# FractalCloud:

A **Fractal**-Inspired Architecture for Efficient Large-Scale Point **Cloud** Processing

## HPCA 2026

Yuzhe Fu, Changchun Zhou, Hancheng Ye, Bowen Duan, Qiyu Huang, Chiyue Wei, Cong Guo, Hai Li, and Yiran Chen

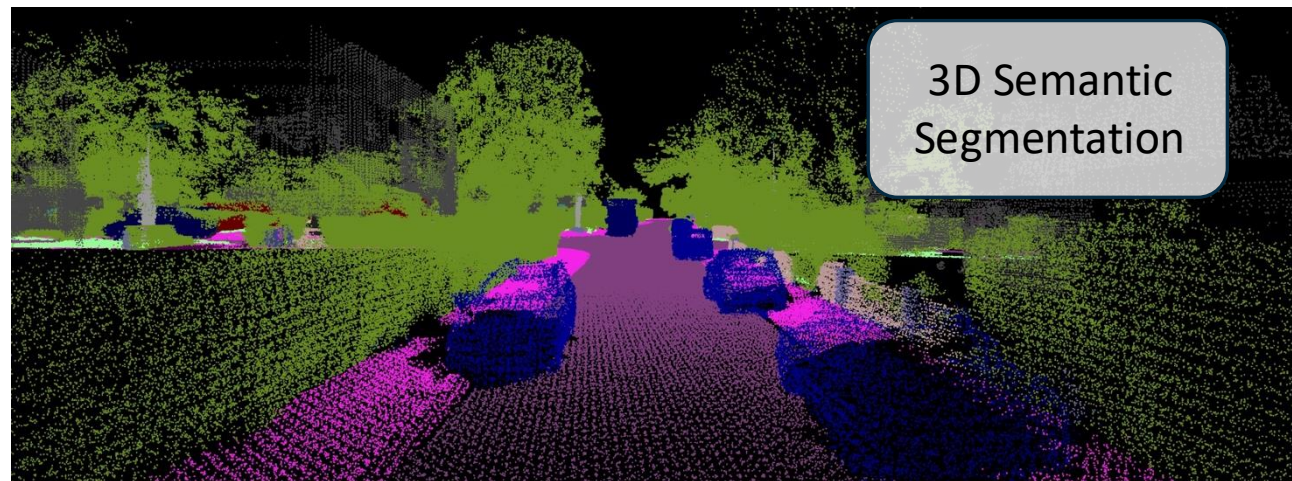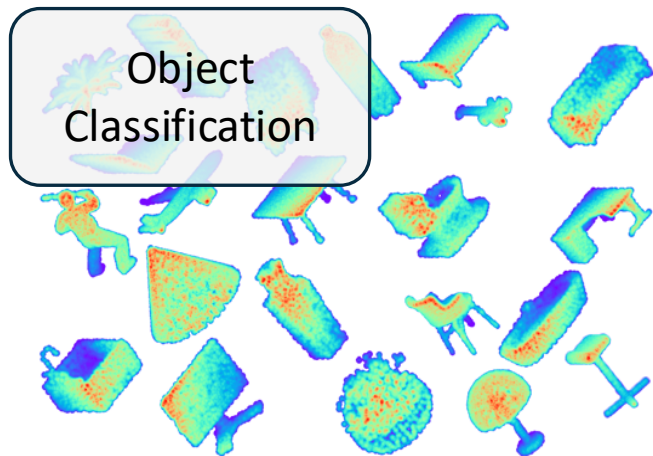*Department of Computer Science, Duke University*
*Department of Electrical and Computer Engineering, Duke University*
*Duke University Center for Computational Evolutionary Intelligence (CEI)*

https://www.geospatialworld.net/news/lizardtech-to-us-patent-lidar-point-cloud-compression-2/

# Point cloud in deep learning: PNN



Object Classification

3D Semantic Segmentation

3D Object Detection

□ : Car    □ : Truck    □ : Pedestrian    □ : Barrier    ■ : Drivable Area    ■ : Lane Divider    ■ : Walkway    ■ : Crosswalk

Qi et al., PointNet: Deep Learning on Point Sets for 3D Classification and Segmentation, CVPR 2017
Vallet et al., TerraMobilita/IQmulus urban point cloud analysis benchmark, CG 2015.
Tang et al., Searching Efficient 3D Architectures with Sparse Point-Voxel Convolution, ECCV 2020.
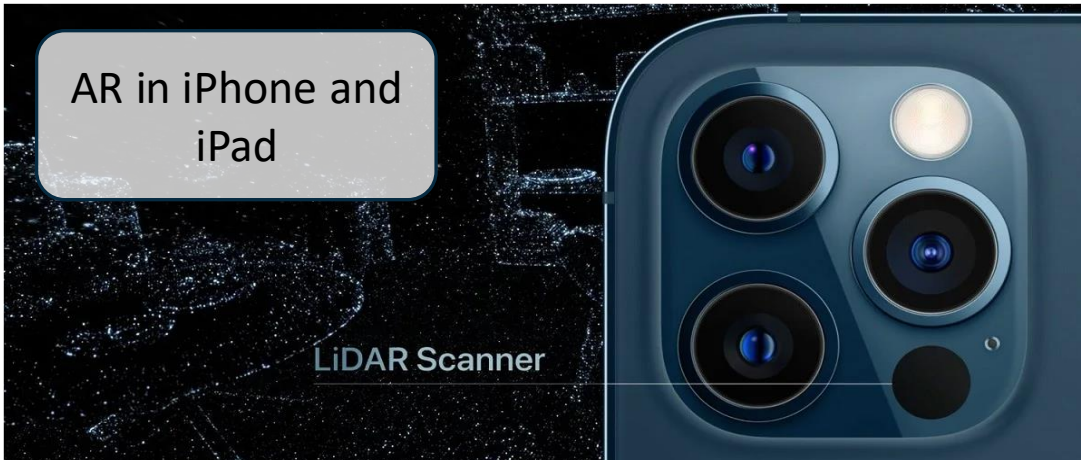
Duke

# Point cloud in daily life

VR Glasses

Autonomous Driving

AR in iPhone and iPad

LiDAR Scanner

Automatic Drones

https://www.softwareone.com/en/insights
iPad - Apple

5 ways LiDAR is transforming the world before our eyes
Tourists scaling the Great Wall of China can now get takeout delivered by drone | CNN Business

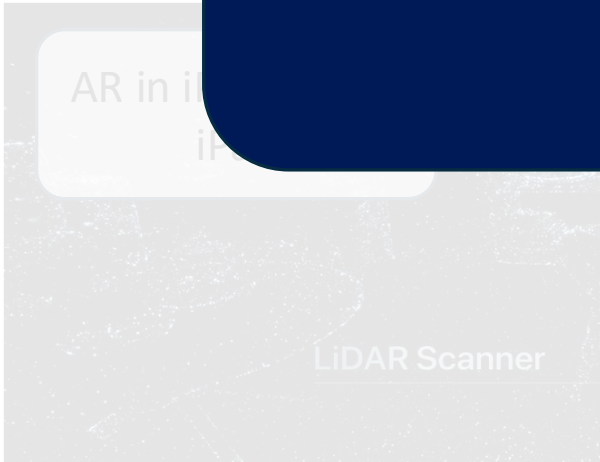# Point cloud in daily life

VR Glasses

Autonomous Driving

AR in iPad

LiDAR Scanner

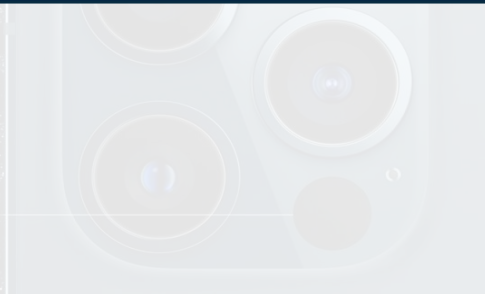**Efficiency is important**

Low latency
Low energy consumption
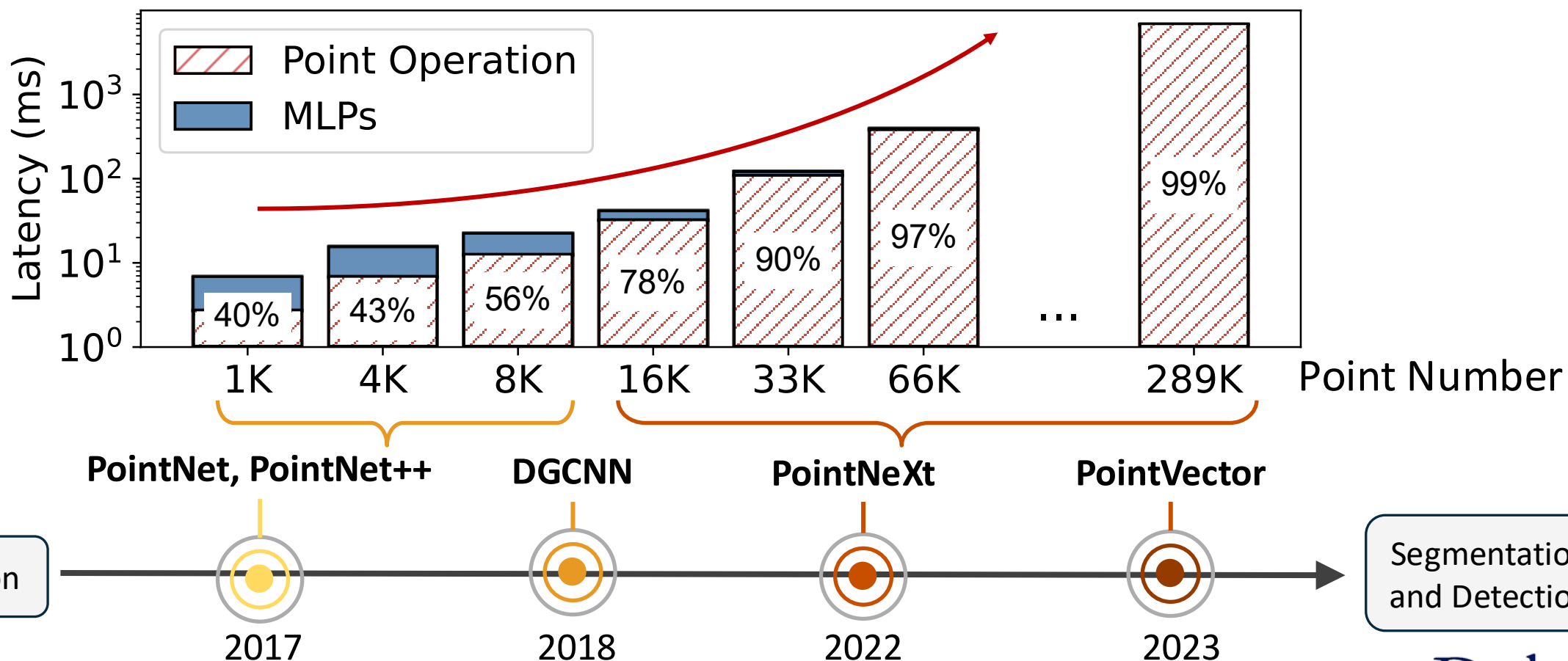
https://www.softwareone.com/en/insights
iPad - Apple

5 ways LiDAR is transforming the world before our eyes
Tourists scaling the Great Wall of China can now get takeout delivered by drone | CNN Business

