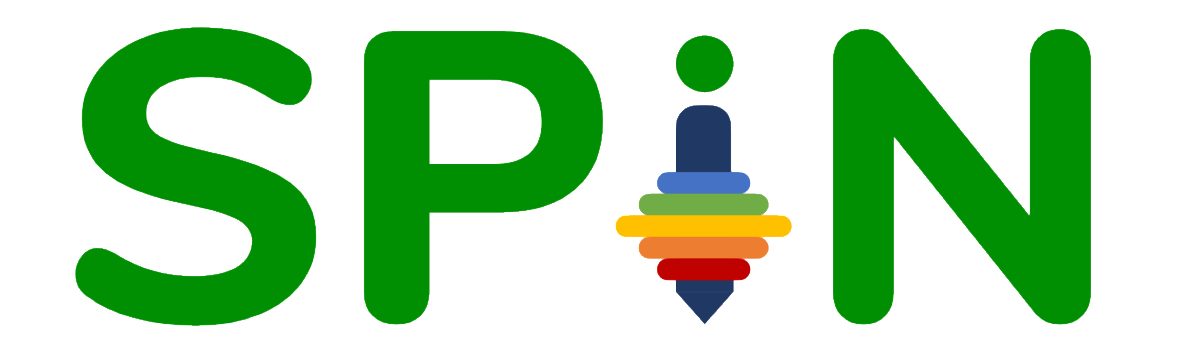


Spira: Exploiting Voxel Data Structural Properties for Efficient Sparse Convolution in Point Cloud Networks



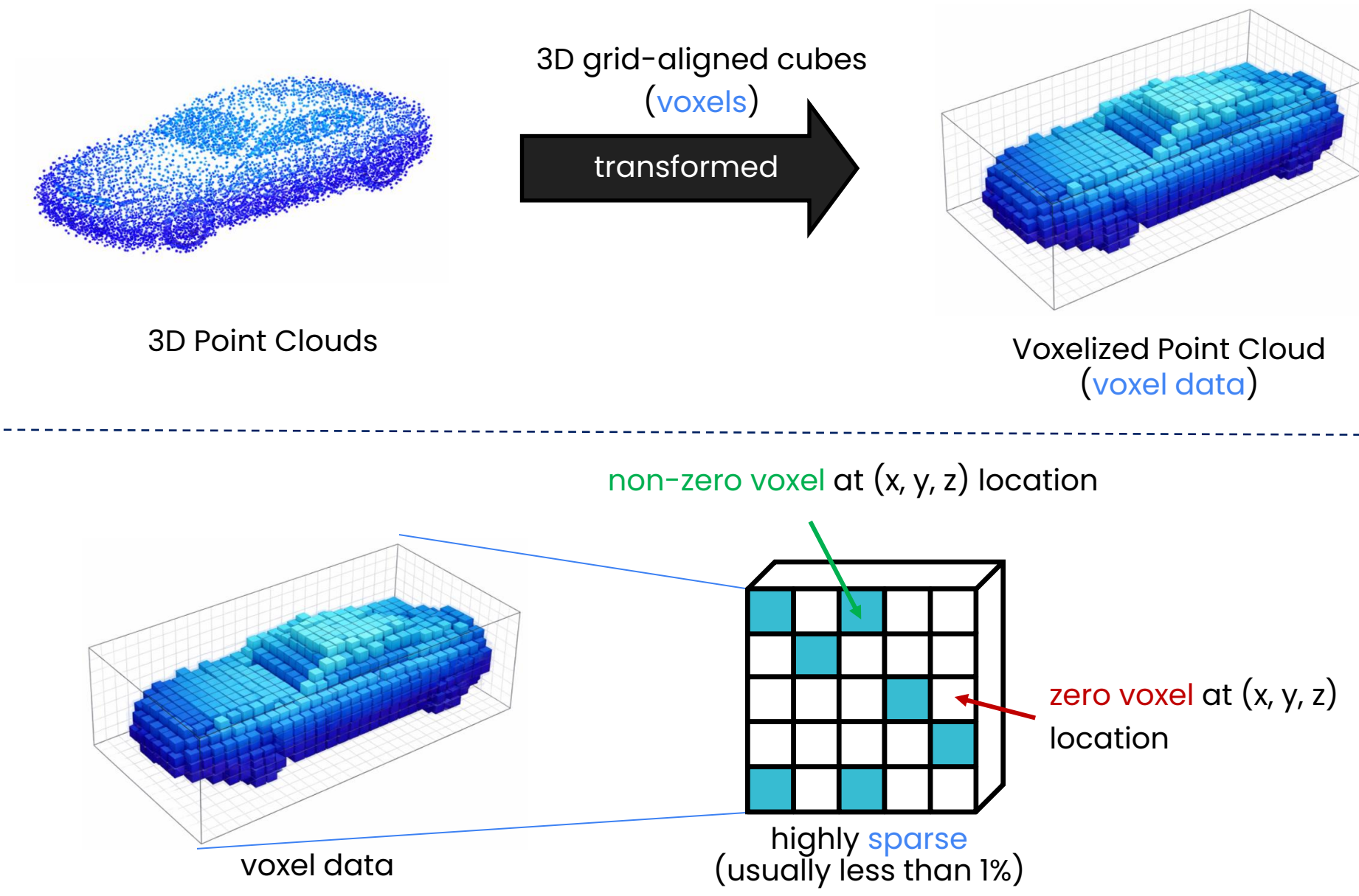
Dionysios Adamopoulos^{1,2} Anastasia Pouloupoulou² Georgios Goumas² Christina Giannoula¹

