

PointCNN++: Performant Convolution on Native Points

Lihan Li^{1,2,*} Haofeng Zhong^{1,2,*} Rui Bu² Mingchao Sun³

Wenzheng Chen^{1†} Baoquan Chen^{1†} Yangyan Li^{2†}

¹ Peking University ² Ant Group ³ AMAP

{hfhzhong, wenzhengchen, baoquan}@pku.edu.cn, {burui.br, yangyan.lyy}@antgroup.com,
 lh_li@stu.pku.edu.cn, sun.mc@outlook.com

Abstract

Existing convolutional learning methods for 3D point cloud data are divided into two paradigms: point-based methods that preserve geometric precision but often face performance challenges, and voxel-based methods that achieve high efficiency through quantization at the cost of geometric fidelity. This loss of precision is a critical bottleneck for tasks such as point cloud registration. We propose PointCNN++, a novel architectural design that fundamentally mitigates this precision-performance trade-off. It **generalizes sparse convolution from voxels to points**, treating voxel-based convolution as a specialized, degraded case of our more general point-based convolution. First, we introduce a point-centric convolution where the receptive field is centered on the original, high-precision point coordinates. Second, to make this high-fidelity operation performant, we design a computational strategy that operates **natively** on points. We formulate the convolution on native points as a Matrix-Vector Multiplication and Reduction (MVMR) problem, for which we develop a dedicated, highly-optimized GPU kernel. Experiments demonstrate that PointCNN++ **uses an order of magnitude less memory and is several times faster** than representative point-based methods. Furthermore, when used as a simple replacement for the voxel-based backbones it generalizes, it **significantly improves point cloud registration and semantic segmentation accuracies while proving both more memory-efficient and faster**. PointCNN++ shows that preserving geometric detail and achieving high performance are not mutually exclusive, paving the way for a new class of 3D learning with high fidelity and efficiency.

Our code is publicly available at: <https://github.com/ant-research/pointintelligence>

1. Introduction

The proliferation of 3D sensing technologies has established point clouds as a primary data modality in domains like autonomous driving [15, 20, 52], robotics [23, 26, 37, 60], and augmented reality [3, 35, 36, 48]. As a direct repre-

sentation of world geometry, a point cloud’s value is intrinsically tied to its high-fidelity spatial information. This very characteristic—inherently sparse and irregular data—however, creates a fundamental challenge for modern computational hardware [2, 22, 32, 33] and software [1, 21, 31, 34] that are heavily optimized for dense and regular data.

To address this challenge, the study of convolution on point clouds has developed two competing paradigms, each embodying a significant compromise. The most prevalent, the voxel-based approach (Figure 1 II), resolves the conflict by **forcefully restricting the convolution on voxel grids**. This approach firstly applies a *global voxelization* to quantize an entire continuous space of interest into a set of sparse voxels, from which sparse convolutional operators are applied. While the performance issue of such convolution has been substantially addressed by leveraging the sparse nature of the data, with representative work from O-CNN [46], SP-Conv [11], to MinkowskiEngine [7], the quantization is an inherently lossy sampling operation, resulting in impaired fine geometric details—representing high-frequency spatial signals. The act of the global voxelization in the beginning of the process establishes an a priori error floor determined by the voxel size, posing a critical bottleneck to tasks like high-precision registration that depend on sub-voxel feature uniqueness.

An alternative philosophy, the point-based paradigm, attempts to preserve the data’s integrity by a relatively more gentle, usually learned, transformation of irregular points into a regular dense tensors, from which the convolution is then applied. Such transform-then-convolve approach (Figure 1 III) is shown to be effective in harvesting the fidelity in point clouds, as shown in the notable representative work of PointCNN [25] and KPConv [44]. However, while the convolution on the dense tensor is indeed computationally efficient, **significant inefficiency arises in the transformation from irregular to regular itself**. Such transformation introduces extra computation, parameters to learn and significant memory access—a major source of inefficiency in point cloud related computation on GPUs [27, 42]. Moreover, the implementation of such point-based methods often involves frequent usage of padding operations, which often intensifies the inefficiency.

This paper posits that this prevailing trade-off should not

#Equal contributions. † Corresponding author. * Work done during internship at Ant Group.

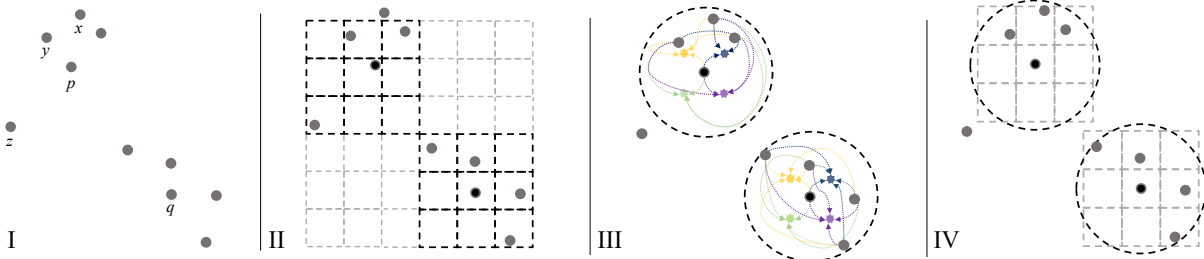


Figure 1. A 2D illustration of convolutional learning for point cloud (I) with voxel-based methods (II), transform-then-convolve methods (III) and convolution on native points (IV). When a voxel center happens to be on an original point (the rare case, as depicted by q), the difference between (II) and (IV) is minimal. However, in the general cases, due to the forceful restricting of computation on voxel grids in (II), several problems arise: 1. misalignment between original points (e.g., point p) and convolution centers, 2. inaccurate neighborhood inclusion (e.g., x , instead of z , should be in the neighborhood of p), and 3. inaccurate convolution kernel usage (e.g., the feature associated with point x is more appropriate for being convolved with the upper-left kernel as shown in (IV), instead of the upper-middle kernel as shown in (II). IV preserves geometric precision as those in III, while avoiding the cumbersome irregular-to-regular “transformation”.

be viewed as a permanent compromise, but rather as a conflict that can be mitigated through holistic computational design. We introduce PointCNN++ (Figure 1 IV), an architectural design that resolves this conflict by advancing a new computational paradigm. Instead of **restricting the convolution** or **introducing extra transformation**, we design the operator and compute kernel as an integrated system, purpose-built for the high-performance processing of native point clouds, without any degradations nor superfluties. Our approach begins by generalizing sparse convolution from discrete voxels to continuous points, centering operations on the true, high-precision coordinates as much as possible. As the last operation, we use a *local and adaptive voxelization* to pair the data to be convolved with the discrete convolution kernels is used as the last operation to minimize the fidelity loss. From this perspective, voxel-based convolution is merely a quantized, degraded special case of our more fundamental point-centric operator. A more general operator, however, is only practical if it is performant. To this end, we formulate the convolution on native points as a Matrix-Vector Multiplication and Reduction (MVMR) problem. With this formulation, we drew inspiration from an efficient algorithm [41] for the Matrix-Vector Multiplication sub-problem to develop a dedicated, highly-optimized GPU kernel for the full problem.

In this paper, we introduce PointCNN++, an architectural design that mitigates the long-standing conflict between precision and performance in convolutional learning for 3D point cloud data. We empirically demonstrate that PointCNN++ not only excels on precision-critical tasks such as registration but also delivers state-of-the-art results on large-scale semantic segmentation benchmarks. Moreover, it is more memory-efficient while being even faster than existing approaches. This work proves that by holistically designing the computational system of point cloud data, achieving high geometric fidelity and high performance is not a mutually exclusive goal.