Duke

# Poor scaling in PNNs

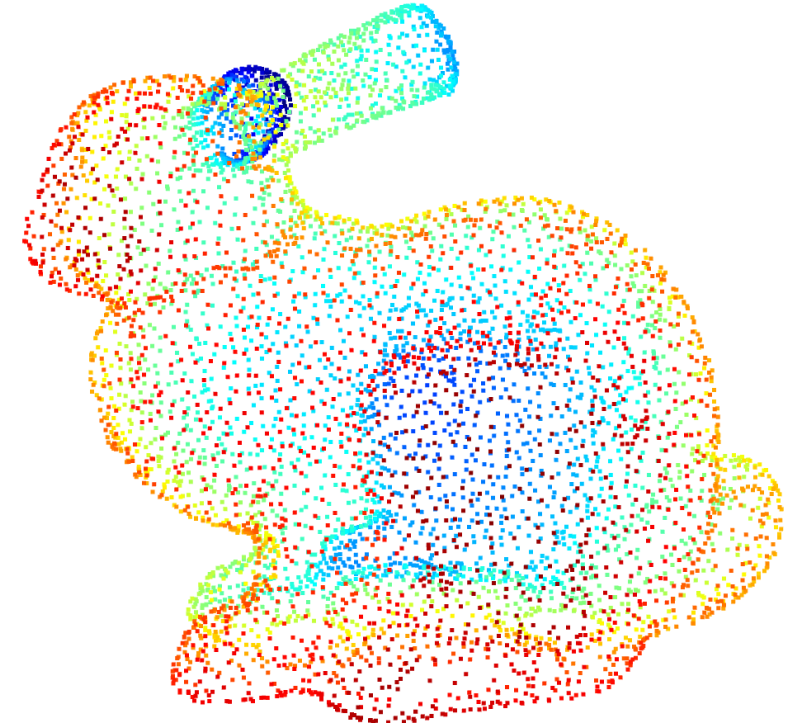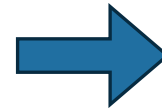## Bottleneck shift: from MLPs to Point Operations

# Point cloud Data



**2D Image**
Pixels: RGB values
Structured in memory

**3D Point Cloud**
Points: (x, y, z), Feature, …
**Unordered** in Memory

Zhou et al., Open3D: A modern library for 3D data processing. arXiv, 2018.
https://pixabay.com/zh/photos/rabbit-nature-wildlife-animal-5469252/

# Point operations in PNNs

## Bottleneck shift: from MLPs to Point Operations

- **Sample**: Centric points

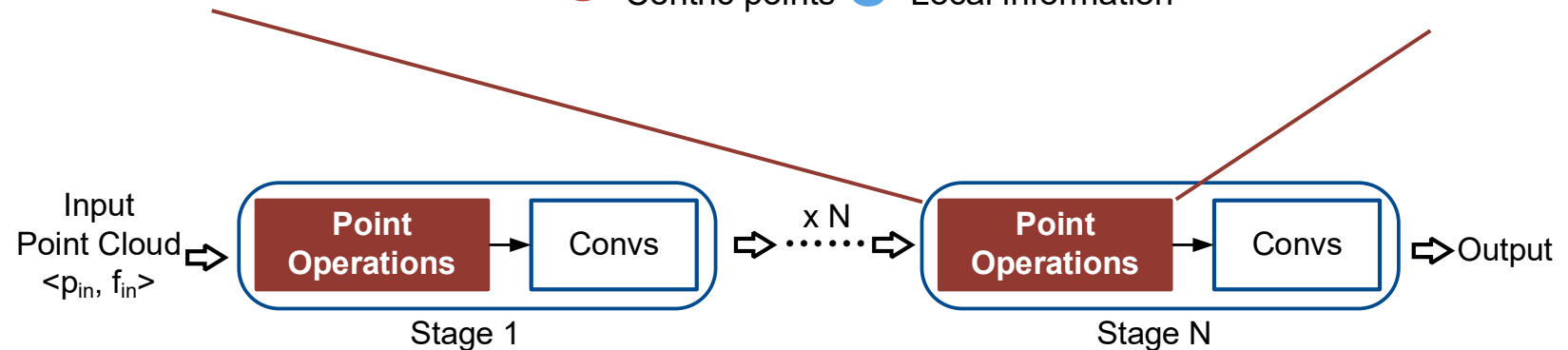- **Neighbor Search**: local information

- **Gather**: Map data from spatial

  domain to feature domain

- **All-to-All Computing**

- **Global Memory Scan**

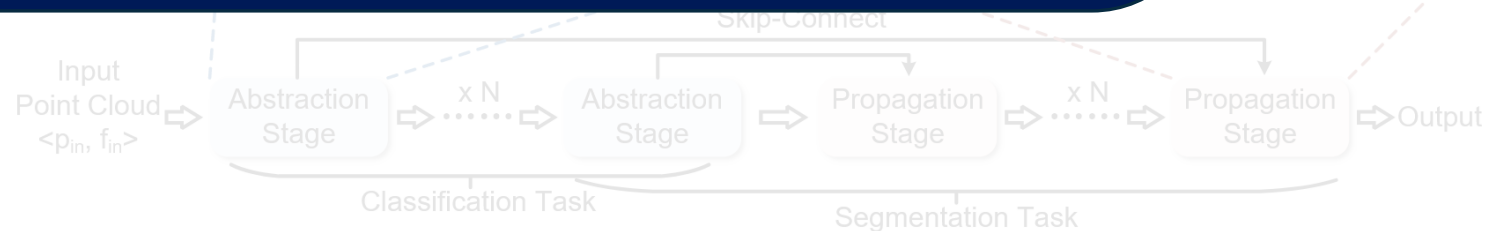- **Iterative Computing**



The backbone of PNNs

# Point operations in PNNs

**Bottleneck shift: from MLPs to Point Operations**

- Samp
- Neig
- Gath

dom

- Irregu
- **Iterative Computing**
- **All-to-All Computing**

**Bottleneck shifts to point operations. Partition can help.**

Input
Point Cloud
$<p_{in}, f_{in}>$

Abstraction
Stage

x N

Abstraction
Stage

Propagation
Stage

x N

Propagation
Stage

Output

Skip-Connect

Classification Task

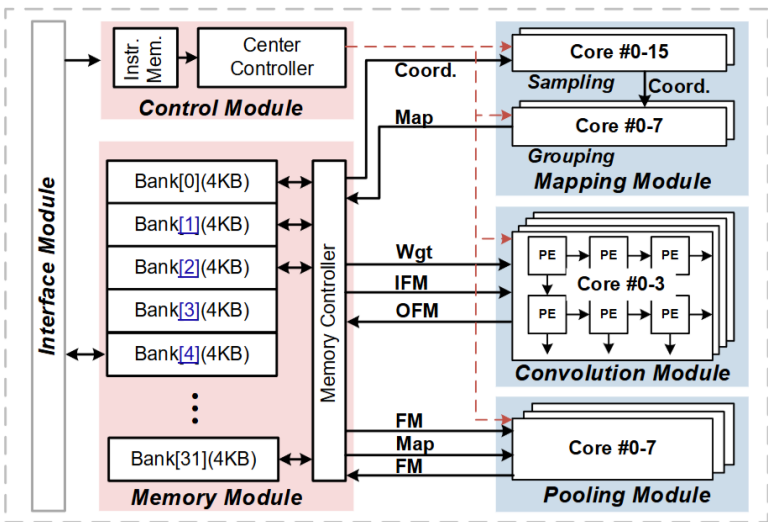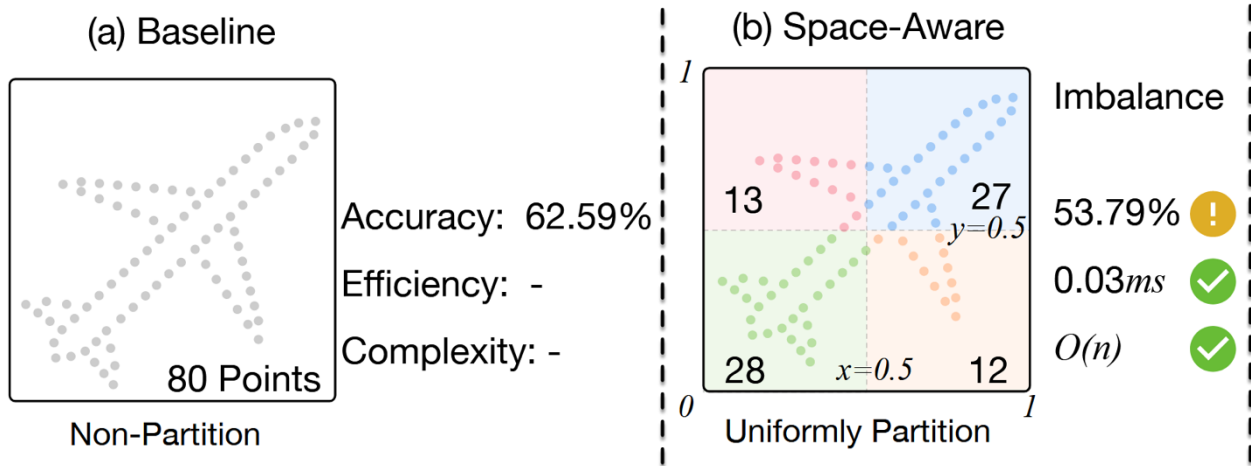Segmentation Task

The backbone of PNNs

Duke

# Current Hardware Architecture

- ## Space-Aware Partition [VLSI'21, ICCAD'23]

  - Example: Uniformly Partition

  - Hardware friendly

  - Streamed memory access

**Imbalanced point distribution**

**Fail to guarantee accuracy**



Kim et al., Pnnpu: A 11.9 tops/w highspeed 3d point cloud-based neural network processor with block-based point processing for regular dram access. VLSI, 2021.
Zhou et al., An Energy-Efficient 3D Point Cloud Neural Network Accelerator with Efficient Filter Pruning, MLP Fusion, and Dual-Stream Sampling. ICCAD, 2023.
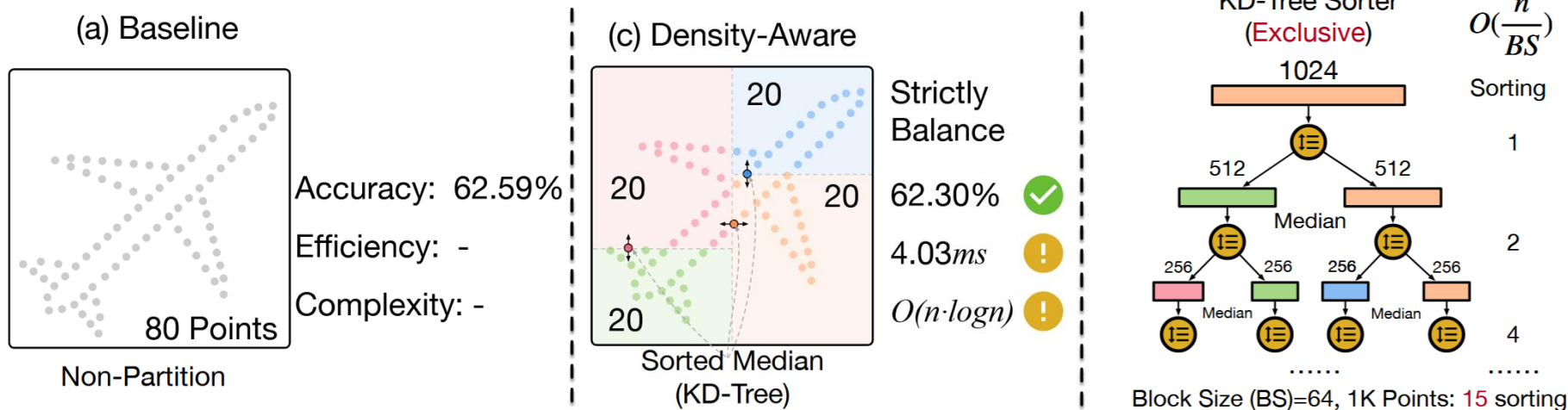
# Current Hardware Architecture