¹Max Planck Institute for Software Systems ²National Technical University of Athens

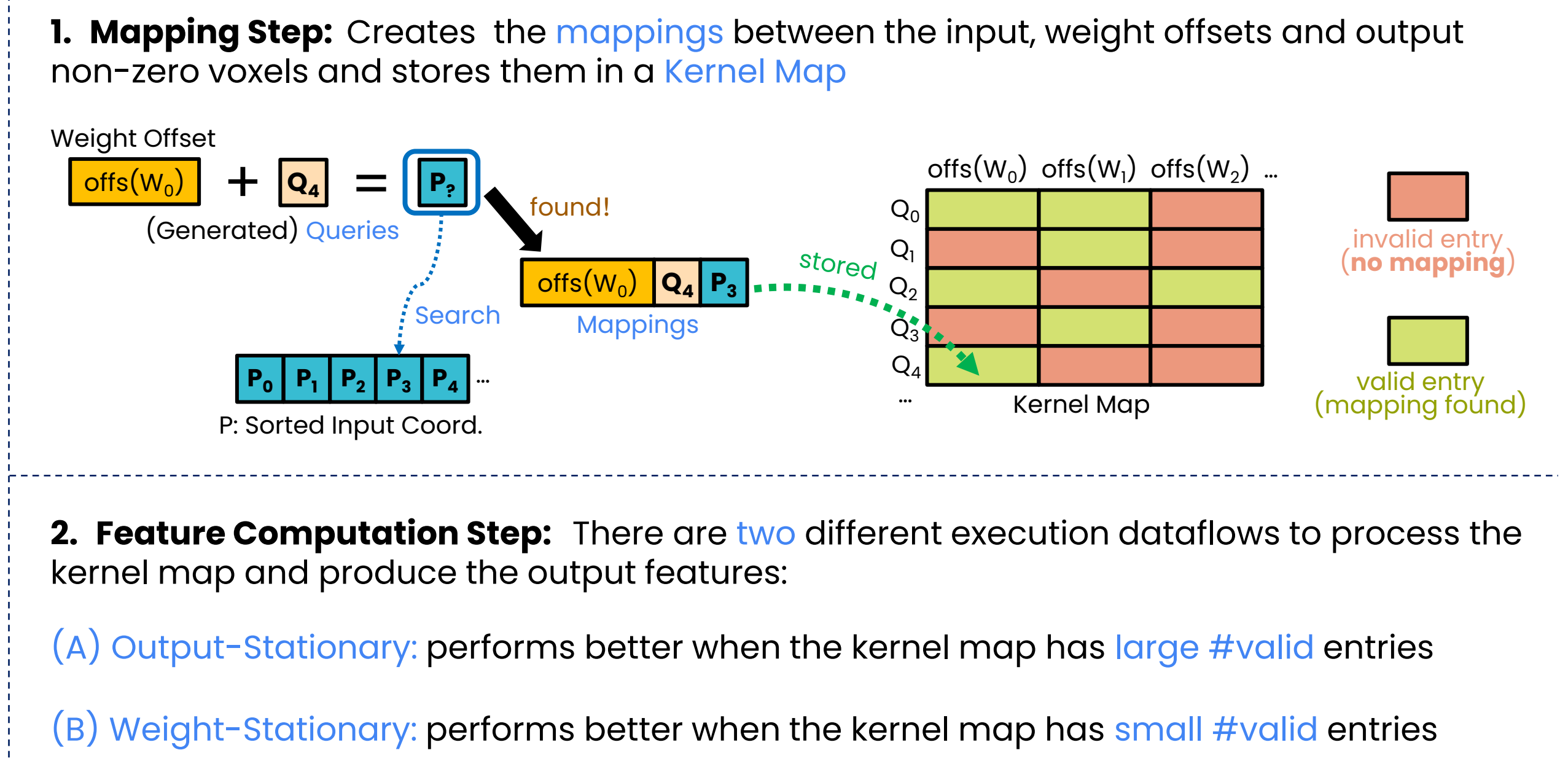
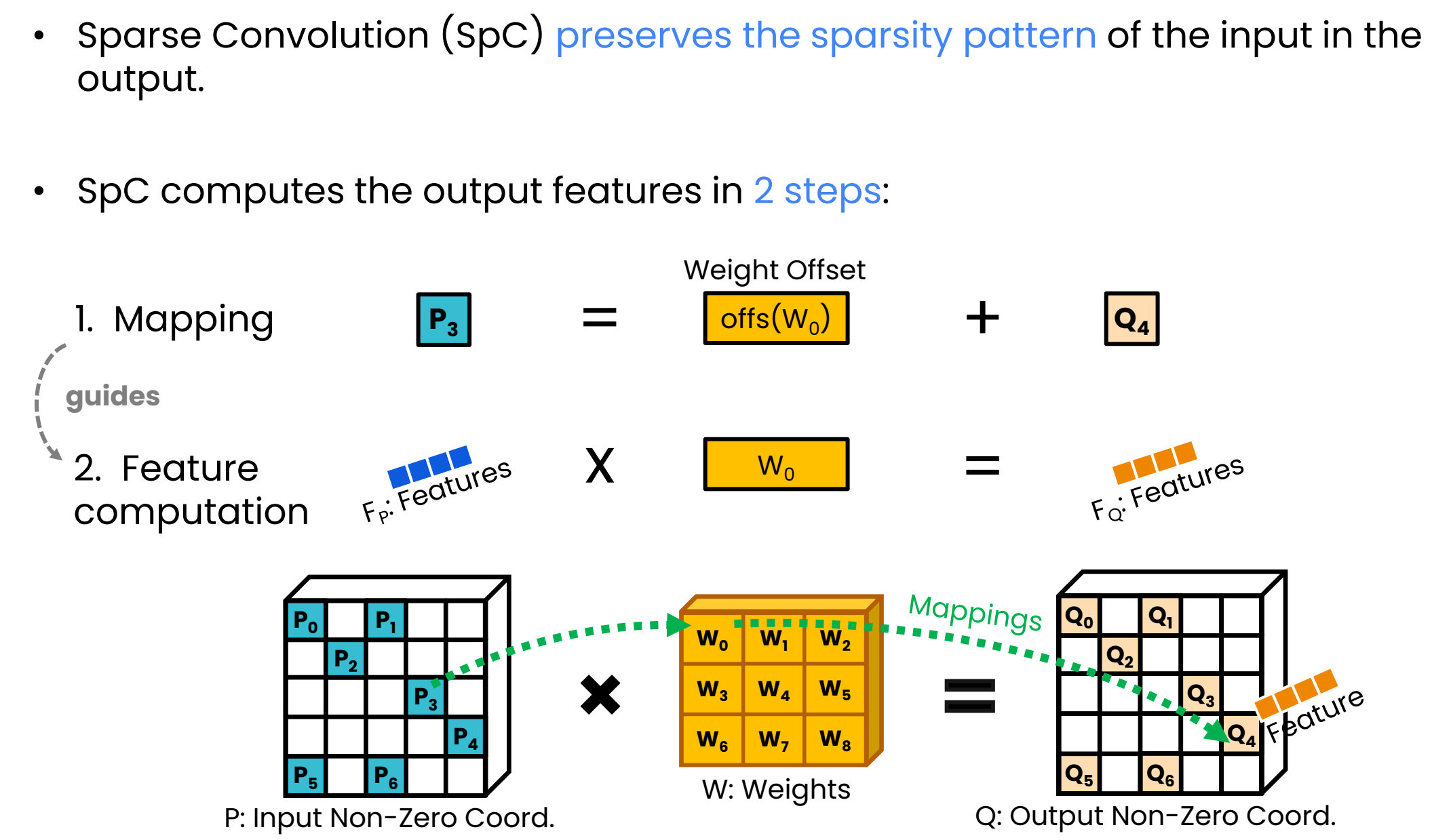
1. Summary

- Sparse Convolution (SpC) is a **key operator in point cloud networks** used in autonomous driving, robotics, and AR/VR.
- SpC operates on voxel data that have inherent structural **properties**; however, existing works **do not exploit** them.
- Spira** is the first SpC engine that **intelligently leverages** structural properties of voxel data, unlocking **significant** performance gains.
- Spira improves inference performance by **1.68x** on average across **six** GPU architectures.