2. Related Work

2.1. Feature Learning for 3D Point Cloud

Deep learning on point clouds has evolved through several operator paradigms. The seminal PointNet [38] processes

each point independently with shared MLPs before global aggregation, while its successor PointNet++ [39] introduces a hierarchical structure by recursively applying this process on nested point subsets. A central line of work focuses on generalizing convolution operator for 2D image data into processing 3D data. Voxel-based methods first discretize space for efficient 3D CNNs; for instance, VoxelNet [58] learns a unified feature representation for points within each voxel. PointCNN [25] proposes a learned \mathcal{X} -Transform to permute local point features into a canonical order before applying a convolution. PointConv [50] learns continuous kernel weights from relative coordinates using an MLP. KPConv [44] uses a set of rigid, learnable kernel points to apply spatially-aware weights. DGCNN [47] constructs dynamic graphs in feature space, and applies convolution on the feature space neighborhoods.

While Transformer has been widely used in processing 3D point cloud data, one thing to note is that they often operate on the features extracted with convolutional backbones [19, 24, 29, 30, 55, 56, 59], without which their effectiveness has not been widely demonstrated. Our work is most closely related to convolutional backbones in general, either used alone or as a part of a larger architecture, for serving general feature learning purposes on point cloud data. We show that with our systematic design, the performance advantage of the voxel-based methods and the precision advantage of point-based methods could be combined.

2.2. Performant Computation on Sparse Voxels

The performance of 3D deep learning is critically dependent on their underlying computational systems. The sparse nature of 3D data is extensively explored to efficiently process data while avoiding computation on empty space. SparseConvNet [14] pioneered the work in this domain, introducing the use of hashmaps to manage the coordinates of active voxels. Building on this, SpConv [51] proposed a highly-optimized grid-based map search and formalized the influential gather-matmul-scatter dataflow. MinkowskiEngine [7] significantly improved upon SparseConvNet’s hashmap implementation to reduce latency and introduced an alternative fetch-on-demand dataflow. More recent frameworks like TorchSparse [42] and TorchSparse++ [43] have continued to

push performance boundaries by developing novel traversal algorithms, hand-tuned CUDA kernels, and co-optimizing data structures and workloads specifically for GPU architectures. *fVDB* [49] integrated these concepts and further developed a flexible framework accepting various popular 3D data representations into the unified voxel-based representation upon which a rich set of operators is based. While existing high-performance methods typically exploit sparsity via voxelization, we show, however, that this is not a necessary coupling. Our approach achieves comparable computational efficiencies by leveraging sparsity directly on native point cloud representations, thereby bypassing the geometric precision loss inherent to voxelization.

3. Method

This section details the design principles of PointCNN++. Our methodology resolves the long-standing precision-performance dichotomy in 3D deep learning through a synergistic co-design of the core operator and its underlying computational system. We begin by establishing a generalized view of convolution and situating prior work within this framework in Sec. 3.1. We then introduce our point-centric convolution, highlighting its properties as a powerful generalization in Sec. 3.2. Finally, we describe the performant system design in Sec. 3.3, from the data representation to the optimized GPU kernel, that makes this high-fidelity operator computationally feasible at scale.

3.1. Preliminaries on Convolution

3.1.1. A Generalized View of Convolution

At its core, a convolution operator computes a feature $\mathbf{F}_i^{\text{out}} \in \mathbb{R}^{C_{\text{out}}}$ at location $\mathbf{P}_i^{\text{out}}$ by aggregating information from neighborhood locations $\{\mathbf{P}_j^{\text{in}}\}$, each associated with a feature $\mathbf{F}_j^{\text{in}} \in \mathbb{R}^{C_{\text{in}}}$, with the help of some learnable kernels $\{\mathbf{W}_k \in \mathbb{R}^{C_{\text{out}} \times C_{\text{in}}}\}$. C_{in} and C_{out} are dimensions, or channels, of the input and output features, respectively. The computational interactions among $\mathbf{F}_i^{\text{out}}$, \mathbf{F}_j^{in} , and \mathbf{W}_k , could be completely depicted by a triplet set $\mathcal{T} = \{(i, j, k)\}$, where a triplet (i, j, k) denotes that for convolution at location $\mathbf{P}_i^{\text{out}}$, location \mathbf{P}_j^{in} is in its neighborhood, and \mathbf{W}_k is the corresponding kernel to use for this neighborhood location. Thus, a generalized view of convolution could be formulated as:

$$\mathbf{F}_i^{\text{out}} = \sum_{(i,j,k) \in \mathcal{T}} \mathbf{W}_k \times \mathbf{F}_j^{\text{in}}. \quad (1)$$

The essence of different convolution types lies in how the triplet set $\mathcal{T} = \{(i, j, k)\}$ is constructed. More specifically, there are three major considerations in the construction of such triplets: ❶ How the locations of output features are placed? — Considerations on defining convolution centers $\{\mathbf{P}_i^{\text{out}}\}$ to use. ❷ How are the neighborhood measured? — Considerations on the pairing of $\{(i, j)\}$ in the triplet. ❸ Which kernel to use for a location in the neighborhood? — Considerations on the pairing of $\{(j, k)\}$ in the triplet.

3.1.2. Convolution on 2D Images

Convolution on 2D images is a specialized instance of Eqn. (1), where the triplet set $\{(i, j, k)\}$ is constructed by

taking fully considerations on the dense and regular nature of the image data that is ubiquitously represented as discrete pixels: ❶ $\{\mathbf{P}_i^{\text{out}}\}$ are the pixels. ❷ $\{(i, j)\}$ are predefined as the pixels in patches around $\{\mathbf{P}_i^{\text{out}}\}$ are taken as neighborhoods, *i.e.*, Chebyshev distance based neighbors. ❸ For $\{(j, k)\}$, the kernels and neighborhoods are often of the same spatial resolution, thus $\{\mathbf{W}_k\}$ are paired with $\{\mathbf{P}_j^{\text{in}}\}$ if they share the same spatial locations.

Significant efforts have been made to achieve efficient GPU algorithms for the specialized triplet set \mathcal{T} of image convolution. In the early development stage of deep learning, image convolution was *not* a highly optimized standard operation on GPU. One way to achieve efficient computation was to materialize \mathbf{F}_j^{in} , thus lower the convolution into an highly optimized general matrix-matrix multiplication (GEMM), as practiced in the early version of Caffe [21]. However, as detailed in cuDNN [6], such materialization inevitably introduces significant amount of extra GPU memory usage, which could be addressed by lazily materializing in on-chip memory only.

3.1.3. Convolution on 3D Sparse Voxels

With a *global voxelization* of the entire 3D space of interest, often in the form of point cloud, the space could be sampled into sparse voxels — a representation that is a generalization of 2D pixels into 3D. By adding one dimension to the spatial coordinates, and relaxing the assumption of dense and regular tensor, convolution operator can be generalized from 2D images into 3D sparse voxels. Yet, convolution on 3D sparse voxels is still a specialized instance of Eqn. (1), as illustrated in Figure 1 II: ❶ $\{\mathbf{P}_i^{\text{out}}\}$ are the discrete voxels. ❷ $\{(i, j)\}$ are constructed by taking the non-empty voxels around $\{\mathbf{P}_i^{\text{out}}\}$ as neighborhoods, *i.e.*, again, Chebyshev distance based neighbors. ❸ Same as that in convolution on 2D images, $\{(j, k)\}$ are constructed by corresponding same spatial locations.