- **Density-Aware Partition [ISCA'22, ASPLOS'25]**

  ◦ Example: KD-Tree

  ◦ Guaranteed accuracy

  ◦ Streamed and balanced memory access

**Exclusive hardware**

**Acceptable when small-scale process**

**New bottleneck for large-scale process**



(a) Baseline

Accuracy: 62.59%
Efficiency: -
Complexity: -

80 Points

Non-Partition

(c) Density-Aware

20      20
20      20

Strictly Balance

62.30% ✓
4.03$ms$ ⚠
$O(n \cdot \log n)$ ⚠

Sorted Median
(KD-Tree)

KD-Tree Sorter
(Exclusive)     $O(\frac{n}{BS})$

1024                    Sorting

512        512          1

Median

256  256  256  256      2

Median    Median

4

Block Size (BS)=64, 1K Points: 15 sorting
BS=256, 289K Points: 2047 sorting

Feng et al., Crescent: taming memory irregularities for accelerating deep point cloud analytics. ISCA, 2022.
Feng et al., StreamGrid: Streaming Point Cloud Analytics via Compulsory Splitting and Deterministic Termination. ASPLOS, 2025.
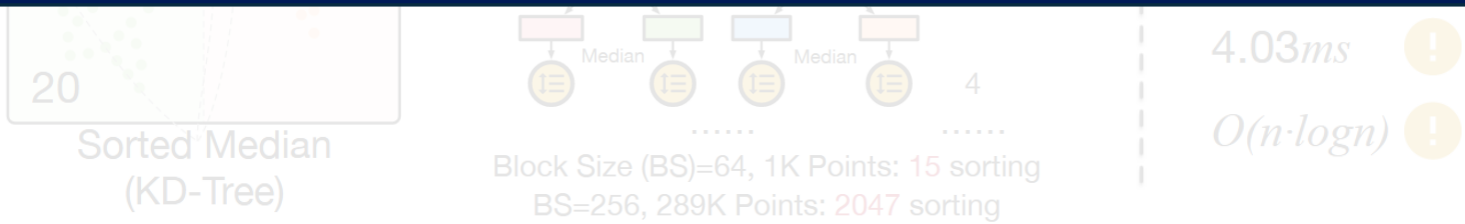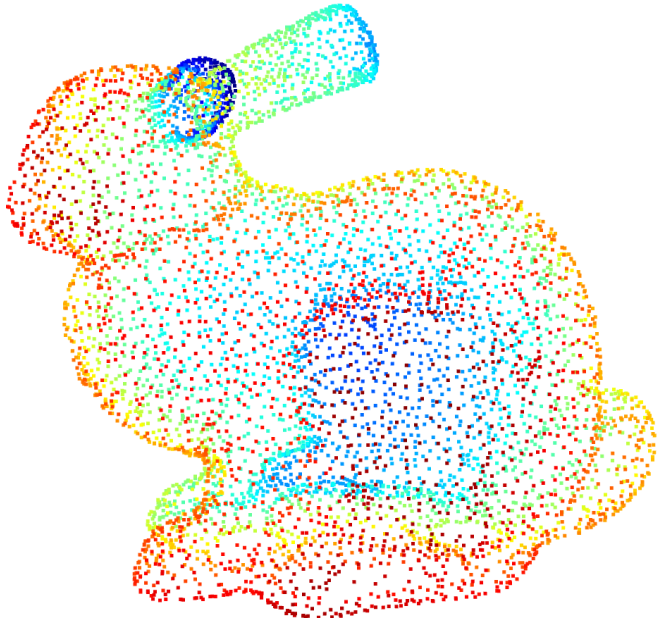
- **Densi**

  ○ Exa

  ○ Str

  ○ Gu

**We hope partition could be**

Accurate **&** Efficient

20

Sorted Median
(KD-Tree)

Median    Median

4

......              ......

Block Size (BS)=64, 1K Points: 15 sorting
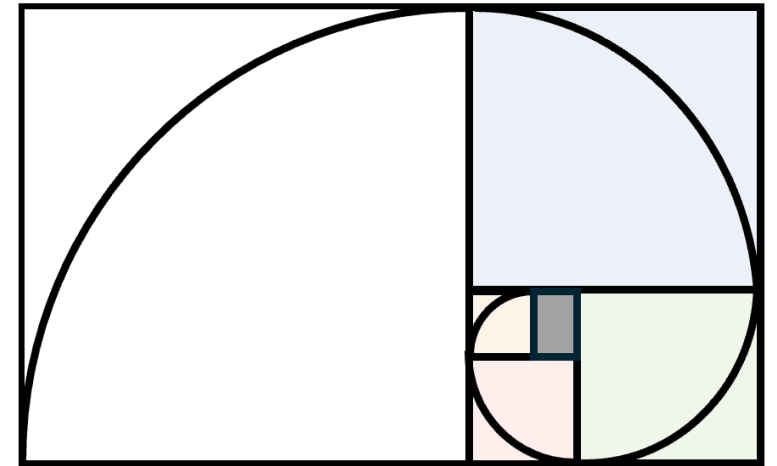BS=256, 289K Points: 2047 sorting

$4.03ms$

$O(n{\cdot}logn)$

Feng et al., Crescent: taming memory irregularities for accelerating deep point cloud analytics. ISCA, 2022.
Feng et al., StreamGrid: Streaming Point Cloud Analytics via Compulsory Splitting and Deterministic Termination. ASPLOS, 2025.

Duke

# **Fractal** Insight



**Real point clouds follows geometry**

**Inspired by fractal geometry**
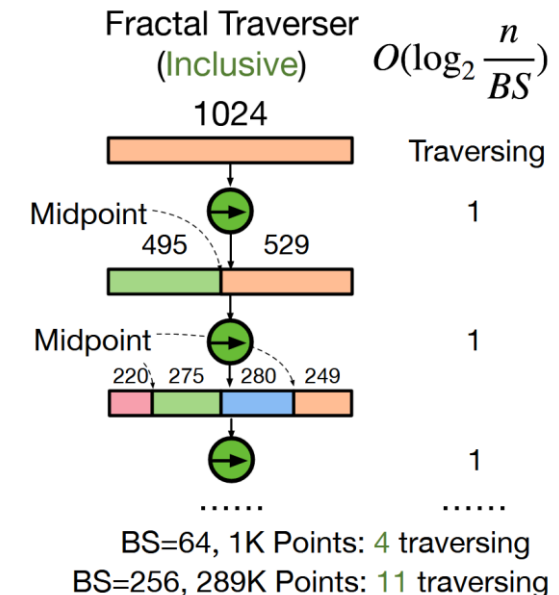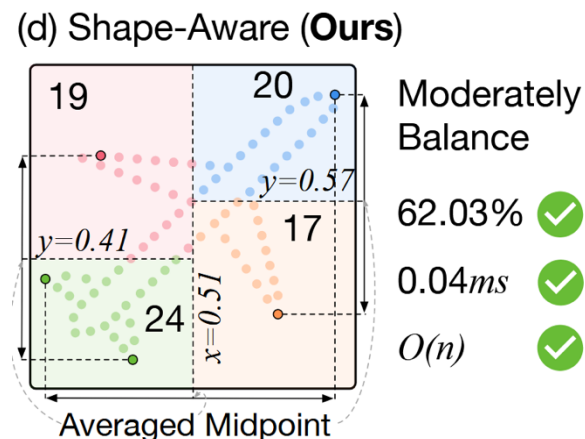
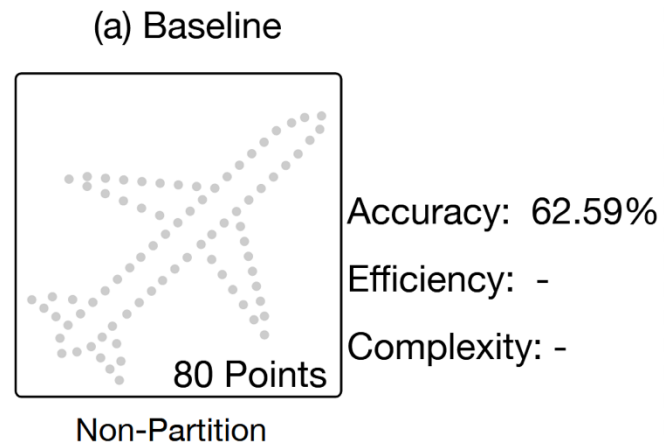◦ Traverse shape, not sort

# Fractal: Accurate and Efficient

- **Shape-Aware Partition**

  ◦ Streamed memory access

  ◦ Guaranteed accuracy

  **Inclusive hardware**

  **Efficient for all-scale process**



(a) Baseline

Accuracy: 62.59%

Efficiency: -

Complexity: -

80 Points

Non-Partition

(d) Shape-Aware (**Ours**)

19   20

$y=0.57$

17

$y=0.41$

24   $x=0.51$

Averaged Midpoint

Moderately Balance

62.03% ✅

$0.04 ms$ ✅

$O(n)$ ✅

Fractal Traverser (Inclusive)   $O(\log_2 \frac{n}{BS})$

1024   Traversing

Midpoint   1

495   529

Midpoint   1

220   275   280   249

1

......   ......

BS=64, 1K Points: 4 traversing

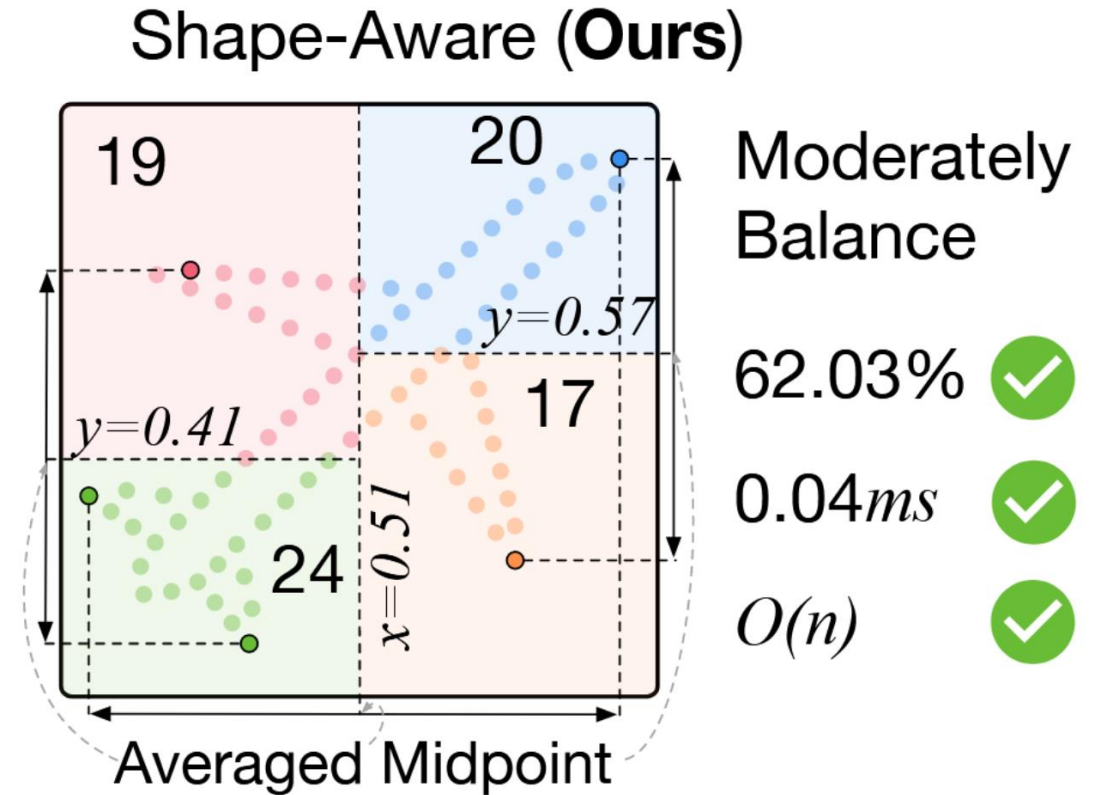BS=256, 289K Points: 11 traversing

# Fractal: Iterative Shape-Aware Partitioning

- **Inputs:**
  - Point cloud
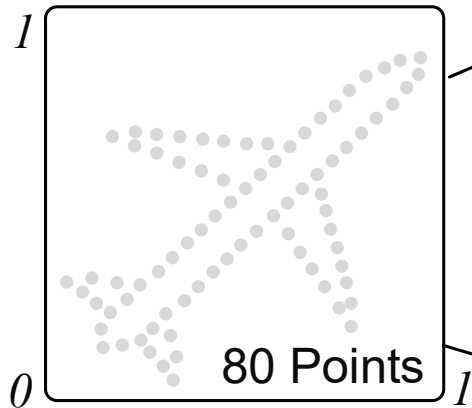  - Threshold (controls block size)