2. Background on Voxel Data



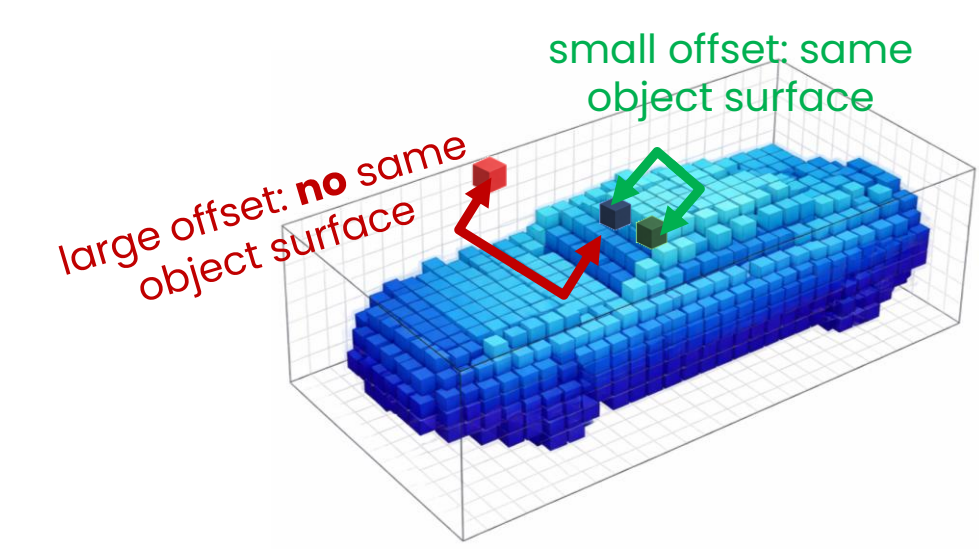
3. Background on Sparse Convolution



4. Motivation: Voxel Data Properties

We identify three important **properties** on real-world voxel data:

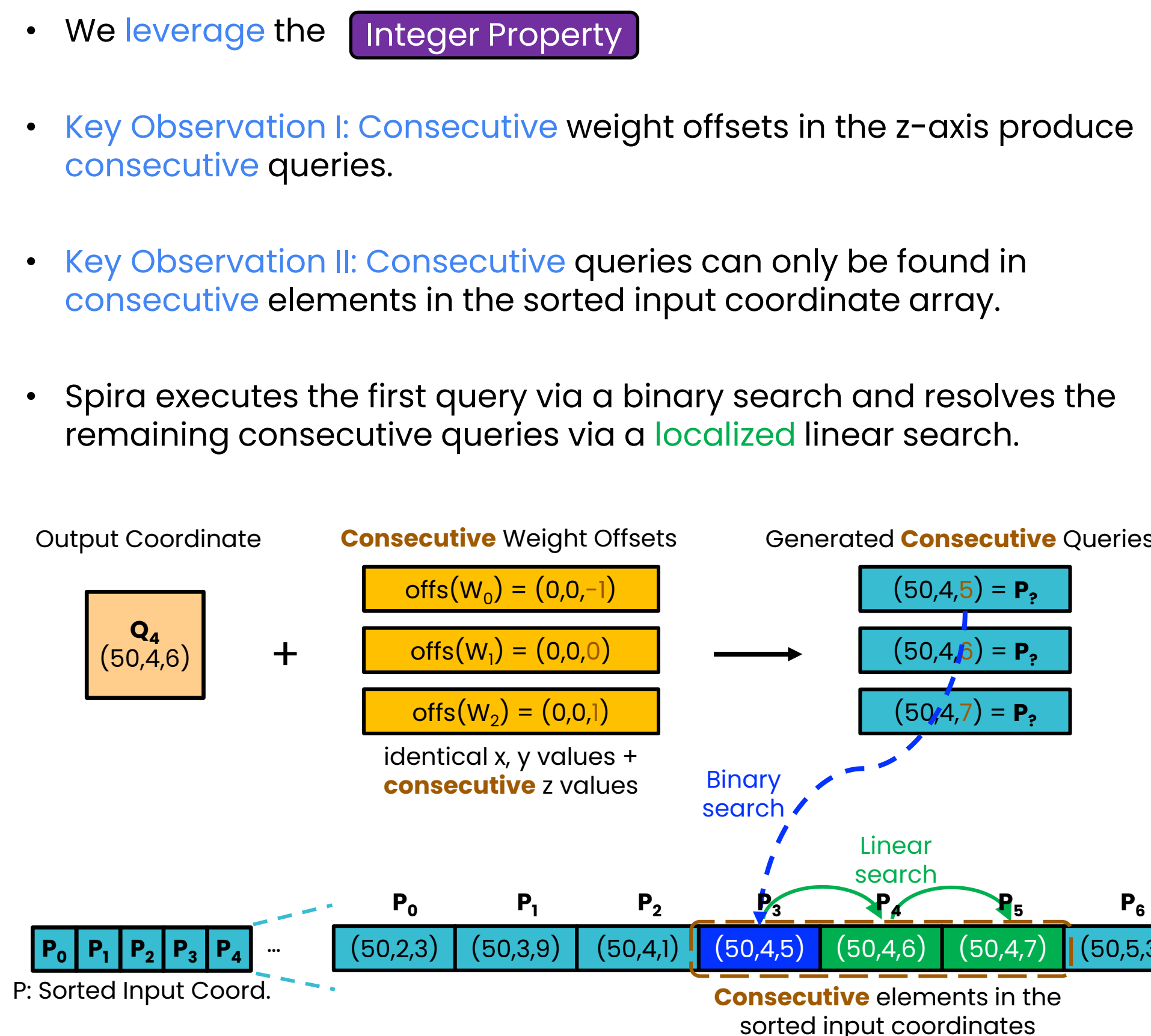
- Integer Property**: Voxel coordinates are **integer-valued**.
- Bounded Property**: Voxel coordinates are **spatially constrained** due to technological sensor limits.
- Neighboring Property**: Neighboring voxel coordinates of the same object surface are **likely** to exist in **small** offset displacements.