While efficient GPU algorithms of voxel-based convolution have been proposed in notable representative work from O-CNN [46], SPConv [11], to MinkowskiEngine [8, 10]. There are inherent drawbacks arise from the definition of voxel-based convolution: (1) placing $\{\mathbf{P}_i^{\text{out}}\}$ at the discrete voxels is at the cost of sacrificing fidelity from the original point cloud; (2) neighborhood construction in Chebyshev distance with imprecise centers intensifies the fidelity erosion; (3) the fineness of kernels are coupled with the fineness of voxelization. A more detailed discussion, and addressing, of such drawbacks are given after the introduction of our method in the next section.

3.2. Convolution on Native Points

Definition. Our design philosophy is to fully harvest the fidelity in the input point cloud. As illustrated in Figure 1 IV, We start by ❶ placing the convolution locations right at the original high-precision points, thus $\{\mathbf{P}_i^{\text{out}} \in \mathbb{R}^3\}$ are true, continuous coordinates in our formulation. Then, ❷ the pairing of $\{(i, j)\}$ could be constructed based on neighborhood search using the precise convolution centers with appropriate distance metrics on continuous space. Finally ❸, for constructing $\{(j, k)\}$, a *local voxelization*, centered right at $\mathbf{P}_i^{\text{out}}$

of the same spatial resolution as the kernels is applied in each neighborhood region, yielding a correspondence between the kernels $\{\mathbf{W}_k\}$ and $\{\mathbf{P}_j^{\text{in}}\}$ in the neighborhood voxels.

Advantages. It is clear that both ❶ and ❷ operates on the full precision of the original point cloud, the quality of the neighborhood regions are well preserved. It is only at the final step, ❸, with the convolution centers align exactly with neighborhood region centers, a *local voxelization* is applied. Such a local voxelization is adaptive to each neighborhood region, resulting in a higher quality mapping between $\{\mathbf{W}_k\}$ and $\{\mathbf{P}_j^{\text{in}}\}$. Note that, with given neighborhood regions, it is the appropriate fineness of the convolution kernels that defines the resolution of the local voxelization, rather than being coupled to the fineness of the global voxelization as that in convolution on sparse voxels.

Voxel Convolution as a Special Case. With the definition of convolution on native points, now we show that convolution on voxels is a special case of convolution on native points, by demonstrating the three degradations have to take that convert convolution on native points into convolution on voxels. First of all, such conversion inevitably introduce a global voxelization to generate voxels that are not necessary for convolution on the native points, but mandatory for convolution on voxels. Degradation ❶, instead of using the original points as the convolution centers, snapping the convolution centers into the centers of the voxels. Degradation ❷, instead of searching neighbors with the original points as centers within other points, searching neighbors in the voxel space. Degradation ❸, instead of choosing appropriate fineness thus resolution of the kernels, using the resolution that is coupled to the size of receptive field in the global voxelized space.

For degradation ❸, in other words, convolution on native points decouples kernel resolution from the receptive field’s physical span, whereas convolution on voxels couples it to the voxel resolution in the receptive field. We demonstrate the difference with some examples. In convolution on native points, given a neighborhood of points, kernels with fineness of either 3^3 or 5^3 could be used up to demand, as the continuous coordinates could result in any of the local voxels. In contrast, in convolution on voxels, if a neighborhood is constructed with 3^3 voxels based on the global voxelization, it does not make sense to use kernels of 5^3 resolution to convolve with such a neighborhood, as the neighborhood has been quantized into the coarser 3^3 voxels, leaving the extra fineness of 5^3 kernels useless.

3.3. Performant Systematic Design

The flexibility of convolution on native points presents a significant computational challenge. To make this high-fidelity design feasible at scale, we introduce a performant systematic design that spans from the fundamental computation abstraction to a highly-optimized GPU kernel.

3.3.1. Computational Abstraction: MVMR

Same as convolution on 2D images, convolution on native points is an instance of Eqn. (1). A tempting strategy for efficient computation of it is to adapt the *im2col* like memory

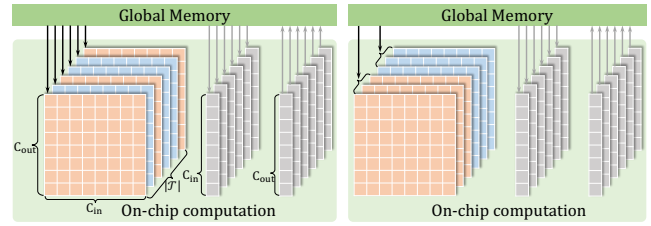


Figure 2. A brute-force MVM computation inefficiently reads $\mathbf{W}_k \in \mathbb{R}^{C_{\text{out}} \times C_{\text{in}}}$ from global memory $|\mathcal{T}|$ times—once for every triplet(left). Sorting the triplets by k optimizes this. Ideally, each unique \mathbf{W}_k is loaded just once into on-chip memory and reused for all its associated computations(right).

materialization technique once was used from image convolution, to lower the computation into a large GEMM problem. The extra memory usage introduced by materialization of the kernels is $K \times C_{\text{out}} \times C_{\text{in}}$, where K is the number of kernel matrices to use, and $K = t^3$ if a t resolution convolution kernel on each spatial axis is used. While extra memory usage for materializing the kernels are manageable, the materialization of neighborhoods would result in extra memory usage of $N_{\text{in}} \times K \times C_{\text{in}}$, where N_{in} is the size of the input point cloud, which is K times larger than the original data. Note that, the materialization of neighborhoods would “fill in” every neighborhood to a uniform maximum size with padding, causing a memory footprint increase far greater than that seen in image convolution and rendering the strategy impractical. Clearly, the strategy that was effective in convolution on 2D images is certainly not acceptable in convolution for native points in terms of extra memory usage, not to mention the latency introduced by the involved memory read and write.

Considering that point cloud is inherently irregular, to the best of our knowledge, there is no special structure to leverage for an alternative strategy in lowering the computation of convolution on native points into a single performant primitive. Inspired by the way cuDNN [6] addresses the extra memory usage problem and its effectiveness, we hypothesize that a computation abstraction that literally follows the general Eqn. (1) formulation has the potential of being performant while achieving zero extra memory usage.

We abstract the computation of Eqn. (1) as a Matrix-Vector Multiplication and Reduction (MVMR) problem, with the observation that MVMR naturally decomposes into two distinct components: the $\mathbf{W}_k \times \mathbf{F}_j^{\text{in}}$ component is a Matrix-Vector Multiplication(MVM), which itself is a popular operation with rich studies on achieving performant computation on GPUs, and the rest component is a Reduction — another example of cases with performant solutions. We develop a unified efficient GPU algorithm by taking inspirations from existing studies on the two sub-problems.

3.3.2. Efficient GPU Algorithm of MVMR

First of all, a native MVM implementation is required because MVMR involves many small MVM computations. Calling a standard library routine for each computation individually from an outer loop would incur prohibitive launch overhead. Second, while the amount of multiplication and

addition is fixed, the key to a performant implementation of MVMR lies in optimizing memory access patterns. GPU performance is critically dependent on: a) data locality, *i.e.*, how much data is accessed from fast, on-chip resources versus the much slower, off-chip global memory. b) coalesced memory access, where threads within a hardware warp access contiguous memory locations to maximize bandwidth.

The output of the convolution is a feature tensor $\mathbf{F}^{\text{out}} \in \mathbb{R}^{N_{\text{out}} \times C_{\text{out}}}$, where N_{out} is the size of the output point cloud. It is computed from three inputs to the convolution operator: 1) the input feature tensor $\mathbf{F}^{\text{in}} \in \mathbb{R}^{N_{\text{in}} \times C_{\text{in}}}$; 2) the weight matrix $\mathbf{W} \in \mathbb{R}^{K \times C_{\text{out}} \times C_{\text{in}}}$; 3) the collection of triplets $\mathcal{T} = \{(i, j, k)\}$, constructed based on the convolution definition, where i, j and k are integers in range $[0, N_{\text{out}}]$, $[0, N_{\text{in}}]$, and $[0, K)$, respectively. The computational interactions among \mathbf{F}^{out} , \mathbf{W} and \mathbf{F}^{in} are completely depicted by the triplets \mathcal{T} . In most cases, $K \ll N_{\text{in}} \simeq N_{\text{out}} \ll |\mathcal{T}|$.