- **Each iteration:**
  - If block size > threshold
    - **Traverse** points along one axis
    - Compute midpoint from min & max
    - Partition
  - Alternate partition axis (x -> y -> z)



Shape-Aware (**Ours**)

19    20

$y=0.57$

$y=0.41$

24    17

$x=0.51$

Averaged Midpoint

Moderately Balance

62.03% ✅

$0.04ms$ ✅

$O(n)$ ✅

# Example for Fractal – 80 Points, threshold 24

**Original Point Cloud**

$1$

80 Points

$0$      $1$

**Start Fractal, with th=24**

**Data Layout in Memory**

| idx | Coordinates |
|-----|-------------|
| 1 | $(x_0, y_0, z_0)$ |
| 2 | $(x_{40}, y_{40}, z_{40})$ |
| 3 | $(x_{63}, y_{63}, z_{63})$ |
| ... | ...... |
| 79 | $(x_8, y_8, z_8)$ |
| 80 | $(x_{56}, y_{56}, z_{56})$ |

**Unordered**

**With Fractal**

19    20

$y=0.57$

$y=0.41$     17

24   $x=0.51$

**After 3 Fractal Iterations**, 4 blocks, all blocks < 24

**Binary Tree Flow**

80

43     37

19   24   17   20

Duke

# Example for Fractal – 80 Points, threshold 24

**Check Fractal**



1

0        80 Points        1

80 > 24, **do Fractal**

**Binary Tree Flow**

**Check**        80

# Example for Fractal – 80 Points, threshold 24



**1ˢᵗ Fractal Iteration**

**Step1: Find min & max** along x-axis

**Binary Tree Flow**

**Process** 80

# Example for Fractal – 80 Points, threshold 24



**1ˢᵗ Fractal Iteration**

**Binary Tree Flow**
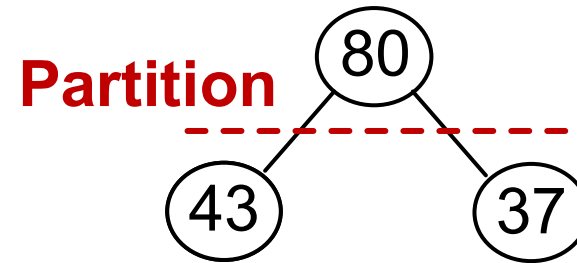
**Process** ⊚80

**Step2: Find midpoint** by (min+max)/2

$x=0.51$

# Example for Fractal – 80 Points, threshold 24

**1st Fractal Iteration**
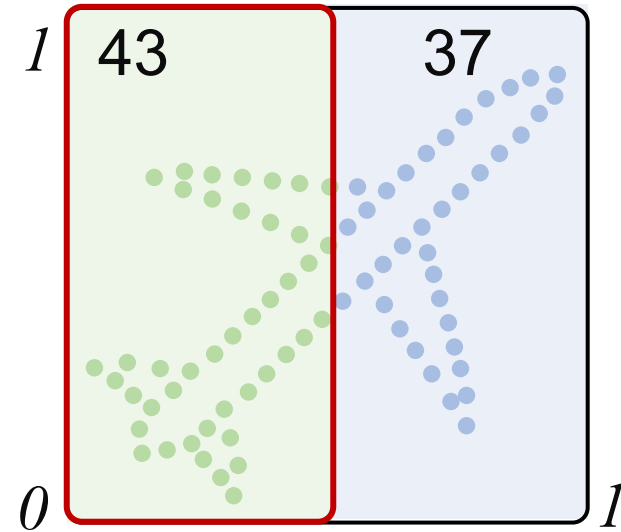


**Binary Tree Flow**

**Step3: Partition 80** into 43- and 37- point blocks

# Example for Fractal – 80 Points, threshold 24

**Check Threshold**



1    43        37

0                1

43 > 24, **do Fractal**

**Binary Tree Flow**

**Check**

80

43        37

Duke

# Example for Fractal – 80 Points, threshold 24
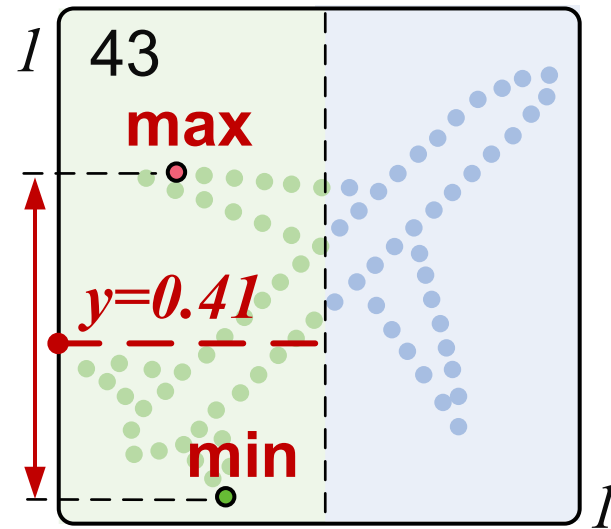


**2ⁿᵈ Fractal Iteration**

**Binary Tree Flow**

**Step1: Find min & max** along y-axis

# Example for Fractal – 80 Points, threshold 24



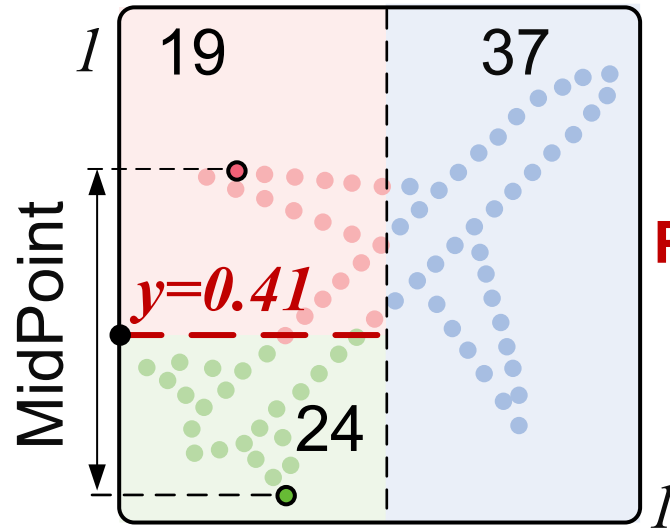**2ⁿᵈ Fractal Iteration**

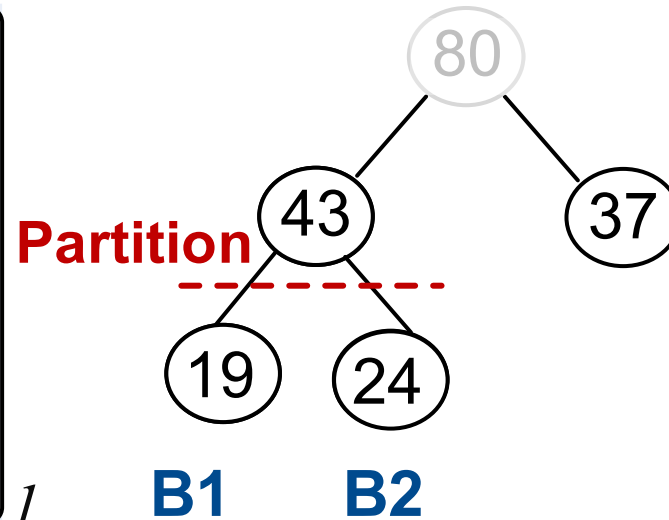**Binary Tree Flow**

**Process**

**Step2: Find midpoint** by (min+max)/2

# Example for Fractal – 80 Points, threshold 24
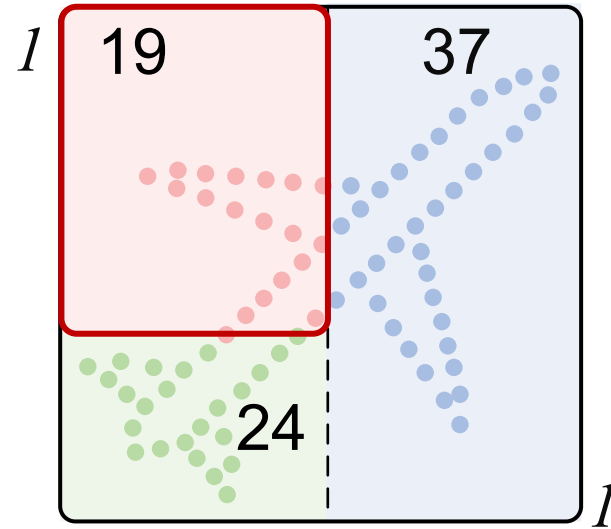


**2$^{nd}$ Fractal Iteration**

**Binary Tree Flow**

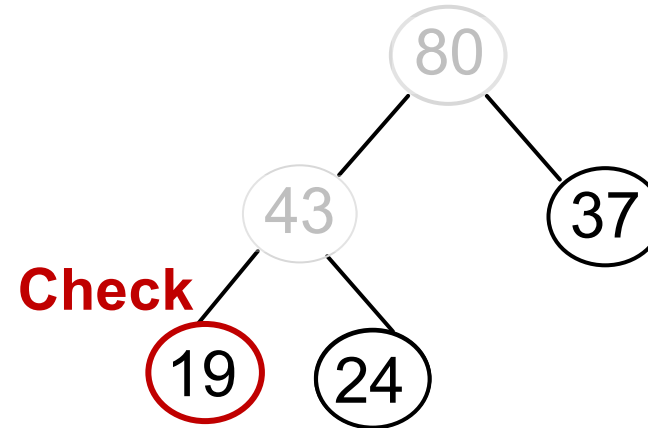**Step3: Partition 43** into 19- and 24- point blocks

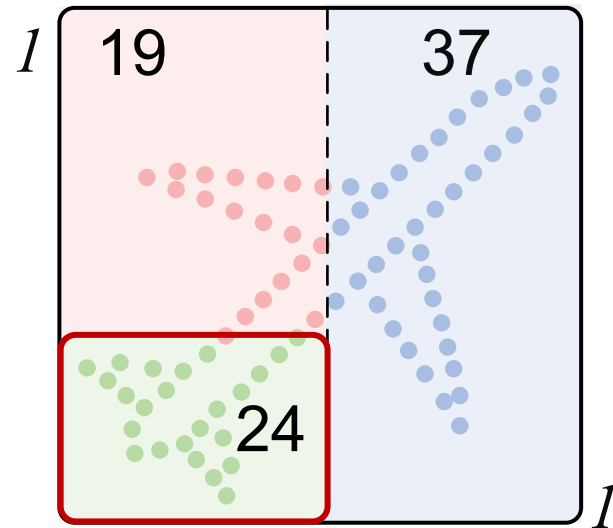# Example for Fractal – 80 Points, threshold 24