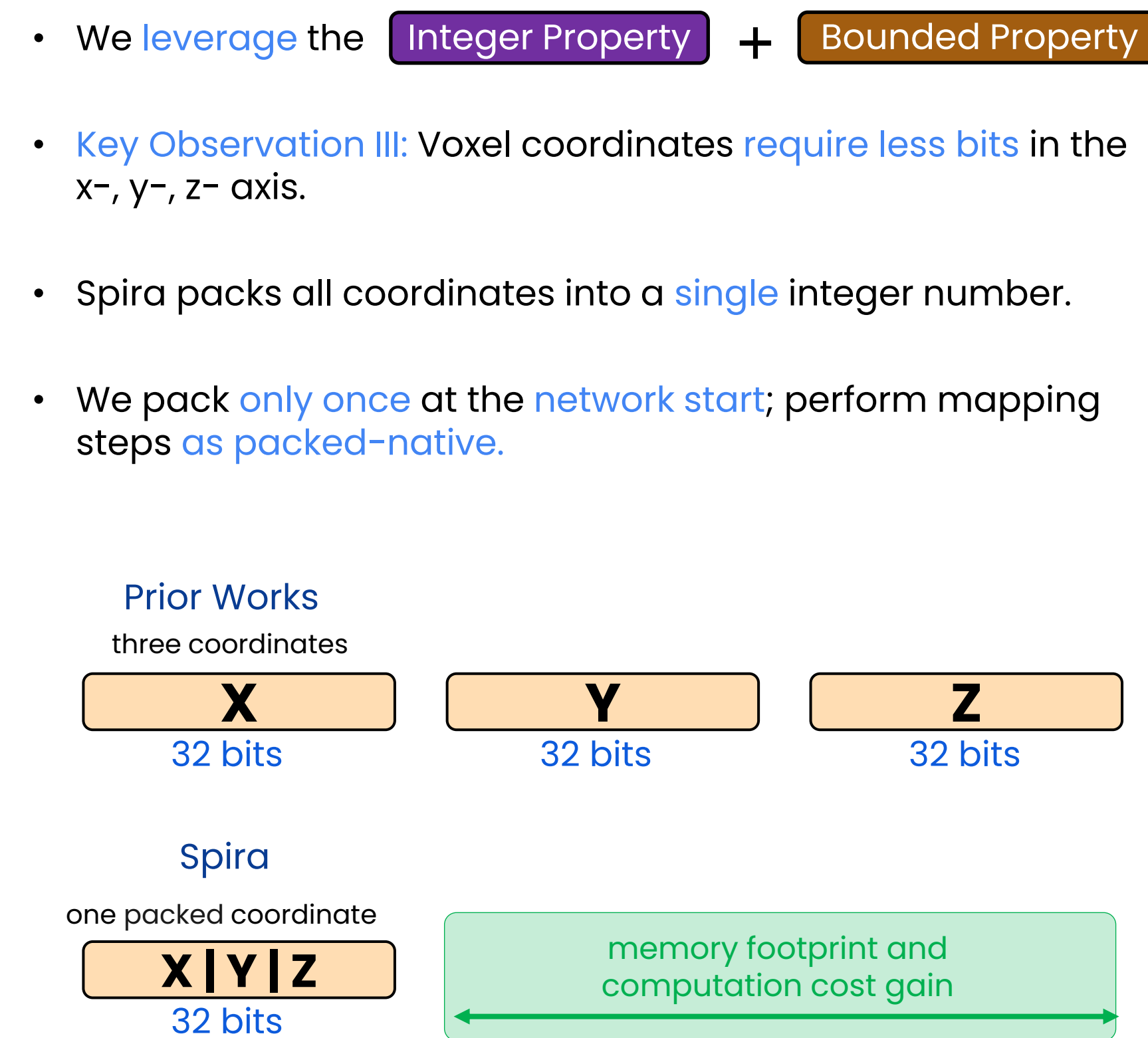
Existing SpC engines [e.g., TorchSparse++²³, Minuet²⁴] **do not leverage** voxel data properties