Note that the process of each triplet is independent of others, as long as the conflict on Reduction write to \mathbf{F}^{out} is governed by atomic primitive. Therefore, there is a straightforward parallelism strategy to accomplish the computation of all the triplets: let each thread handle the computation of one triplet. While this strategy seems ideal in terms of parallelism, significant amount of slow global memory access is involved, as illustrated in Figure 2 (left):

$$|\mathcal{T}| \times \left(\underbrace{C_{\text{out}} \times C_{\text{in}}}_{\text{Read from } \mathbf{W}} + \underbrace{C_{\text{in}}}_{\text{Read from } \mathbf{F}^{\text{in}}} + \underbrace{C_{\text{out}}}_{\text{Atomic Write to } \mathbf{F}^{\text{out}}} \right). \quad (2)$$

Clearly, the global memory read from \mathbf{W} is a bottleneck, and therefore, this access pattern must be addressed to improve performance. Mathematically, the final results are invariant to the processing order of the triplets¹, but a sorting of the triplets could significantly change the memory access patterns, as illustrated in Figure 2 (right). More specifically, given the triplet list $\mathcal{T} = \{(i, j, k)\}$ sorted by k , if it is split into consecutive groups of length L , based on Pigeonhole Principle, *most* of the groups contain triplets that share the same k value. More specifically, in the case of $K \ll |\mathcal{T}|$, a practical approximations of the expectation number of unique k in each group is $1 + \frac{L \times K}{|\mathcal{T}|}$ ², which approaches 1, when $L \times K \ll |\mathcal{T}|$. Due to this reason, when each of such group is processed together, *almost* only one global memory read of \mathbf{W} is required for each group, and the overall global memory access could be effectively reduce to $\mathcal{O}(\frac{1}{L})$ of that in Eqn. (2):

$$\frac{|\mathcal{T}|}{L} \times \left(\underbrace{C_{\text{out}} \times C_{\text{in}}}_{\text{Read from } \mathbf{W}} + \underbrace{L \times C_{\text{in}}}_{\text{Read from } \mathbf{F}^{\text{in}}} + \underbrace{L \times C_{\text{out}}}_{\text{Atomic Write to } \mathbf{F}^{\text{out}}} \right). \quad (3)$$

Beyond the option of sorting \mathcal{T} by k , there are two alternatives: sorting by i and sorting by j . Following the analysis of the saving introduced from sorting by k , they would

¹Practically, the results may vary with order due to numerical effects (*e.g.*, floating-point arithmetic).

²This is a classic problem that combines order statistics with a variant of the coupon collector's problem. The expected number of unique values in a single group is $1 + (L - 1) \left[\frac{K(1 - (1 - \frac{1}{K})^{|\mathcal{T}|}) - 1}{|\mathcal{T}| - 1} \right]$.

Algorithm 1 MVMR Kernel for Computing Eqn. (1)

Inputs: $\mathbf{W}, \mathbf{F}^{\text{in}}, \mathcal{T}^L$.

Output: \mathbf{F}^{out} .

```

1:  $(\tilde{i}, \tilde{j}, \tilde{k}) \leftarrow \mathcal{T}_0^L$  ▷ initialize with the first triplet.
2:  $\tilde{W} \leftarrow \mathbf{W}_{\tilde{k}}, \tilde{F}^{\text{in}} \leftarrow \mathbf{F}_{\tilde{j}}^{\text{in}}$  ▷ read from global memory.
3:  $\tilde{F}^{\text{out}} \leftarrow \tilde{W} \times \tilde{F}^{\text{in}}$  ▷ fast on-chip computation.
4: for all  $(i, j, k)$  in  $\mathcal{T}_{[1, 2, \dots, L-1]}^L$  do
5:   if  $j \neq \tilde{j}$  then
6:      $\tilde{j} \leftarrow j, \tilde{F}^{\text{in}} \leftarrow \mathbf{F}_j^{\text{in}}$  ▷ read only if necessary.
7:   end if
8:   if  $k \neq \tilde{k}$  then
9:      $\tilde{k} \leftarrow k, \tilde{W} \leftarrow \mathbf{W}_{\tilde{k}}$  ▷ read only if necessary.
10:  end if
11:  if  $i \neq \tilde{i}$  then
12:    atomicAdd $(\tilde{F}^{\text{out}}, \mathbf{F}_i^{\text{out}})$  ▷ only if necessary.
13:     $\tilde{i} \leftarrow i, \tilde{F}^{\text{out}} \leftarrow \tilde{W} \times \tilde{F}^{\text{in}}$  ▷ on-chip, fast.
14:  else
15:     $\tilde{F}^{\text{out}} \leftarrow \tilde{F}^{\text{out}} + \tilde{W} \times \tilde{F}^{\text{in}}$  ▷ on-chip, fast.
16:  end if
17: end for
18: atomicAdd $(\tilde{F}^{\text{out}}, \mathbf{F}_{\tilde{i}}^{\text{out}})$ 

```

result in saving read from \mathbf{F}^{in} and atomic write to \mathbf{F}^{out} , respectively. However, both of these savings are less effective, as: 1) the amount to load from \mathbf{F}^{in} or atomic write to \mathbf{F}^{out} is often orders of magnitude smaller than the read from \mathbf{W} ; 2) the expectation number of unique i and j in each group is $1 + \frac{L \times N_{\text{out}}}{|\mathcal{T}|}$ and $1 + \frac{L \times N_{\text{in}}}{|\mathcal{T}|}$, respectively, — not as ideal as the $1 + \frac{L \times K}{|\mathcal{T}|}$ in the sorting by k case, as $K \ll N_{\text{in}} \simeq N_{\text{out}}$. Our analysis coincides with our profiling, thus we choose to sort by k . Nevertheless, the above analysis is based on typical configurations, in terms of unknown configurations, an auto tuning mechanism could be used to select the index to be sorted by, before the massive executions.

While the strategy of grouped processing on the triplet list sorted by k is effective in addressing the critical data locality issue, a naive parallelism strategy by letting each thread handle the computation of one group is impractical as there is often not enough on-chip memory to afford an entire $C_{\text{out}} \times C_{\text{in}}$ matrix for each thread. We discern that the MVM sub-problem in MVMR is an extreme case of an already specialized Tall-and-Skinny Matrix-Matrix Multiplication problem, by viewing a vector as a one column matrix, for which an efficient GPU algorithm is proposed in [41]. We refer the readers to [41] for the design considerations and details that could be easily adopted into our MVM sub-problem. Heavily based on this, a practical and coalesced memory access friendly parallelism strategy that integrates the aforementioned grouping logic is proposed. This strategy divides each of the $C_{\text{out}} \times C_{\text{in}}$ kernel matrix into $B_{\text{out}} \times B_{\text{in}}$ blocks, respectively, and assigns the computation related to each block of a group's MVM computation to a warp of threads, rather than assigning the entire computation to a single thread.