**Check Threshold**

1  19   37

24

19 < 24, **no Fractal**

**Binary Tree Flow**

80
43      37
**Check**
19   24

# Example for Fractal – 80 Points, threshold 24

**Check Threshold**



1 | 19     37

24

1

24 == 24, **no Fractal**

**Binary Tree Flow**



80

43     37

19   24   **Check**

Duke

# Example for Fractal – 80 Points, threshold 24

**Check Threshold**

**Binary Tree Flow**

*1* | 19 | 37 | *1*

**Same flow** for all Fractal Iterations

Not sort, **only Linear Traverse**

80

43          37
             **Check**

19   24

# Example for Fractal – 80 Points, threshold 24

**With Fractal**



**Binary Tree Flow**



**After 3 Fractal Iterations**, 4 blocks, all blocks < 24

# Example for Fractal – 80 Points, threshold 24



**With Fractal**

**Binary Tree Flow**

**Depth First Traversal**

**DFT to determine the block order**

# Example for Fractal – 80 Points, threshold 24



**After Fractal**

**Data Layout in Memory**

Original Point Cloud

**Four Point Blocks**

**Spatially Orgnized**

| idx | Coordinates | |
|-----|-------------|---|
| 1 | $(x_0, y_0, z_0)$ | B1 |
| ... | ...... | |
| 20 | $(x_{32}, y_{32}, z_{32})$ | B2 |
| ... | ...... | |
| 44 | $(x_{40}, y_{40}, z_{40})$ | B3 |
| ... | ...... | |
| 80 | $(x_{56}, y_{56}, z_{56})$ | B4 |

# Extend Fractal from Points to Operations

- **Fractal is cheap and scalable.**

- **Blocks are mutually independent.**

**Block-level parallelism**
◦ Local computation and memory access



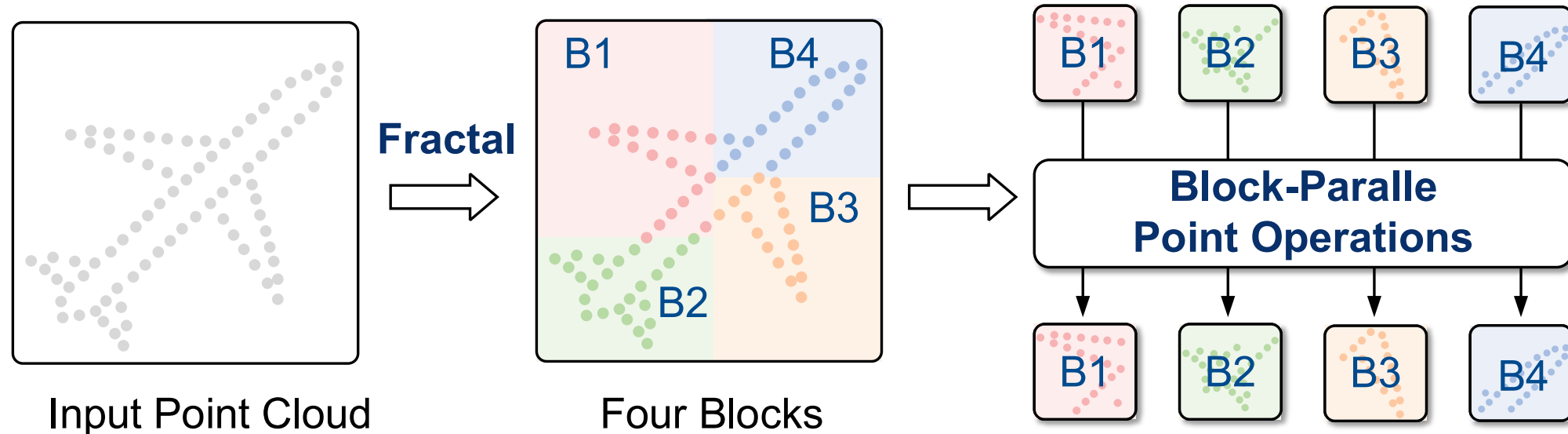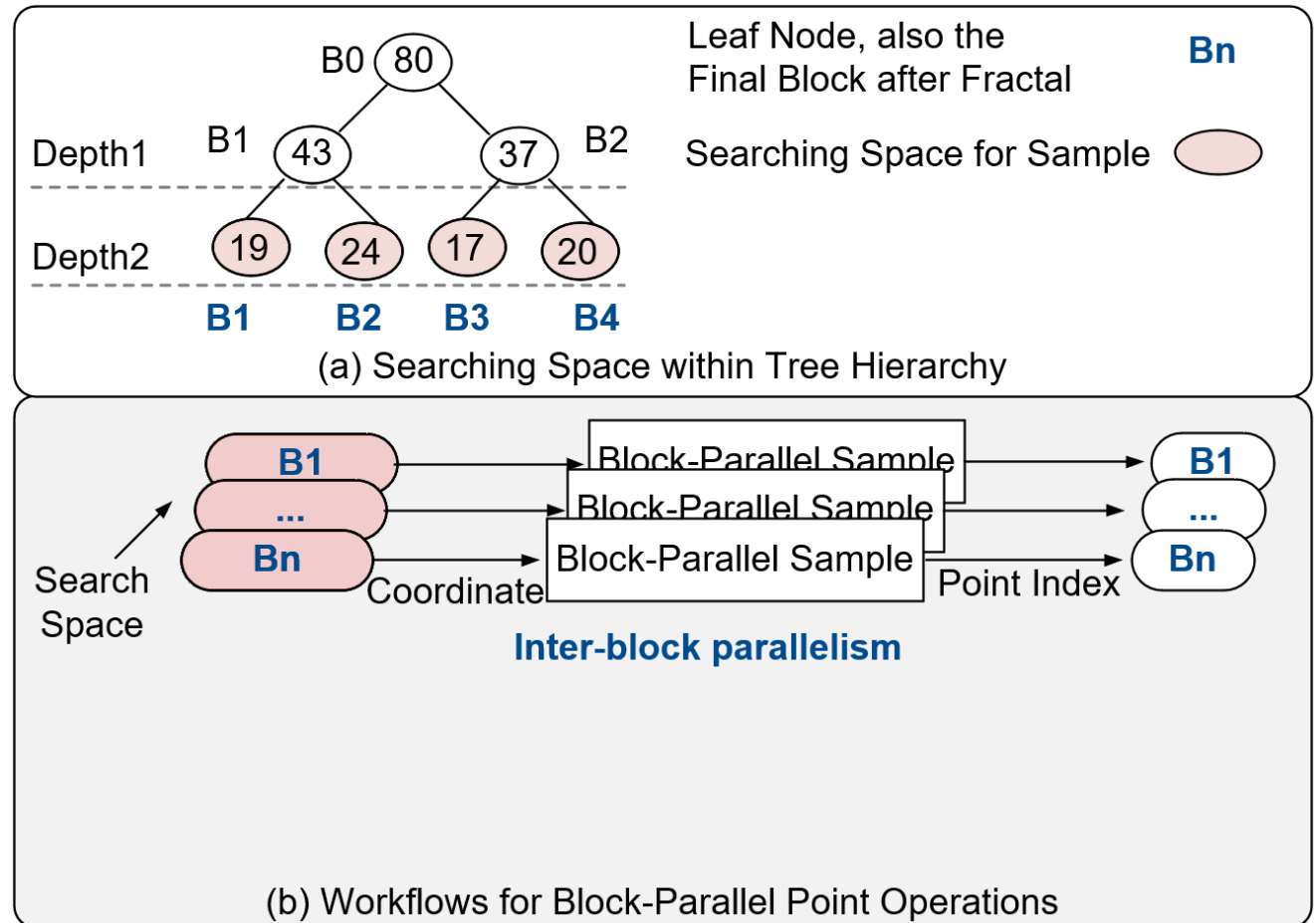Input Point Cloud     Four Blocks

**Block-Paralle
Point Operations**

# Block-Parallel Point Operations

- **Block-wise Sample**

  ◦ Process within current block

  ◦ Inter-block parallelism



Leaf Node, also the Final Block after Fractal — **Bn**

Searching Space for Sample

B0 80

Depth1  B1 43    37  B2

Depth2  19  24  17  20

**B1    B2    B3    B4**

(a) Searching Space within Tree Hierarchy

**B1**
**...**
**Bn**

Search Space

Block-Parallel Sample
Block-Parallel Sample
Block-Parallel Sample

Coordinate                  Point Index

**B1**
**...**
**Bn**

**Inter-block parallelism**

(b) Workflows for Block-Parallel Point Operations

Duke

# Block-Parallel Point Operations

- **Block-wise Neighbor Search**

  ○ Expend searching to parent node

  ○ One parent level is sufficient



(a) Searching Space within Tree Hierarchy

(b) Workflows for Block-Parallel Point Operations

# Block-Parallel Point Operations

- **Block-wise Gather**

  - Same rules as neighbor search



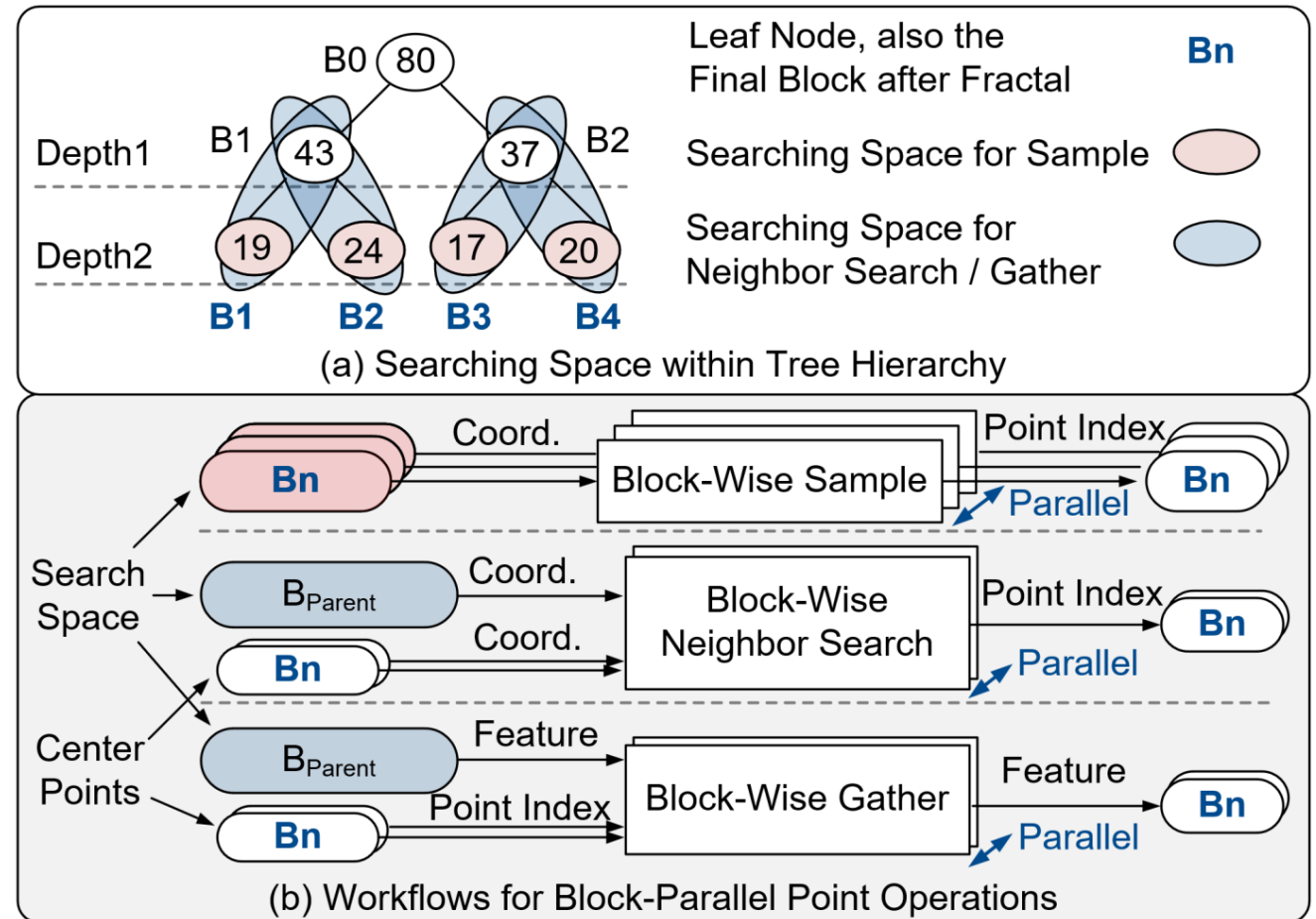(a) Searching Space within Tree Hierarchy

