$  from each element in **S**
- m<sub>i</sub> is the maximum value of each row in S
- Maximum value causes data dependencies
  - m<sub>i</sub> calculated after all k<sub>i</sub> multiplications
  - Exponential calculation waits for m<sub>i</sub>
  - Exponential sum computed before Matmul with V



$$output_{i} = softmax(\mathbf{S}_{i}) \cdot \mathbf{V} = \frac{\sum_{j=1}^{n} exp(s_{ij} - m_{i})\mathbf{v}_{j}}{\sum_{j=1}^{n} exp(s_{ij} - m_{i})}$$

Data DependencyData DependencyData Dependency $S_i = \mathbf{q}_i \cdot K^T$  $m_i = MAX(s_i)$  $e_i = EXP(s_i - m_i)$  $d_i = e_i/SUM(e_i)$  $out_i = d_i \cdot V$ 

### Inhibit efficient pipeline for computation of different kj, vj

#### > Use dynamic maximum value m<sub>i</sub>

- Decouble the computation of each  $\mathbf{k}_i$ ,  $\mathbf{v}_i$  in  $\mathbf{K}$ ,  $\mathbf{V}$
- Use the m<sub>i</sub> of the scores s<sub>ii</sub> already calculated

### > Block-wise parallel attention calculation

- · Partition keys & values into blocks along the sequence dimension
- Parallelize attention computation within different blocks





Algorithm: Fine-Grained Attention Mechanism **INPUT:**  $\mathbf{q}_1, \dots, \mathbf{q}_n, \mathbf{k}_1, \dots, \mathbf{k}_n$  and  $\mathbf{v}_1, \dots, \mathbf{v}_n$  of size  $1 \times d_k$  as rows of **Q**, **K** and **V** of size  $n \times d_k$ **OUTPUT: output<sub>1</sub>, ..., output<sub>n</sub>** of size  $1 \times d_k$  as rows of attention output matrix of size  $n \times d_{k}$ Parallel for i in range(n) do **Parallel for**  $B_{\gamma}$  in  $\{B_1, B_2, ..., B_Y\}$  do Initialize  $\mathbf{o}_{i}^{(B_{y})} = (0)_{1 \times d_{k}}, l_{i}^{(B_{y})} = 0, m'_{i}^{(B_{y})} = -\infty$ for *j* in  $B_{\nu}$  do  $s_{ii} = \mathbf{q}_i \cdot \mathbf{k}_i^T$  $m_i^{(B_y)} = max\left(m_i^{(B_y)}, s_{ij}\right)$  $\widetilde{s}_{ij} = exp\left(s_{ij} - m_i^{(B_y)}\right)$  $\widetilde{m}_i^{(B_y)} = exp\left(m_i^{(B_y)} - m_i^{(B_y)}\right)$  $l_i^{(B_y)} = \widetilde{m}_i^{(B_y)} \widetilde{l}_i^{(B_y)} + \widetilde{s}_{ij}$  $\mathbf{o}_i^{(B_y)} = \widetilde{m}_i^{(B_y)} \cdot \mathbf{o}_i^{(B_y)} + \widetilde{s}_{ij} \cdot \mathbf{v}_j$  $u'_{\cdot}^{(B_{\mathcal{Y}})} = m_{\cdot}^{(B_{\mathcal{Y}})}$  $M_i = MAX^{\iota}(m_i^{(B_1)}, m_i^{(B_2)}, ..., m_i^{(B_Y)})$  $\sum_{y} exp\left(m_{i}^{\left(B_{y}\right)} - M_{i}\right) \cdot \mathbf{o}_{i}^{\left(B_{y}\right)}$  $output_i =$  $\sum_{y} exp\left(m_{i}^{\left(B_{y}\right)} - M_{i}\right) l_{i}^{\left(B_{y}\right)}$ 

#### Dataflow of $PE_{i,y}$

		S <sub>ip</sub>			Vec Mu			
Stage 2		$m_i^{(By)}$	$m_i^{(By)}$	•••	$m_i^{(By)}$	Dyna	amic Max	
								-
> Sta	de 2							
	3							

The stored m<sub>i</sub> is compared with s<sub>ii</sub> to get a new m<sub>i</sub>