We summarize our algorithm for efficient computation of MVMR in Algorithm 1, where $\mathcal{T}^L = \{(i, j, k)\}$ are the L triplets assigned to a warp of threads, and on-chip resources are denoted as $\tilde{W} \in \mathbb{R}^{B_{\text{out}} \times B_{\text{in}}}$, $\tilde{F}^{\text{in}} \in \mathbb{R}^{B_{\text{in}}}$, and $\tilde{F}^{\text{out}} \in \mathbb{R}^{B_{\text{out}}}$. Note that effective global memory access saving are sup-

ported for all the three sorting options in Algorithm 1. An implementation based on Triton [45] is provided in supplementary materials. There are three hyperparameters in our algorithm, L , B_{out} and B_{in} , we use 128, 32 and 32, respectively, in all of our experiments.

Note that there is zero extra memory usage, as all the computation are native on the input tensors, without resorting to any intermediate global memory — an achieved goal similar to that in cuDNN [6] on the implementation of efficient convolution for dense and regular data. The one and only “extra” computation in our algorithm is the sorting of the triplets by k , which is a highly optimized parallel algorithm on GPU, the latency of which is negligible.

3.3.3. Efficient Gradient Computation

An efficient backward pass is critical for performant model training and finetuning. By applying the chain rule to our generalized convolution as described in Eqn. (1), two required gradients in distinct computational patterns are revealed, both of which demand efficient computation.

First, the gradient with respect to the input feature \mathbf{F}_j^{in} is an accumulation of transformed output gradients:

$$\nabla_{\mathbf{F}_j^{\text{in}}} \mathcal{L} = \sum_{(i,j,k) \in \mathcal{T}} \mathbf{W}_k^T \times \nabla_{\mathbf{F}_i^{\text{out}}} \mathcal{L}. \quad (4)$$

This operation mirrors the structure of the forward pass, where a set of matrices (now the transposed weights, \mathbf{W}_k^T) are multiplied by a set of vectors (the incoming gradients, $\nabla_{\mathbf{F}_i^{\text{out}}} \mathcal{L}$) and reduced. Consequently, this computation can be framed as an MVMR problem, allowing us to leverage the very same highly-optimized kernel designed for the forward pass, ensuring high efficiency.

Second, computing the gradient with respect to a weight kernel \mathbf{W}_k involves summing the outer products of the upstream output gradients and the corresponding input feature vectors over all associated triplets:

$$\nabla_{\mathbf{W}_k} \mathcal{L} = \sum_{(i,j,k) \in \mathcal{T}} \nabla_{\mathbf{F}_i^{\text{out}}} \mathcal{L} \otimes \mathbf{F}_j^{\text{in}}. \quad (5)$$

This structure follows a different pattern, which we abstract as Vector-Vector Outer product and Reduction (VVOR). For this, we develop a second dedicated GPU kernel, following similar principles used to develop the kernel for MVMR, which efficiently computes these outer products and reduces them into the final weight gradients. The algorithm details, as well as its Triton based implementation, of VVOR, are provided in the supplementary materials.

By implementing the backward pass with these two specialized, performant kernels—reusing MVMR for feature gradients and introducing VVOR for weight gradients—we ensure that the entire training process is computationally efficient, completing our holistic system design.

4. Experiments

Our experiments are designed to validate that PointCNN++ delivers both high performance and high fidelity. We first

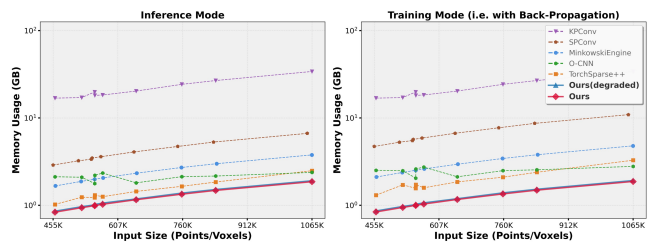


Figure 3. Memory usage comparison of one convolution layer.

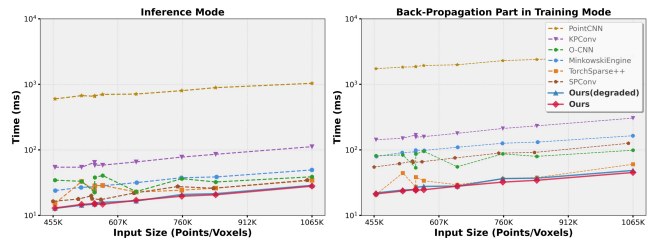


Figure 4. Performance comparison of one convolution layer.

detail our implementation in Sec. 4.1, then use micro-benchmarks to quantify its GPU memory usage and computational efficiency with extensive comparisons against the performance prioritized voxel-based methods in Sec. 4.2. Finally, we demonstrate its superior accuracy and generalization through two downstream tasks: the geometrically sensitive task of point cloud registration in Sec. 4.3 and the large-scale scene understanding task of semantic segmentation in Sec. 4.4.

4.1. Implementation Details

We represent point clouds as jagged tensors same as that in [49], and also highlight specific architectural refinements. While not the core novelty of our work, these choices diverge from common practices and offer notable benefits for quality and robustness, which we hope will prove valuable to the community. We opt for a fixed radius search over a fixed-number (K-Nearest Neighbors, or KNN) search, as its spatially-local receptive field is better suited for convolutional learning, whereas KNN is often a choice imposed by architectural limitations. For sampling, we employ a voxel-based downsampling that, by not snapping points to voxel centers, better preserves thin structures and sparsely captured regions than random downsampling with negligible latency (contradicting the claim in [18] that random sampling is essential for efficiency), and we use the original pre-downsampled points for efficient upsampling. While these operations can be lowered into spatial lookups as in Open3D [57] with the underlying ASH [12] engine based on hash map, our empirical findings led us to implement a more robust and competitively fast mechanism built upon two highly-optimized GPU primitives: sorting and searching. The implementation of these operations will be open sourced alongside our core convolution operator.

4.2. Performance Study

This section is dedicated to analyzing the computational performance of our method. We establish its efficiency and scal-

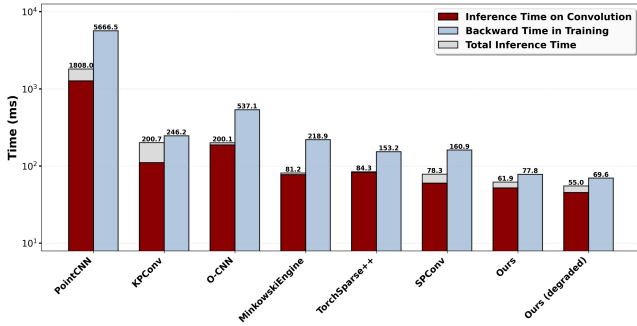


Figure 5. Performance of various convolutional backbones for 3D learning assembled in an identical ResNet-18 [17] architecture.

Table 1. Quantitative comparison for point cloud registration on the KITTI. The best and second-best results are in red and blue. Architecture abbreviations: MkEngine (MinkowskiEngine), KP+Attn (KPConv with Attention), KP+GCN (KPConv with Graph Convolutional Network), and PNpp+Attn (PointNet++ with Attention), with task-specific designs of each method list in *italic*.

Method	Arch.	RTE(m) ↓		RRE(°) ↓		Recall(%) ↑ @0.2m, 1°	Param
		Mean	Std	Mean	Std		
FCGF (2019)	MkEngine <i>ResUNet</i>	0.36	0.133	0.110	0.41	85.2	8.75M
DGR (2020)	MkEngine <i>IR prediction</i>	0.35	0.072	0.090	0.36	92.1	244.68M
CoFiNet (2021)	KP+Attn <i>Course2Fine</i>	0.39	0.065	0.091	0.39	89.4	5.48M
Predator (2021)	KP+GCN <i>Overlap pred.</i>	0.35	0.063	0.081	0.27	93.9	158.42M
GeoTrans (2022)	KP+Attn <i>Geom. emb.</i>	0.38	0.068	0.096	0.29	85.6	25.50M
Regformer (2023)	PNpp+Attn <i>Proj-aware</i>	0.22	0.077	0.058	0.28	94.6	3.12M
UMEReg (2024)	MkEngine <i>UME + SEM</i>	0.49	0.023	0.490	0.21	79.6	7.17M
Ours (2025)	PointCNN++ <i>ResUNet</i>	0.19	0.03	0.060	0.10	99.8	8.75M

ability by benchmarking its memory footprint and latency against other representative methods on the isolated convolution operator level, as well as the timing for the end-to-end forward and backward passes³.