(b) Workflows for Block-Parallel Point Operations
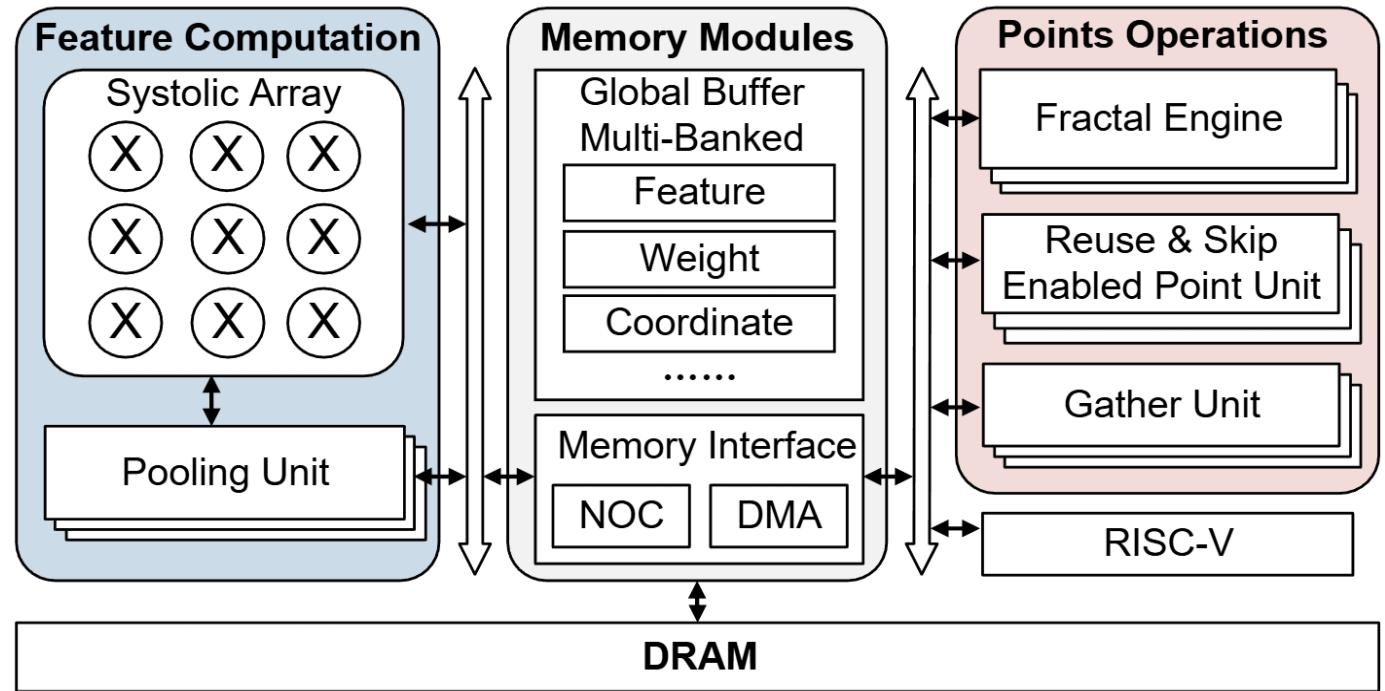
# Block-Parallel Point Operations

**Block-wise Sample**
**Block-wise Neighbor Search**
**Block-wise Gather**

- **Eliminate all-to-all computing**

- **Unlock block-level parallelism**

- **On-chip feasible**



(a) Searching Space within Tree Hierarchy

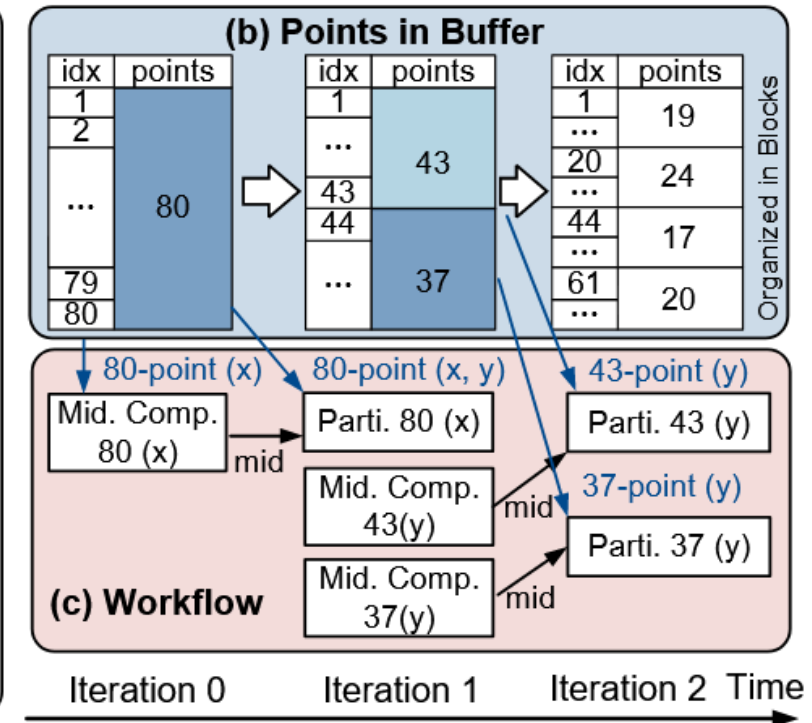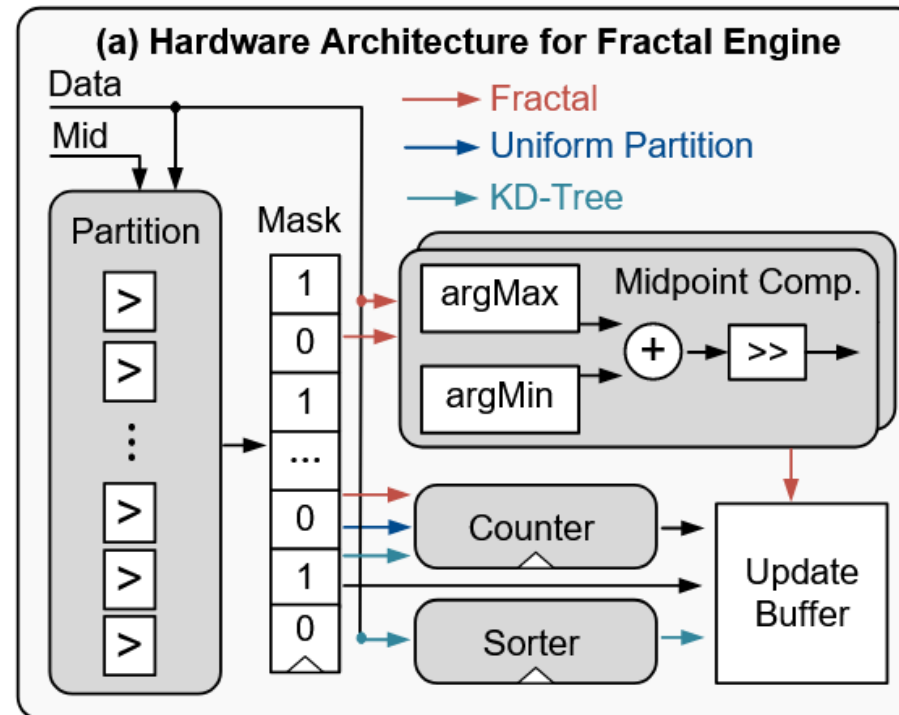(b) Workflows for Block-Parallel Point Operations

# FractalCloud: Point Cloud Accelerator

- Systolic Array

- Network on Chip (NOC)

- Direct Memory Access (DMA)

- RISC-V MCU

- SRAM (274KB)

- **Fractal Engine**

- **Reuse-Skip Enabled Point Unit (RSPU)**

# Fractal Engine

- Reconfigurable structure for multiple partitions:
  - Fractal, KD-Tree, uniform partition.

- Fractal:
  - Simple Hardware
  - Inclusive
  - Fully pipelined

# Reuse Skip Enabled Point Unit (RSPU)

- **Unified module for all point operations**

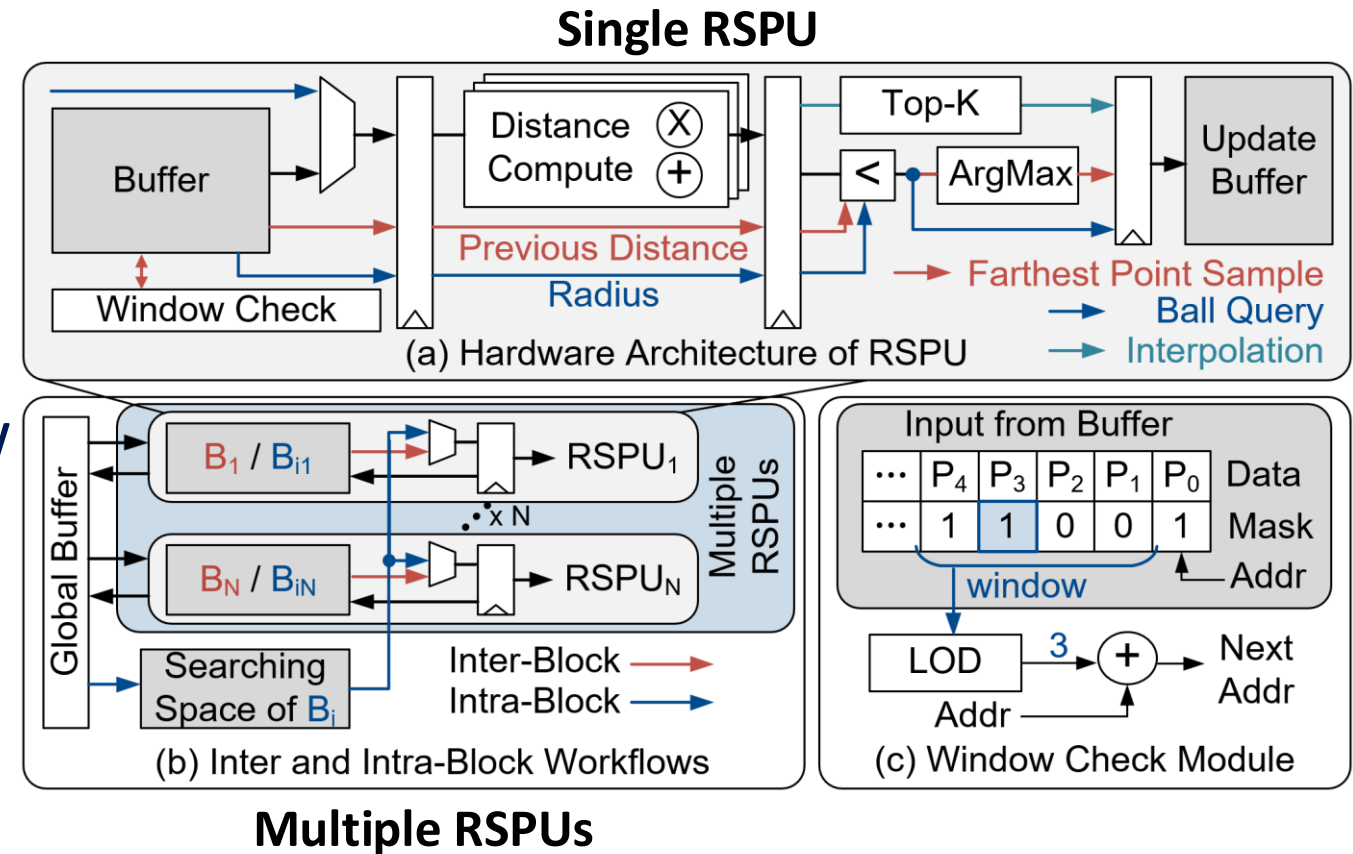  ○ FPS, Ball Query, KNN (Interpolation)