Algorithm: Fine-Grained Attention Mechanism **INPUT:**  $\mathbf{q}_1, \dots, \mathbf{q}_n, \mathbf{k}_1, \dots, \mathbf{k}_n$  and  $\mathbf{v}_1, \dots, \mathbf{v}_n$  of size  $1 \times d_k$  as rows of **Q**, **K** and **V** of size  $n \times d_k$ **OUTPUT: output<sub>1</sub>, ..., output<sub>n</sub>** of size  $1 \times d_k$  as rows of attention output matrix of size  $n \times d_{k}$ Parallel for i in range(n) do **Parallel for**  $B_y$  in  $\{B_1, B_2, ..., B_Y\}$  do **Initialize**  $\mathbf{o}_{i}^{(B_{y})} = (0)_{1 \times d_{k}}, \ l_{i}^{(B_{y})} = 0, \ m_{i}^{(B_{y})} = -\infty$ for *j* in  $B_{y}$  do  $s_{ij} = \mathbf{q}_i \cdot \mathbf{k}_j^T$  $m_i^{(B_y)} = max\left(m_i^{(B_y)}, s_{ij}\right)$  $\begin{aligned} \widetilde{s}_{ij} &= exp\left(s_{ij} - m_i^{(B_y)}\right) \\ \widetilde{m}_i^{(B_y)} &= exp\left(m'_i^{(B_y)} - m_i^{(B_y)}\right) \end{aligned}$  $l_i^{(B_y)} = \widetilde{m}_i^{(B_y)} \widetilde{l}_i^{(B_y)} + \widetilde{s}_{ij}$  $\mathbf{o}_{i}^{(B_{y})} = \widetilde{m}_{i}^{(B_{y})} \cdot \mathbf{o}_{i}^{(B_{y})} + \widetilde{s}_{ij} \cdot \mathbf{v}_{j}$  $m_i^{(B_y)} = m_i^{(B_y)}$  $M_i = MAX \left( m_i^{(B_1)}, m_i^{(B_2)}, ..., m_i^{(B_Y)} \right)$  $\sum_{y} exp\left(m_{i}^{(B_{y})} - M_{i}\right) \cdot \mathbf{o}_{i}^{(B_{y})}$  $output_i =$  $\sum_{y} exp\left(m_i^{(B_y)} - M_i\right) l_i^{(B_y)}$ 



• Subtracting the new m<sub>i</sub>

Algorithm: Fine-Grained Attention Mechanism **INPUT:**  $\mathbf{q}_1, \dots, \mathbf{q}_n, \mathbf{k}_1, \dots, \mathbf{k}_n$  and  $\mathbf{v}_1, \dots, \mathbf{v}_n$  of size  $1 \times d_k$  as rows of **Q**, **K** and **V** of size  $n \times d_k$ **OUTPUT: output<sub>1</sub>, ..., output<sub>n</sub>** of size  $1 \times d_k$  as rows of attention output matrix of size  $n \times d_{k}$ **Parallel for** *i* **in** *range*(*n*) **do Parallel for**  $B_{\gamma}$  in  $\{B_1, B_2, ..., B_Y\}$  do Initialize  $\mathbf{o}_{i}^{(B_{y})} = (0)_{1 \times d_{k}}, l_{i}^{(B_{y})} = 0, m'_{i}^{(B_{y})} = -\infty$ for *j* in  $B_{y}$  do  $s_{ij} = \mathbf{q}_i \cdot \mathbf{k}_i^T$  $m_i^{(B_y)} = max\left(m_i^{(B_y)}, s_{ij}\right)$  $\tilde{s}_{ij} = exp\left(s_{ij} - m\right)$  $(-m_i^{(B_y)})$  $\widetilde{m}_{i}^{(B_{y})} = exp\left(m_{i}^{(B_{y})}\right)$  $l_i^{(B_y)} = \widetilde{m}_i^{(B_y)} \widetilde{l}_i^{(B_y)} + \widetilde{s}_{ij}$  $\mathbf{o}_i^{(B_y)} = \widetilde{m}_i^{(B_y)} \cdot \mathbf{o}_i^{(B_y)} + \widetilde{s}_{ij} \cdot \mathbf{v}_j$  $=m^{(B_y)}$  $M_i = MAX(m_i^{(B_1)}, m_i^{(B_2)}, ..., m_i^{(B_Y)})$  $\sum_{y} exp\left(m_{i}^{\left(B_{y}\right)}\right)$  $output_i =$  $\sum_{v} \overline{exp\left(m_{i}^{\left(B_{y}\right)} - M_{i}\right)l_{i}^{\left(B_{y}\right)}}$ 



- Multiplication of  $s_{ij}$  and  $v_j$  is conducted
- The I<sub>i</sub>, o<sub>i</sub> results updated by the difference between the stored and new m<sub>i</sub>
- The stored m<sub>i</sub> is updated by the new one