Experimental Setup. To ensure a fair and controlled comparison, all backbones used in this study are built upon a ResNet-18 [17] architecture. For the operator-level benchmarks, we use a standard convolution configuration ($C_{in} = 64, C_{out} = 128, K = 3^3$). We benchmark our method against several representative point-based [25, 44] and voxel-based [7, 11, 43, 46] backbones. Crucially, we also include a variant named ‘Ours (degraded)’, as described in Sec. 3.2. This allows us to directly quantify the performance difference between our approach and a traditional voxel-based paradigm within the same framework. For both latency and memory analysis, we use 10 large-scale scenes from the S3DIS dataset [4]. Timings and peak GPU memory consumption are recorded on a single NVIDIA RTX 4090 GPU. To ensure stability, all reported results are averaged over 10 independent runs per scene.

Memory Analysis. The memory usage results in Figure 3 provide compelling validation for the ‘nativeness’ of our

³See supplementary material for performance benchmarks on various GPUs and scalability limits.

method. Our operator demonstrates unparalleled efficiency, consistently consuming the least GPU memory by a significant margin than existing voxel-based methods and over an order of magnitude less than existing point-based approaches⁴. This advantage is a direct consequence of our kernel’s design, which operates with a zero extra memory footprint. It fundamentally bypasses the need for the memory-intensive padding or tensor materialization that inflates the memory consumption of competing approaches. These results serve as a powerful testament to the efficacy and superior architectural design of our native operator.

Latency Analysis. Our analysis begins with the core convolution operator (Figure 4), where our method substantially outperforms leading point-based and voxel-based alternatives. As shown in Figure 5, with a high-quality implementation of other utilities, as described in Sec. 4.1, the convolution operations dominate network latency in our method, thus the operator-level advantage translates directly to powerful end-to-end performance, which is faster than all baselines in both forward and backward passes. Surprisingly, our method, when degraded into a voxel-based variant as described in Sec. 3.2, demonstrated an additional speed improvement over the full method.⁵

4.3. Downstream Task: Point Cloud Registration

Our primary goal here is to investigate a simple yet critical question: *can solely replacing a standard voxel-based backbone with our geometrically faithful PointCNN++ counterpart yield significant performance gains?*

To this end, we deliberately avoid any task-specific tuning. We take the classic FCGF [8] architecture and simply replace its MinkowskiEngine-based backbone with our PointCNN++ implementation. All other components remain identical to the original FCGF. We then compare this straightforwardly upgraded model against a series of recent state-of-the-art methods, including DGR [9], Predator [19], CoFiNet [53], GeoTrans [40], Regformer [28], and UMEReg [16]. These methods often achieve high performance through highly-specific mechanisms and sophisticated designs.

In contrast, our approach brings only a better engine to a vintage chassis, yet, as we will show, proves to be remarkably competitive. Exploring the fusion of our backbone with advanced modules is left for future work.

We conduct evaluations on two standard benchmarks that represent distinct environments. KITTI Odometry [13] is a widely-used outdoor dataset from an autonomous driving platform, featuring sparse LiDAR scans. 3DMatch [54] is a

⁴The implementation of PointCNN [25] is based on TensorFlow, while all the rest are based on PyTorch. We exclude the memory usage comparison with it, as the differences on the memory allocation strategy of these two frameworks prevent a fair direct comparison.

⁵We did not expect this, as we anticipated statistically equivalent performance. While the definitive cause for this speed-up remains unknown, our analysis showed that the variant’s kernel matrix usage is more concentrated on spatially ‘central’ kernels. We conjecture that this access pattern could lead to improved cache performance, though this requires further investigation.

Table 2. Quantitative comparison on 3DMatch across different input point densities. We report Registration Recall (RR \uparrow), Feature Matching Recall (FMR \uparrow), and Inlier Ratio (IR \uparrow), all in percentages (%). The best results are in red bold, and the second best in blue bold.

Method	5000 Pts			2500 Pts			1000 Pts			500 Pts			250 Pts		
	RR	FMR	IR	RR	FMR	IR	RR	FMR	IR	RR	FMR	IR	RR	FMR	IR
FCGF(2019)	85.1	97.4	52.8	84.7	97.3	51.1	83.3	97.0	46.7	81.6	96.7	41.5	71.4	96.6	34.1
CoFiNet(2021)	89.3	98.1	49.8	88.9	98.3	51.2	88.4	98.1	51.9	87.4	98.2	52.2	87.0	98.3	52.2
Predator(2021)	89.0	96.6	58.0	89.9	96.6	58.4	90.6	96.5	57.1	88.5	96.3	54.1	86.6	96.5	49.3
GeoTrans(2023)	91.4	97.9	70.5	91.1	98.0	72.9	92.0	98.1	75.2	91.7	98.2	79.8	91.2	98.3	84.6
Ours	90.3	99.3	58.2	90.2	99.3	61.4	89.2	99.3	62.5	89.1	99.3	63.4	88.3	99.1	64.1

large-scale indoor dataset composed of RGB-D scans from various scenes. Following UMEReg [16], we report Relative Translation Error (RTE), Relative Rotation Error (RRE) and Registration Recall (RR) on the KITTI. On the 3DMatch, we follow the protocol of CoFiNet [53] and use Registration Recall (RR), Feature Matching Recall (FMR), and Inlier Ratio (IR) for evaluation.

Results on KITTI. As shown in Table 1, the results on the KITTI dataset powerfully validate the effectiveness of our method in the zero-task specific tuning “plug-and-play” setting. By simply replacing the backbone in the classic FCGF architecture, we achieve state-of-the-art performance, nearly halving its RTE to 0.19m and boosting its RR to a near-perfect 99.8%. Beyond raw accuracy, our model demonstrates unparalleled registration stability, achieving the lowest standard deviations in both translation and rotation by a wide margin. This combination of top-tier accuracy and exceptional consistency proves that our high-fidelity operator produces superior features, leading to quantifiably better and more reliable registrations.

Results on 3DMatch. As shown in Table 2, our method demonstrates compelling results on the 3DMatch benchmark. This is particularly noteworthy as these results are achieved using an old backbone, which is considerably simpler than the bespoke architectures of competitors like GeoTrans [40]. It indicates that our operator is powerful enough on its own to elevate a conventional architecture to produce features of state-of-the-art quality.

4.4. Downstream Task: Semantic Segmentation

To evaluate the generalization and versatility of PointCNN++, we further benchmark it on semantic segmentation, a task that requires capturing both fine-grained local geometry and large-scale contextual information.

Experimental Setup. We conduct a head-to-head comparison on the challenging nuScenes [5] semantic segmentation benchmark. Following the methodology in Section 4.3, we adopt a “drop-in” replacement strategy to ensure fairness. We use a ResUNet2 backbone (with MSC pre-training and a batch size of 3) as the reference architecture. The core convolution operators of the backbone are replaced while keeping the rest of the architecture, including the loss functions and training pipeline, identical.