5. Spira Design

One-Shot Z-Delta Search in Mapping

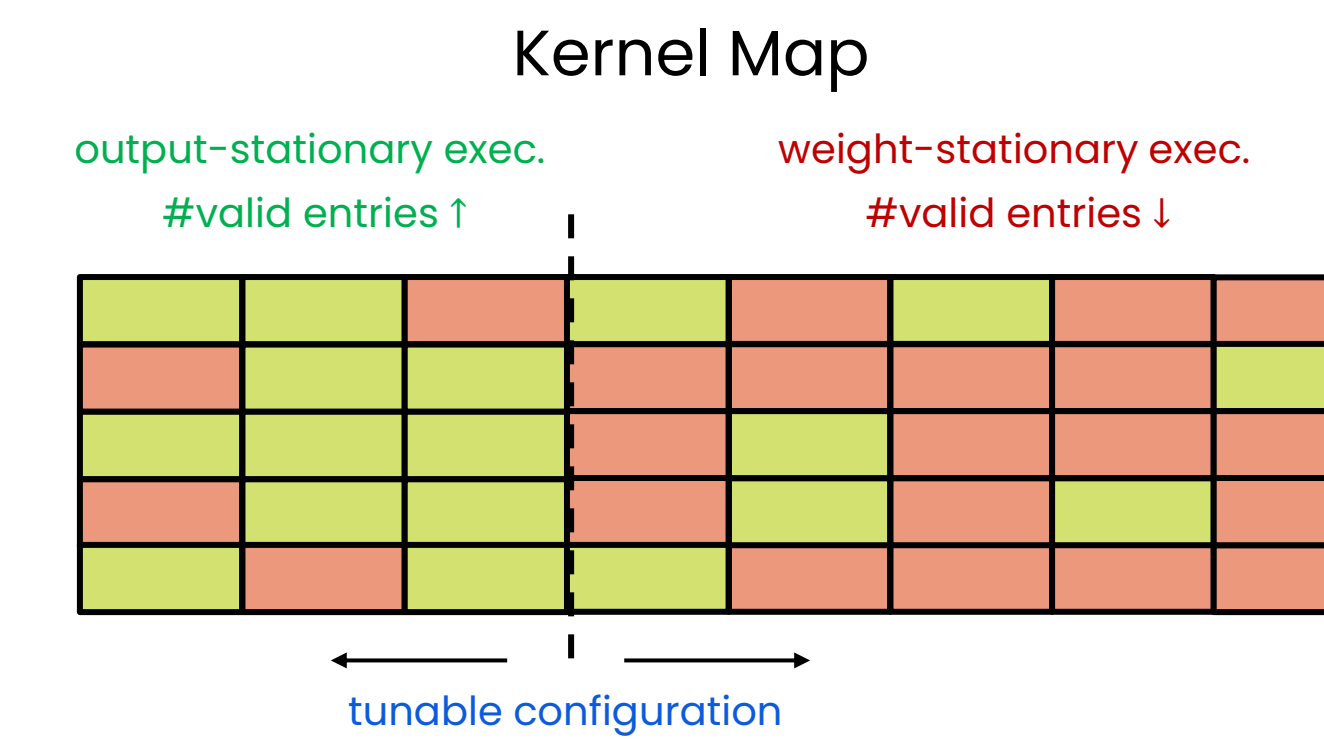


Packed-Native Mapping Step



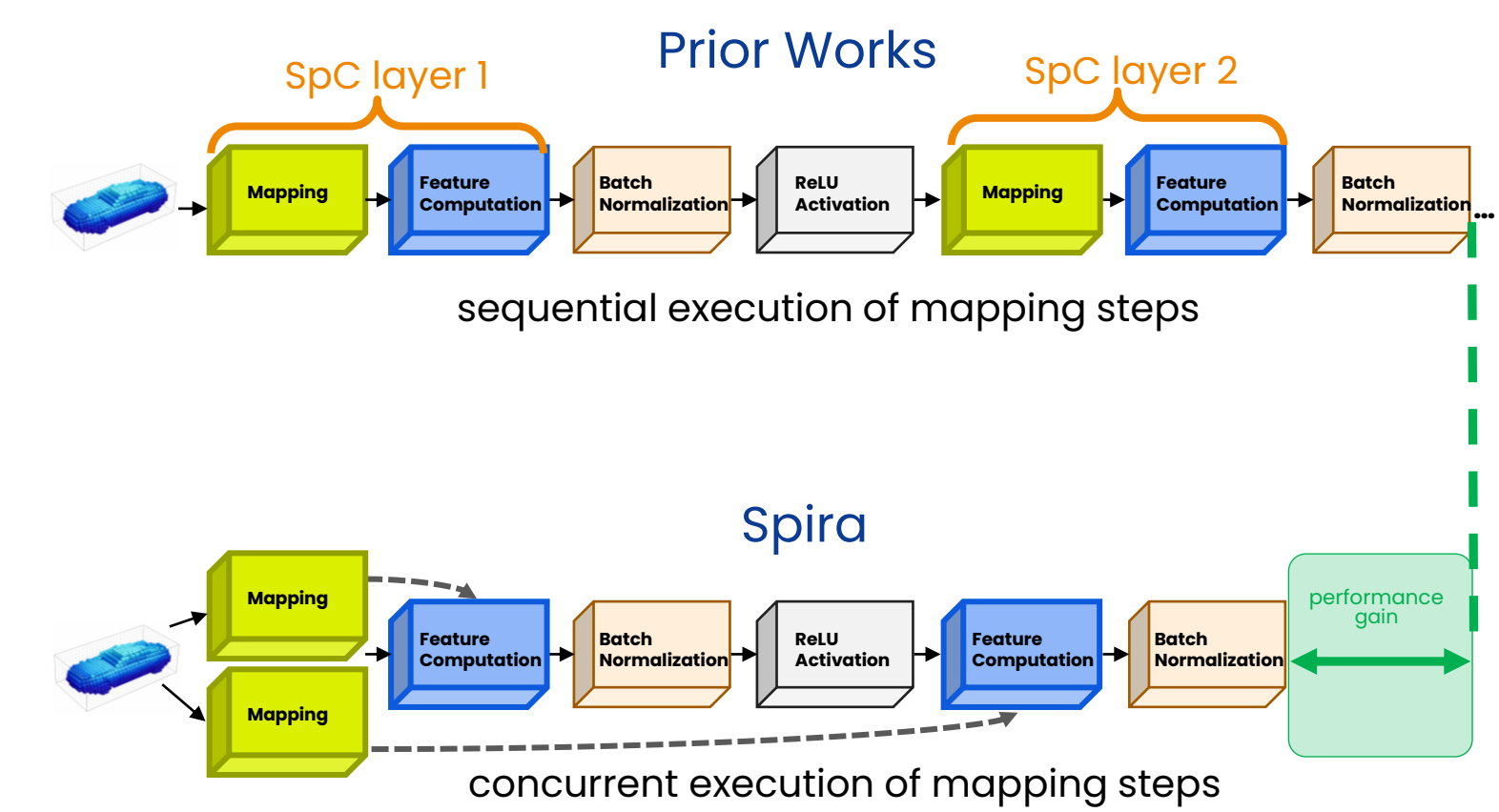
Adaptive Hybrid Dataflow Execution

- We **leverage** the **Neighboring Property**.
- Key Observation IV**: Kernel map columns of weight offsets that are associated with **smaller** offset displacements have **larger #valid entries**.