Algorithm: Fine-Grained Attention Mechanism **INPUT:**  $\mathbf{q}_1, \dots, \mathbf{q}_n, \mathbf{k}_1, \dots, \mathbf{k}_n$  and  $\mathbf{v}_1, \dots, \mathbf{v}_n$  of size  $1 \times d_k$  as rows of **Q**, **K** and **V** of size  $n \times d_k$ **OUTPUT: output<sub>1</sub>, ..., output<sub>n</sub>** of size  $1 \times d_k$  as rows of attention output matrix of size  $n \times d_{k}$ **Parallel for** i **in** range(n) **do Parallel for**  $B_{y}$  in  $\{B_1, B_2, ..., B_Y\}$  do Initialize  $\mathbf{o}_{i}^{(B_{y})} = (0)_{1 \times d_{k}}, l_{i}^{(B_{y})} = 0, m'_{i}^{(B_{y})} = -\infty$ for *j* in  $B_{y}$  do  $s_{ij} = \mathbf{q}_i \cdot \mathbf{k}_i^T$  $m_i^{(B_y)} = max\left(m_i^{(B_y)}, s_{ij}\right)$  $\tilde{s}_{ij} = exp\left(s_{ij} - m_i^{(B_{ij})}\right)$  $\tilde{m}_i^{(B_{ij})} = exp\left(m_i^{(B_{ij})}\right)$  $= \widetilde{m}_{i}^{(B_{y})} l_{i}^{(B_{y})} + \widetilde{s}_{ii}$  $+ \tilde{s}_{ii} \cdot \mathbf{v}_{i}$  $M_i = MAX (m_i^{(B_1)}, m_i^{(B_2)})$  $({}^{2}), ..., m_{i}^{(B_{Y})}$  $\sum_{y} exp\left(m_{i}^{\left(B_{y}\right)}\right)$  $output_i =$  $\sum_{y} exp\left(m_{i}^{(B_{y})} - M_{i}\right) l_{i}^{(B_{y})}$ 



### Sum Stage

- Obtain the global maximum value M<sub>i</sub> of all local m<sub>i</sub>
- Updating I<sub>i</sub>, o<sub>i</sub> from each block with the difference between their m<sub>i</sub> and M<sub>i</sub>
- The overall attention result output<sub>i</sub> for q<sub>i</sub> is achieved

Algorithm: Fine-Grained Attention Mechanism **INPUT:**  $\mathbf{q}_1, ..., \mathbf{q}_n, \mathbf{k}_1, ..., \mathbf{k}_n$  and  $\mathbf{v}_1, ..., \mathbf{v}_n$  of size  $1 \times d_k$  as rows of **Q**, **K** and **V** of size  $n \times d_k$ **OUTPUT: output<sub>1</sub>, ..., output<sub>n</sub>** of size  $1 \times d_k$  as rows of attention output matrix of size  $n \times d_{k}$ **Parallel for** i **in** range(n) **do Parallel for**  $B_{y}$  in  $\{B_1, B_2, ..., B_Y\}$  do Initialize  $\mathbf{o}_{i}^{(B_{y})} = (0)_{1 \times d_{k}}, l_{i}^{(B_{y})} = 0, m'_{i}^{(B_{y})} = -\infty$ for *j* in  $B_{\nu}$  do  $s_{ij} = \mathbf{q}_i \cdot \mathbf{k}_j^T$  $m_i^{(B_y)} = max\left(m_i^{(B_y)}, s_{ij}\right)$  $\widetilde{s}_{ij} = exp\left(s_{ij} - m_i^{(B_y)}\right)$  $\widetilde{m}_i^{(B_y)} = exp\left(m_i^{(B_y)} - m_i^{(B_y)}\right)$  $l_i^{(B_y)} = \widetilde{m}_i^{(B_y)} l_i^{(B_y)} + \widetilde{s}_{ij}$  $\mathbf{o}_{i}^{(B_{y})} = \widetilde{m}_{i}^{(B_{y})} \cdot \mathbf{o}_{i}^{(B_{y})} + \widetilde{s}_{ij} \cdot \mathbf{v}_{j}$  $m_i^{(B_y)} = m_i^{(B_y)}$  $M_i = MAX(m_i^{(B_1)}, m_i^{(B_2)}, ..., m_i^{(B_Y)})$  $\sum_{y} exp\left(m_{i}^{\left(B_{y}\right)} - M_{i}\right) \cdot \mathbf{o}_{i}^{\left(B_{y}\right)}$  $output_i =$  $\sum_{y} exp\left(m_{i}^{\left(B_{y}\right)}-M_{i}\right) l_{i}^{\left(B_{y}\right)}$ 

### **Overview of APTO Architecture**



# Serialization & Mapping Core



#### Serialization

- Encodes and sorts points' **z-curve indices**
- ① IE converts coordinates into z-curve indices
- ② N/2 sorter sorts N/2 indices
- ③ NMM combines two sorted arrays into one