Baselines. We compare our operator against representative point-based and voxel-based paradigms. Specifically, we include KPConv [44] as the point-based baseline and MinkowskiEngine [7] as the voxel-based baseline. This

Table 3. Head-to-head comparison on nuScenes semantic segmentation using a ResUNet2 backbone.

Operator	mIoU (%) \uparrow	mAcc (%) \uparrow	Mem (GB) \downarrow	Time (per iter.) \downarrow
KpConv (point-based)	-	-	17.03	23.82
MinkowskiEngine (voxel-based)	74.4	81.7	2.61	0.131
PointCNN++ (Ours)	78.2	85.3	2.43	0.102

setup allows us to directly measure the efficiency and accuracy gains provided by the PointCNN++ operator within a consistent network framework.

Results and Analysis. The results are summarized in Table 3. Under the same experimental setting, training KPConv would require approximately 116 days, rendering it impractical for real-world deployment. We therefore only report its memory and per-iteration time measured during the training phase. PointCNN++ achieves the highest performance with an mIoU of 78.2% and an mAcc of 85.3%, outperforming the voxel-based MinkowskiEngine by 3.8% in mIoU. Crucially, PointCNN++ delivers these improvements with a smaller memory footprint (2.43 GB vs. 2.61 GB) and faster execution (0.102s vs. 0.131s per iteration) than the voxel-based baseline. Compared to the point-based KPConv, PointCNN++ reduces memory usage by nearly an order of magnitude and is over 200 \times faster per iteration. These results underscore the dual advantage of our method: it preserves the geometric fidelity typically associated with point-based methods while surpassing the computational efficiency of optimized voxel-based frameworks.

5. Conclusion

We present PointCNN++, a solution to the persistent trade-off in 3D deep learning between the performance of voxel-based grids and the precision of point-based operations. Our work dismantles this compromise by introducing a computational system designed specifically for irregular point clouds. By framing convolution on this data as an MVMR problem, we developed a dedicated GPU kernel that executes natively without the inefficient overheads. Our results confirm that sub-voxel accuracy could be retained without the performance penalties of prior methods, demonstrating that computational efficiency and geometric fidelity can be achieved in unison, enabling powerful, faithful geometric learning.

Acknowledgments

This work was supported by Ant Group Research Intern Program. And we express our sincere gratitude to Jiaheng Li from Peking University for his assistance in setting up the Ant Group infrastructure used in the comparison experiments.

References

- [1] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. TensorFlow: a system for Large-Scale machine learning. In *12th USENIX symposium on operating systems design and implementation (OSDI 16)*, 2016.
- [2] Advanced Micro Devices (AMD). AMD CDNA 2 Architecture. Whitepaper, 2021. Available at: <https://www.amd.com/content/dam/amd/en/documents/instinct-business-docs/white-papers/amd-cdna2-white-paper.pdf>.
- [3] Evangelos Alexiou, Evgeniy Upenik, and Touradj Ebrahimi. Towards subjective quality assessment of point cloud imaging in augmented reality. In *2017 IEEE 19th international workshop on multimedia signal processing (MMSP)*, 2017.
- [4] Iro Armeni, Ozan Sener, Amir R. Zamir, Helen Jiang, Ioannis Brilakis, Martin Fischer, and Silvio Savarese. 3d semantic parsing of large-scale indoor spaces. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [5] Holger Caesar, Varun Bankiti, Alex H. Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multi-modal dataset for autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11621–11631, 2020.
- [6] Sharan Chetlur, Cliff Woolley, Philippe Vandermersch, Jonathan Cohen, John Tran, Bryan Catanzaro, and Evan Shelhamer. cudnn: Efficient primitives for deep learning. *arXiv preprint arXiv:1410.0759*, 2014.
- [7] Christopher Choy, JunYoung Gwak, and Silvio Savarese. 4d spatio-temporal convnets: Minkowski convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3075–3084, 2019.
- [8] Christopher Choy, Jaesik Park, and Vladlen Koltun. Fully convolutional geometric features. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 8958–8966, 2019.
- [9] Christopher Choy, Wei Dong, and Vladlen Koltun. Deep global registration. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020.
- [10] Christopher Choy, Junha Lee, Rene Ranftl, Jaesik Park, and Vladlen Koltun. High-dimensional convolutional networks for geometric pattern recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020.
- [11] Spconv Contributors. Spconv: Spatially sparse convolution library. <https://github.com/traveller59/spconv>, 2022.
- [12] Wei Dong, Yixing Lao, Michael Kaess, and Vladlen Koltun. ASH: A modern framework for parallel spatial hashing in 3d perception. *IEEE transactions on pattern analysis and machine intelligence*, 2022.
- [13] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *2012 IEEE conference on computer vision and pattern recognition*, 2012.
- [14] Benjamin Graham and Laurens Van der Maaten. Sub-manifold sparse convolutional networks. *arXiv preprint arXiv:1706.01307*, 2017.
- [15] An Guo, Yang Feng, and Zhenyu Chen. Lirtest: augmenting lidar point clouds for automated testing of autonomous driving systems. In *Proceedings of the 31st ACM SIGSOFT International Symposium on Software Testing and Analysis*, 2022.
- [16] Yuval Haitman, Amit Efraim, and Joseph M Francos. Umeregrobust-universal manifold embedding compatible features for robust point cloud registration. In *European Conference on Computer Vision*, 2024.
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016.
- [18] Qingyong Hu, Bo Yang, Linhai Xie, Stefano Rosa, Yulan Guo, Zhihua Wang, Niki Trigoni, and Andrew Markham. Randla-net: Efficient semantic segmentation of large-scale point clouds. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020.
- [19] Shengyu Huang, Zan Gojcic, Mikhail Usvyatsov, Andreas Wieser, and Konrad Schindler. Predator: Registration of 3d point clouds with low overlap. In *Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition*, 2021.
- [20] Xinyu Huang, Xinjing Cheng, Qichuan Geng, Binbin Cao, Dingfu Zhou, Peng Wang, Yuanqing Lin, and Ruigang Yang. The apollo-scape dataset for autonomous driving. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2018.
- [21] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio Guadarrama, and Trevor Darrell. Caffe: Convolutional architecture for fast feature embedding. In *Proceedings of the 22nd ACM international conference on Multimedia*, 2014.
- [22] Norm Jouppi, George Kurian, Sheng Li, Peter Ma, Rahul Nagarajan, Lifeng Nai, Nishant Patil, Suvinay Subramanian, Andy Swing, Brian Towles, et al. Tpu v4: An optically reconfigurable supercomputer for machine learning with hardware support for embeddings. In *Proceedings of the 50th annual international symposium on computer architecture*, 2023.
- [23] Pileun Kim, Jingdao Chen, and Yong K Cho. Slam-driven robotic mapping and registration of 3d point clouds. *Automation in Construction*, 89:38–48, 2018.
- [24] Xin Lai, Jianhui Liu, Li Jiang, Liwei Wang, Hengshuang Zhao, Shu Liu, Xiaojuan Qi, and Jiaya Jia. Stratified transformer for 3d point cloud segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8500–8509, 2022.
- [25] Yangyan Li, Rui Bu, Mingchao Sun, Wei Wu, Xinhan Di, and Baoquan Chen. PointCNN: Convolution on \mathcal{X} -transformed points. *Advances in neural information processing systems*, 2018.
- [26] Jiong Lin, Lechen Zhang, Kwansoo Lee, Jialong Ning, Judah Goldfeder, and Hod Lipson. Autourdf: Unsupervised robot