- Different** weight offsets in kernel map can be processed with **either** output- or weight-stationary **dataflow** execution.
- The configuration is tunable, allowing a **best-fit** execution.



Network-Wide Mapping

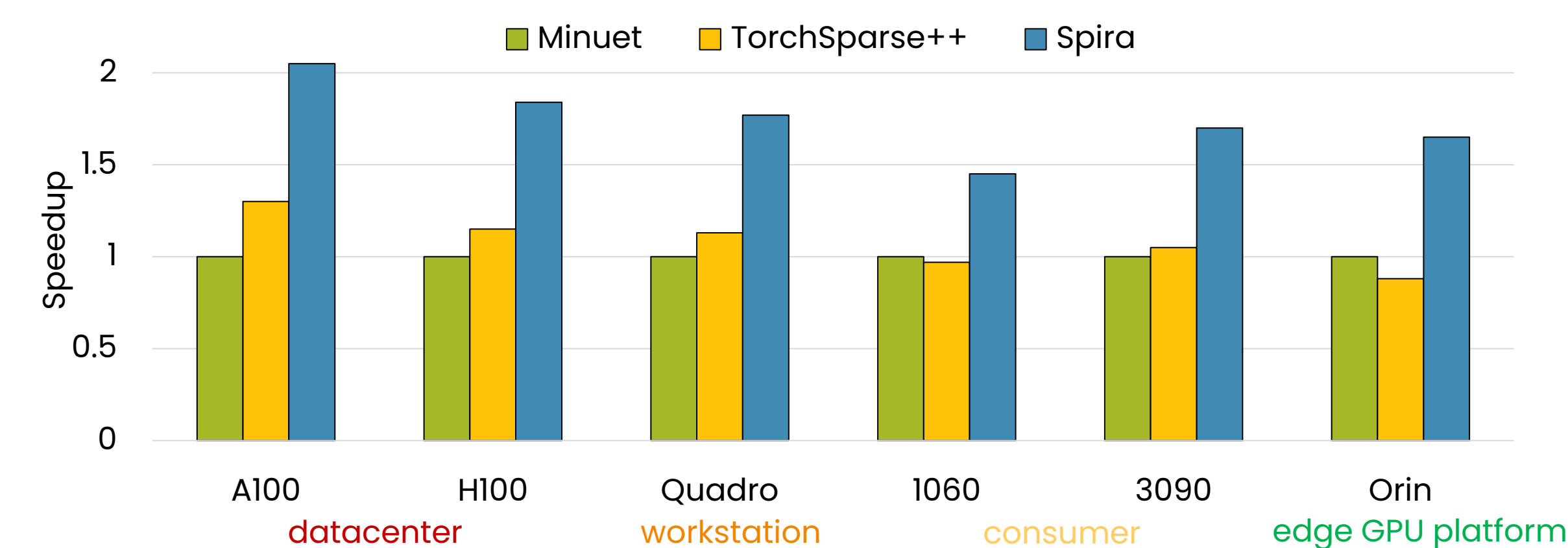
- Key Observation V**: Mapping steps have **no true dependencies** neither with feature computation steps nor between them.
- Spira executes all mapping steps **concurrently** at **network start** and significantly **improves** GPU SM utilization.