  ○ Blocks run with DFT order

- **Flexible Block-Parallel Workflow**
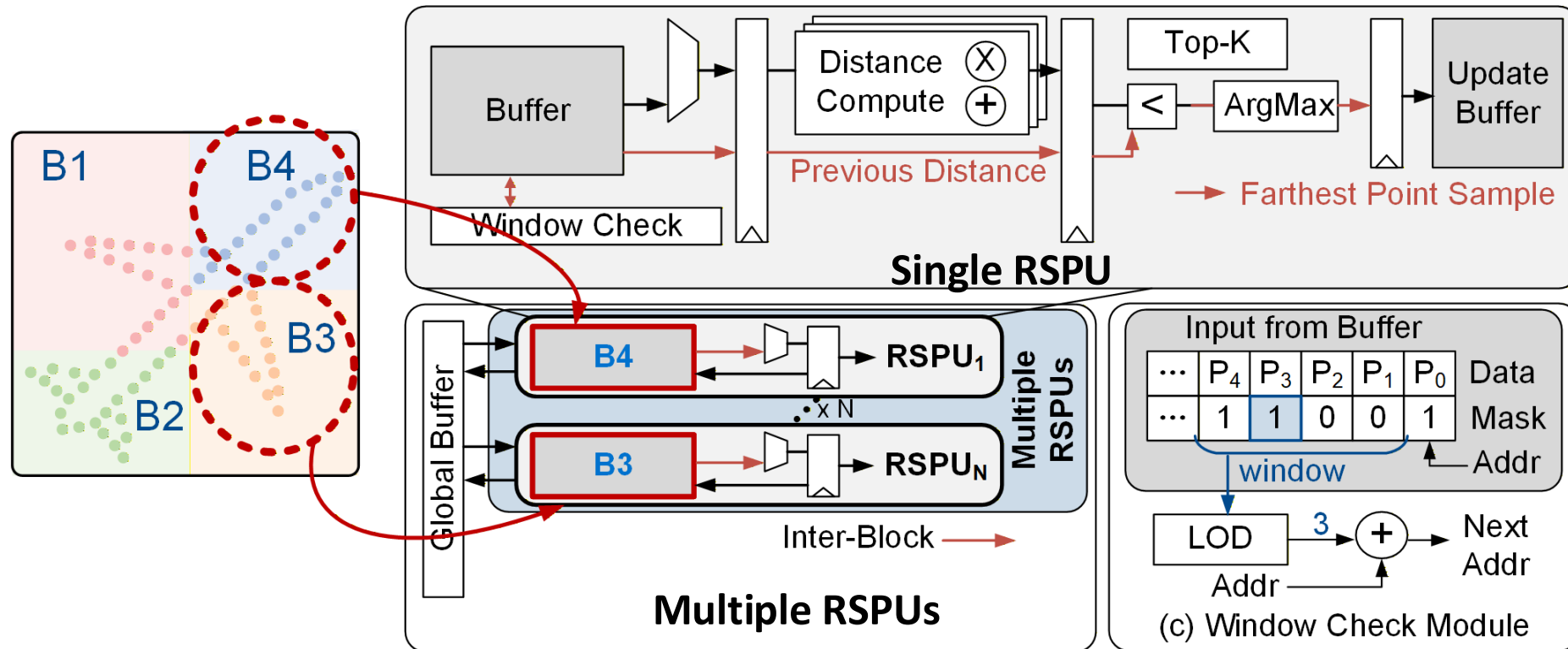
Block-Wise Sample

  ○ **Inter-block parallelism**

Block-Wise Neighbor Search

  ○ **Intra-block parallelism**

**Single RSPU**
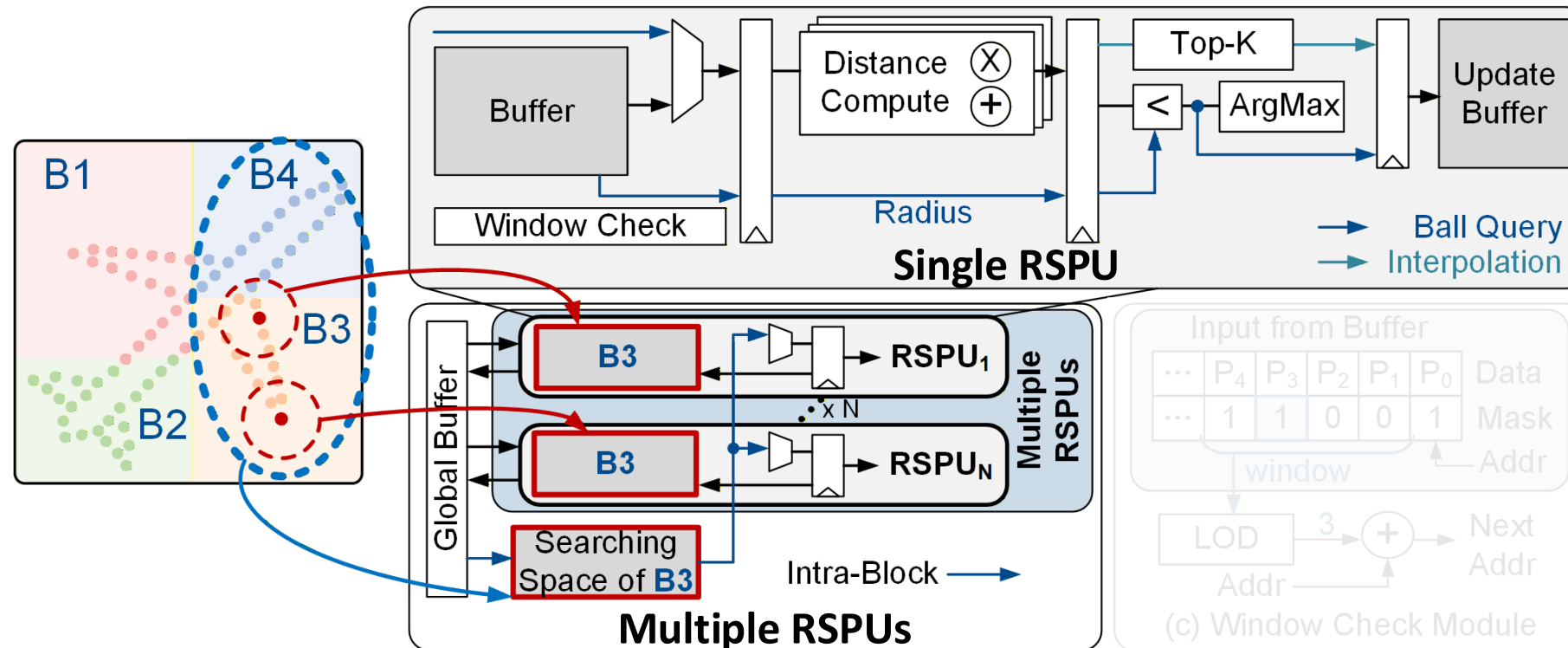


**Multiple RSPUs**

# Flexible Block-Parallel for Multiple RSPUs

- **Block-Wise Sample: inter-block parallelism**
  - Each RSPU handles one FPS within one block

- Window check: **Skip redundant computation**

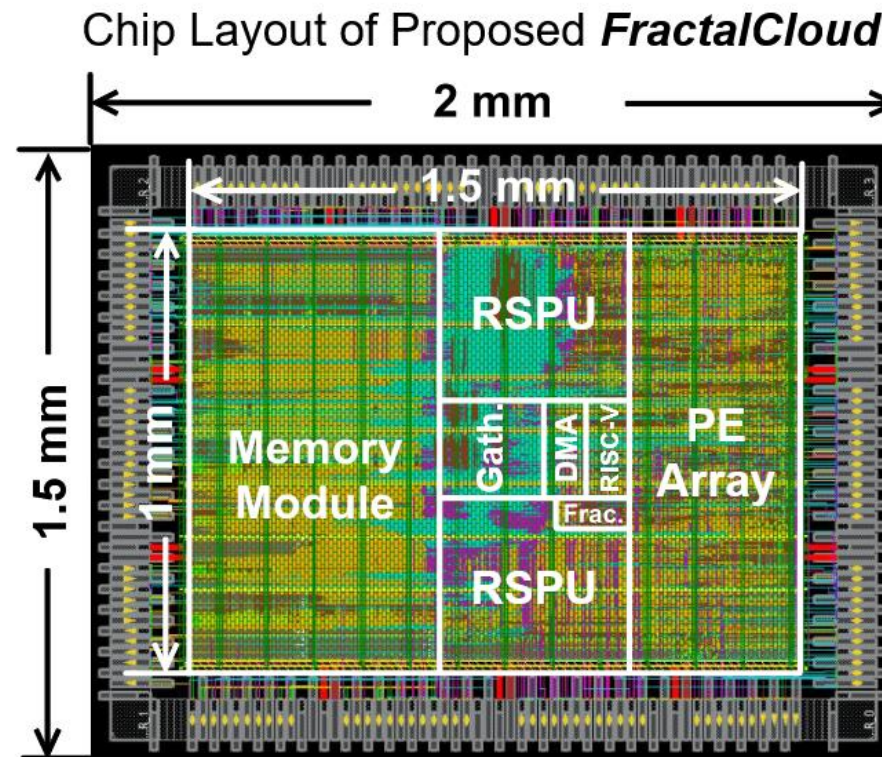# Reuse-and-Skip-enabled Point Unit (RSPU)

- **Block-Wise Neighbor Search: intra-block parallelism**
  - Each RSPU process different centric points in same block
- **Data reusing from parent node**
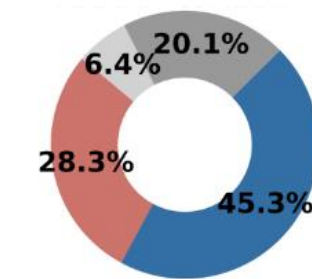
# HW Implementation

- **Small hardware:**
  - TSMC 28nm
  - Core Area: 1.5 mm$^2$
  - Power: 0.58 W
  - Frequency: 1 GHz