- modeling from point cloud frames using cluster registration. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, 2025.
- [27] Yujun Lin, Zhekai Zhang, Haotian Tang, Hanrui Wang, and Song Han. Pointacc: Efficient point cloud accelerator. In *MICRO-54: 54th Annual IEEE/ACM International Symposium on Microarchitecture*, 2021.
- [28] Jiuming Liu, Guangming Wang, Zhe Liu, Chaokang Jiang, Marc Pollefeys, and Hesheng Wang. Regformer: An efficient projection-aware transformer network for large-scale point cloud registration. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023.
- [29] Zhijian Liu, Xinyu Yang, Haotian Tang, Shang Yang, and Song Han. Flatformer: Flattened window attention for efficient point cloud transformer, 2023.
- [30] Dening Lu, Qian Xie, Kyle Gao, Linlin Xu, and Jonathan Li. 3dctn: 3d convolution-transformer network for point cloud classification. *IEEE Transactions on Intelligent Transportation Systems*, pages 1–12, 2022.
- [31] John Nickolls, Ian Buck, Michael Garland, and Kevin Skadron. Scalable parallel programming with cuda: Is cuda the parallel programming model that application developers have been waiting for? *Queue*, 6(2):40–53, 2008.
- [32] NVIDIA. NVIDIA A100 Tensor Core GPU Architecture. Whitepaper, 2020. Available at: <https://www.nvidia.com/content/dam/en-zz/Solutions/Data-Center/nvidia-ampere-architecture-whitepaper.pdf>.
- [33] NVIDIA. NVIDIA H100 Tensor Core GPU Architecture. Whitepaper, 2022. Available at: <https://resources.nvidia.com/en-us-h100-whitepaper/h100-whitepaper-12-0922>.
- [34] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*, 2019.
- [35] Jette J Peek, Xucong Zhang, Klaus Hildebrandt, Samuel Alexander Max, Amir H Sadeghi, AJJC Bogers, and EAF Mahtab. A novel 3d image registration technique for augmented reality vision in minimally invasive thoracoscopic pulmonary segmentectomy. *International journal of computer assisted radiology and surgery*, 20(4):787–795, 2025.
- [36] Alessio Pierluigi Placitelli and Luigi Gallo. Low-cost augmented reality systems via 3d point cloud sensors. In *2011 Seventh International Conference on Signal Image Technology & Internet-Based Systems*, 2011.
- [37] François Pomerleau, Francis Colas, Roland Siegwart, et al. A review of point cloud registration algorithms for mobile robotics. *Foundations and Trends® in Robotics*, 4(1):1–104, 2015.
- [38] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. PointNet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017.
- [39] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. PointNet++: Deep hierarchical feature learning on point sets in a metric space. *Advances in neural information processing systems*, 2017.
- [40] Zheng Qin, Hao Yu, Changjian Wang, Yulan Guo, Yuxing Peng, Slobodan Ilic, Dewen Hu, and Kai Xu. Geotransformer: Fast and robust point cloud registration with geometric transformer. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- [41] Cody Rivera, Jieyang Chen, Nan Xiong, Jing Zhang, Shuaiwen Leon Song, and Dingwen Tao. Tsm2x: High-performance tall-and-skinny matrix–matrix multiplication on gpus. *Journal of Parallel and Distributed Computing*, 151: 70–85, 2021.
- [42] Haotian Tang, Zhijian Liu, Xiuyu Li, Yujun Lin, and Song Han. Torchsparse: Efficient point cloud inference engine. *Proceedings of Machine Learning and Systems*, 2022.
- [43] Haotian Tang, Shang Yang, Zhijian Liu, Ke Hong, Zhongming Yu, Xiuyu Li, Guohao Dai, Yu Wang, and Song Han. Torchsparse++: Efficient training and inference framework for sparse convolution on gpus. In *IEEE/ACM International Symposium on Microarchitecture (MICRO)*, 2023.
- [44] Hugues Thomas, Charles R Qi, Jean-Emmanuel Deschaud, Beatriz Marcotegui, François Goulette, and Leonidas J Guibas. KPConv: Flexible and deformable convolution for point clouds. In *Proceedings of the IEEE/CVF international conference on computer vision*, 2019.
- [45] Philippe Tillet, Hsiang-Tsung Kung, and David Cox. Triton: an intermediate language and compiler for tiled neural network computations. In *Proceedings of the 3rd ACM SIGPLAN International Workshop on Machine Learning and Programming Languages*, pages 10–19, 2019.
- [46] Peng-Shuai Wang, Yang Liu, Yu-Xiao Guo, Chun-Yu Sun, and Xin Tong. O-cnn: Octree-based convolutional neural networks for 3d shape analysis. *ACM Transactions On Graphics (TOG)*, 36(4):1–11, 2017.
- [47] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E Sarma, Michael M Bronstein, and Justin M Solomon. Dynamic graph cnn for learning on point clouds. *ACM Transactions on Graphics (tog)*, 2019.
- [48] Zeyu Wang, Cuong Nguyen, Paul Asente, and Julie Dorsey. Pointshopar: Supporting environmental design prototyping using point cloud in augmented reality. In *Proceedings of the 2023 CHI conference on human factors in computing systems*, 2023.
- [49] Francis Williams, Jiahui Huang, Jonathan Swartz, Gergely Klar, Vijay Thakkar, Matthew Cong, Xuanchi Ren, Ruilong Li, Clement Fuji-Tsang, Sanja Fidler, et al. fvd: A deep-learning framework for sparse, large scale, and high performance spatial intelligence. *ACM Transactions on Graphics (TOG)*, 2024.
- [50] Wenxuan Wu, Zhongang Qi, and Li Fuxin. Pointconv: Deep convolutional networks on 3d point clouds. In *Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition*, 2019.
- [51] Yan Yan, Yuxing Mao, and Bo Li. Second: Sparsely embedded convolutional detection. *Sensors*, 2018.
- [52] Zetong Yang, Li Chen, Yanan Sun, and Hongyang Li. Visual point cloud forecasting enables scalable autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024.
- [53] Hao Yu, Fu Li, Mahdi Saleh, Benjamin Busam, and Slobodan Ilic. Cofinet: Reliable coarse-to-fine correspondences for robust pointcloud registration. *Advances in Neural Information Processing Systems*, 2021.
- [54] Andy Zeng, Shuran Song, Matthias Nießner, Matthew Fisher, Jianxiong Xiao, and Thomas Funkhouser. 3dmatch: Learning local geometric descriptors from rgb-d reconstructions. In

Proceedings of the IEEE conference on computer vision and pattern recognition, 2017.

- [55] Cheng Zhang, Haocheng Wan, Shengqiang Liu, Xinyi Shen, and Zizhao Wu. Pvt: Point-voxel transformer for 3d deep learning. *arXiv preprint arXiv:2108.06076*, 2021.
- [56] Hengshuang Zhao, Li Jiang, Jiaya Jia, Philip HS Torr, and Vladlen Koltun. Point transformer. In *Proceedings of the IEEE/CVF international conference on computer vision*, 2021.
- [57] Qian-Yi Zhou, Jaesik Park, and Vladlen Koltun. Open3D: A modern library for 3d data processing. *arXiv preprint arXiv:1801.09847*, 2018.
- [58] Yin Zhou and Oncel Tuzel. Voxelnet: End-to-end learning for point cloud based 3d object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018.
- [59] Zixiang Zhou, Xiangchen Zhao, Yu Wang, Panqu Wang, and Hassan Foroosh. Centerformer: Center-based transformer for 3d object detection. In *ECCV*, 2022.
- [60] Haoyi Zhu, Yating Wang, Di Huang, Weicai Ye, Wanli Ouyang, and Tong He. Point cloud matters: Rethinking the impact of different observation spaces on robot learning. *Advances in Neural Information Processing Systems*, 2024.