6. Methodology

- 3x** Point Cloud Networks: ResNet, Large ResNet, UNet
- 3x** Point Cloud Datasets: Waymo, SemanticKITTI, ScanNet
- 6x** GPUs:
 - A100, H100: **datacenter**
 - Quadro RTX 5000: **workstation**
 - GTX 1060, RTX 3090: **consumer**
 - Jetson Orin AGX: **edge GPU platform**
- Comparison Points
 - **TorchSparse++** [Tang et al., MICRO'23]
 - **Minuet** [Yang et al., EuroSys'24]

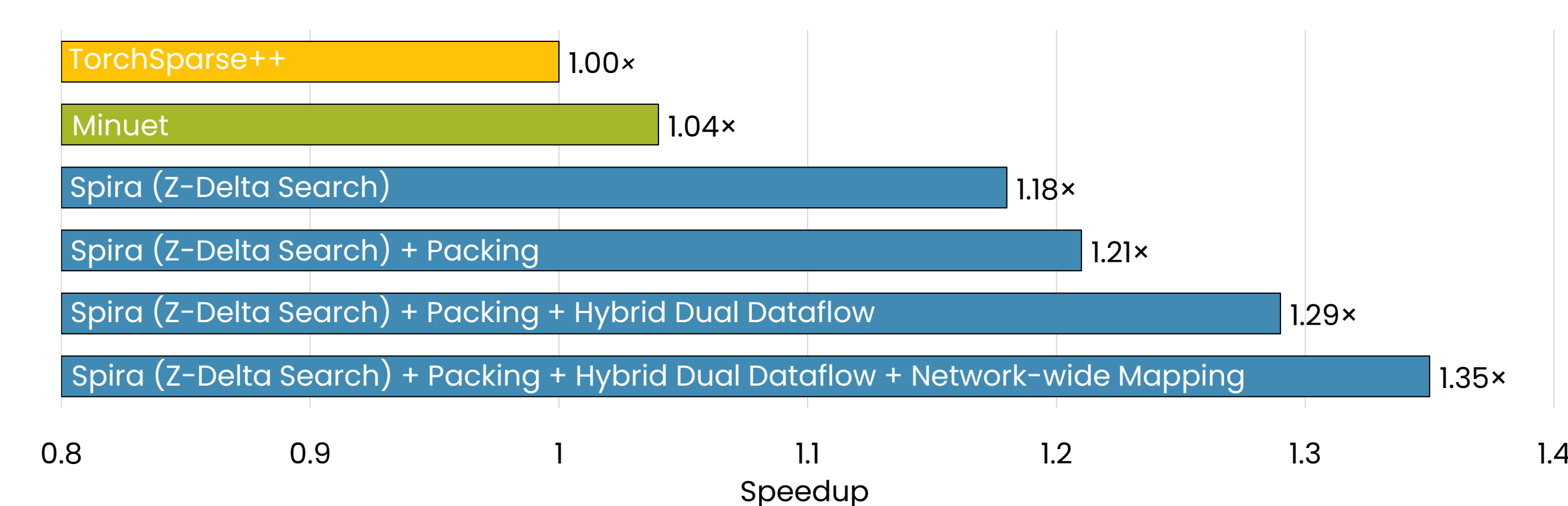
Averaged inference performance across 3 networks and 3 datasets



Spira significantly outperforms existing works by **1.68x** averaged across six GPUs

7. Evaluation

Performance benefits of Spira's key ideas in UNet



Largest performance gain comes from Spira's Z-Delta search algorithm (**1.18x** over TorchSparse++)

Officially artifact evaluated as available, functional and reproducible.