Chip Layout of Proposed *FractalCloud*



Detailed Specifications

| Technology | 28nm |
| --- | --- |
| Die Area | 3 mm$^2$ |
| Core Area | 1.5 mm$^2$ |
| SRAM Size | 274 KB |
| Frequency | 1 GHz |
| Ave. Power | 0.58 W |

**Area Breakdown**

6.4%  20.1%
28.3%
45.3%

**Energy Breakdown**

1.2%
9.7%  27.4%
61.6%

PE Array   Memory   RSPU   Other

# Evaluation

- **Network Benchmarks**
  - Inputs scale from 1K to 289K
  - Three PNNs
  - Three Tasks
  - Three Datasets

- **Hardware Architectures**
  - Same PE cores
  - Fixed Frequency
  - Equal DRAM Bandwidth
  - ......

| Model | Notation | Task | Dataset | Scene |
|---|---|---|---|---|
| PointNet++ | PN++ (c) | Classification | ModelNet40 | Object |
| PointNeXt | PNXt (c) | | | |
| PointNet++ | PN++ (ps) | Part Segmentation | ShapeNet | Object |
| PointNeXt | PNXt (ps) | | | |
| PointNet++ | PN++ (s) | Segmentation | S3DIS | Indoor |
| PointNeXt | PNXt (s) | | | |
| PointVector | PVr (s) | | | |

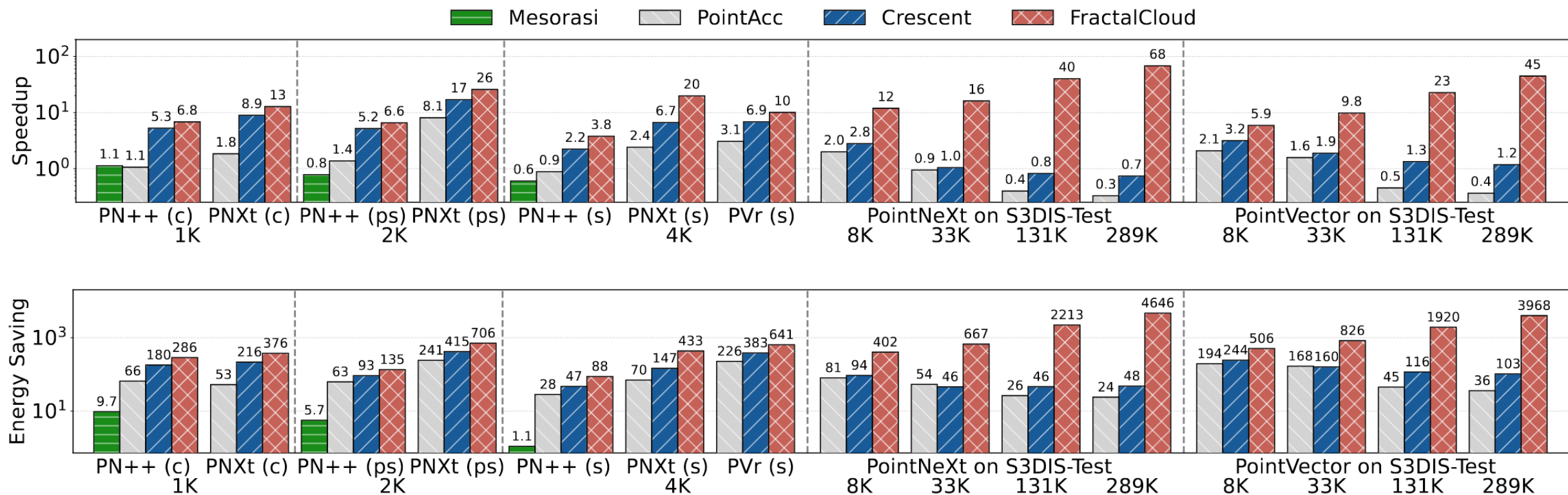| Accelerator | Mesorasi [27] | PointAcc [28] | Crescent [29] | *FractalCloud* |
|---|---|---|---|---|
| Cores | 16x16 | 16x16 | 16x16 | 16x16 |
| SRAM (KB) | 1624 | 274 | 1622.8 | 274 |
| Frequency | 1GHz | 1GHz | 1GHz | 1GHz |
| Area (mm2) | 4.59 | 1.91 | 4.75 | 1.5 |
| DRAM Bandwidth | DDR4-2133 17GB/s | DDR4-2133 17GB/s | DDR4-2133 17GB/s | DDR4-2133 17GB/s |
| Technology | 28nm | 28nm | 28nm | 28nm |
| Peak Performance | 512 GOPS | 512 GOPS | 512 GOPS | 512 GOPS |

# Network Accuracy



**Guaranteed accuracy:**

**Less than 0.7%** accuracy loss for all models
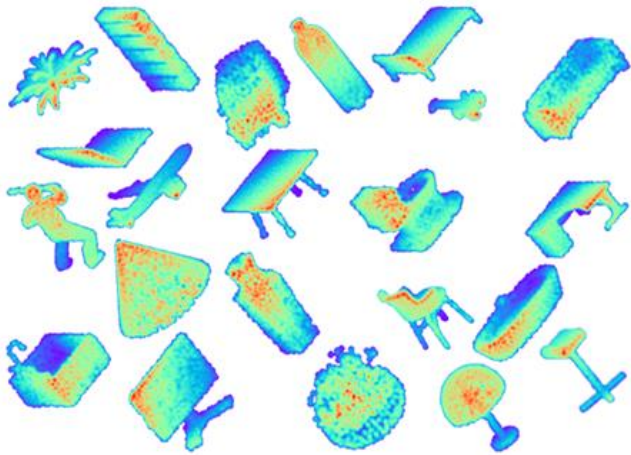
Better performance than SOTA works

# Performance Gain over SOTA accelerators
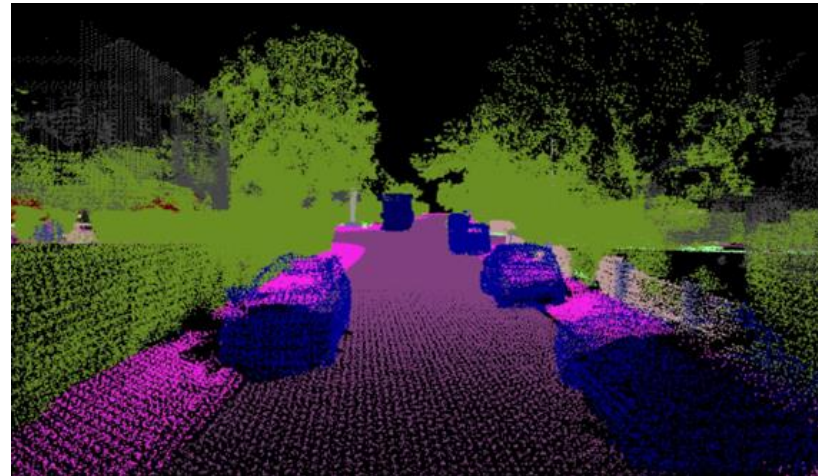


**Huge performance:**

Average 21.7x speedup

Average 27x energy saving
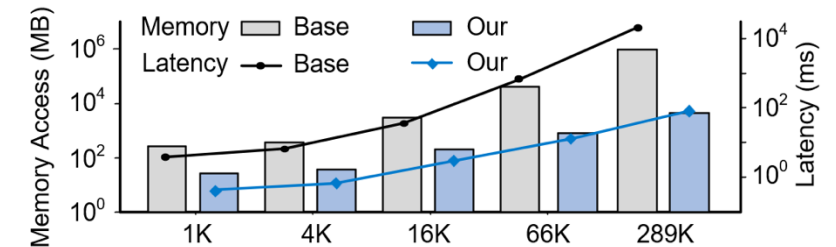
# FractalCloud for Efficient PNN Acceleration

- **Application:** AR/VR, automatic drive, drones, …
- **From small to large input processing**
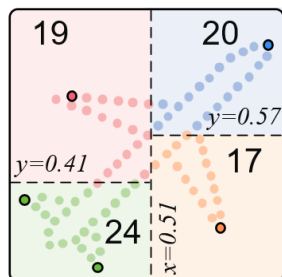


**1K @ 2017 (Simple)**
Object Classification

**300K @ 2024 (Complex)**
Semantic Segmentation

**FractalCloud Optimization**
**21.7x speedup**

**Accuracy and Efficiency**

With Fractal

| idx | Coordinates | |
|-----|-------------|---|
| 1 | $(x_0, y_0, z_0)$ | B1 |
| ... | ...... | |
| 20 | $(x_{32}, y_{32}, z_{32})$ | B2 |
| ... | ...... | |
| 44 | $(x_{40}, y_{40}, z_{40})$ | B3 |
| ... | ...... | |
| 80 | $(x_{56}, y_{56}, z_{56})$ | B4 |

Spatially Orgnized

**Fractal: Shape-Aware Partition**

**Local Computation**

(a) Searching Space within Tree Hierarchy

(b) Workflows for Block-Parallel Point Operations
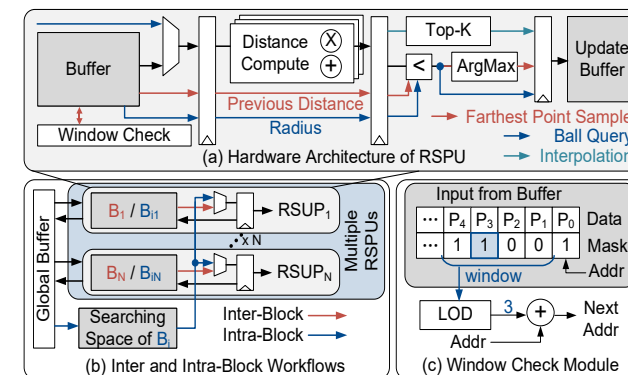
**Block-Parallel Point Operation**

# FractalCloud for Efficient PNN Acceleration

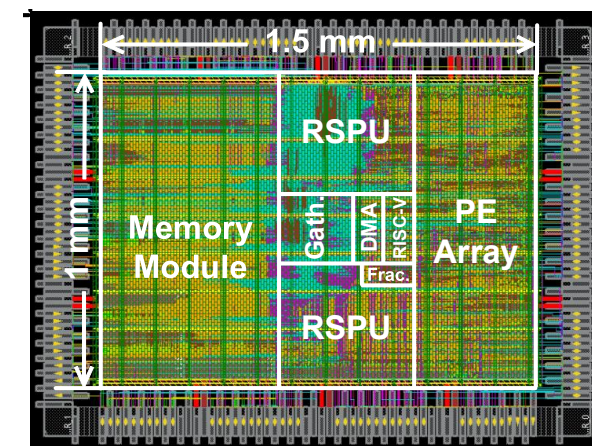Structured Memory, Local Search

Dedicated Architecture, Data Reuse

**21.7× Speedup**
**27× Energy Save**

**Block-Parallel Hardware**

(a) Hardware Architecture of RSPU

(b) Inter and Intra-Block Workflows

(c) Window Check Module

**Reuse-and-Skip-Enabled Point Unit**

**Low latency & low energy cost**

**Area: 1.5mm², Power 0.58W**

# Acknowledgements

Center of Computational Evolutionary Intelligence (CEI)

# FractalCloud
## HPCA 2026



**Authors:**

Yuzhe Fu, Changchun Zhou,
Hancheng Ye, Bowen Duan, Qiyu
Huang, Chiyue Wei, Cong Guo,
Hai Li, and Yiran Chen

# Thanks for Listening.

Codes are open-sourced at

https://github.com/Yuzhe-Fu/FractalCloud



**Duke**

Center of Computational Evolutionary Intelligence (CEI)