- Kernel map creation
  - Supports parallel neighbor searching
  - ① Divides voxels into 8 sets based on  $\phi_0$
  - ② NMM merges the shifted input and output voxels
  - ③ Intersection detector finds voxels in the same position

### Enables parallel processing based on the z-curve indices

# SpConv Core



#### Gather control unit

- ① Appends output coordinates to the ID FIFO
- ② Retrives entries based on the FIFO's front element

#### > Counter

- ① Resets to 0 when starts to load an output's entries
- Counts the total number of entries
- ③ Sets the bitmap element to 1 if exceeds the threshold

### Performs SpConv & uses neighbor counts to guide pruning

### **Attention Core**

> Each unit in the attention core performs 4



- > Parallel processing along the sequence dimension
  - **q**<sub>i</sub> is sent to each row of the array
  - K, V is sent to each column in blocks
- stages in the fine-grained attention dataflow > Support the position-aware pruning strategy
  - Pruned k<sub>i</sub> & v<sub>i</sub> are skipped according to the bitmap

High-throughput pipeline & efficient utilization of on-chip buffers

# **Experimental Setup**

- Baseline
  - Octformer & PTv3
  - Implemented in PyTorch 1.12.1 & CUDA 11.8

### Dataset

- ModelNet40
  - Small-scale, for object classification
- ScanNet & S3DIS
  - Large-scale, for scene semantic segmentation

### Hardware platform

- NVIDIA RTX 4090 GPU
- State-of-the-art point cloud accelerators
  - PointAcc & SpOctA

#### Table 1: Models and datasets.

Application	Dataset	Model	#Points
Classification	ModelNet40	Octformer- MN	1k
Segmentation	ScanNet	Octformer- SN PTv3- SN	100k
	S3DIS	PTv3- S3D	1M

[1] Dongxu Lyu, Zhenyu Lil, Yuzhou Chen, Jinming Zhang, Ningyi Xu, and Guanghui He. 2023. SpOctA: A 3D Sparse Convolution Accelerator with Octree-Encoding Based Map Search and Inherent Sparsity-Aware Processing. In 2023 IEEE/ACM International Conference on Computer Aided Design (ICCAD). IEEE, 1–9.

# Effectiveness of the pruning strategy

> Compared to BigBird, ours achieves **higher accuracy** under the same sparsity



Table 2: Impacts of the position-aware pruning strategy.

Model	Result	Loss	Sparsity	Window Size
Octformer- MN	91.75%	0.43%	32.20%	32
Octformer- SN	73.90%	0.52%	33.52%	32
PTv3- SN	76.63%	0.64%	34.02%	1024
PTv3-S3D	72.12%	0.30%	41.35%	128

### Achieves 35.27% sparsity with less than 1% accuracy loss

[1] Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. 2020. Big bird: Transformers for longer sequences. *Advances in neural information processing systems* 33 (2020), 17283–17297.

### **End-to-End Performance**

- Compared with RTX 4090 GPU, PointAcc & SpOctA
  - Speedups of 10.22x, 3.53x, 2.70x
  - Energy savings of **153.59x**, **8.75x**, **7.25x**



Meets real-time processing requirements

### **Attention Optimization**

- > Attention acceleration as window size increases
  - Window size is 32: 4.24x over PointAcc & 2.60x over SpOctA
  - Window size is 1024: **7.95x** over PointAcc & **7.35x** over SpOctA



Efficiently accelerate attention calculation in large-size windows

# **Ablation Study**

- > Optimizations for serialization-based models
  - PNS (Parallel Neighbor Search): 1.14x
  - PAP (Position-Aware Pruning): **1.52x**
  - FAD (Fine-Grained Attention Dataflow): **2.26x**



Majority of optimization effects stems from PAP and FAD



APTO optimizes SpConv's mapping stage by <u>searching neighbor voxels in</u> <u>parallel</u> based on the lowest 3 bits of their z-curve indices

APTO takes a position-aware attention pruning strategy to reduce redundant attention computations, which prunes tokens generated from voxels

# This is the first work to accelerate serialization-based point cloud transformers

APTO adopts a <u>me-grained attention dataflow</u> with dynamic maximum values to remove data dependencies and splits K, V into blocks for parallel processing

APTO achieves average 10.22×, 3.53× and 2.70× speedups over RTX 4090 GPU, PointAcc, and SpOctA, with 153.59×, 8.57× and 7.25× energy savings

### 30th Asia and South Pacific Design Automation Conference (ASPDAC '25) Tokyo, Japan

Thank you Q & A









