

Sparse PointPillars: Maintaining and Exploiting Input Sparsity to Improve Runtime on Embedded Systems

Kyle Vedder* and Eric Eaton†

Abstract—Bird’s Eye View (BEV) is a popular representation for processing 3D point clouds, and by its nature is fundamentally sparse. Motivated by the computational limitations of mobile robot platforms, we create a fast, high-performance BEV 3D object detector that maintains *and* exploits this input sparsity to decrease runtimes over non-sparse baselines and avoids the tradeoff between pseudoimage area and runtime. We present results on KITTI, a canonical 3D detection dataset, and Matterport-Chair, a novel Matterport3D-derived chair detection dataset from scenes in real furnished homes. We evaluate runtime characteristics using a desktop GPU, an embedded ML accelerator, and a robot CPU, demonstrating that our method results in significant detection speedups ($2\times$ or more) for embedded systems with only a modest decrease in detection quality. Our work represents a new approach for practitioners to optimize models for embedded systems by maintaining *and* exploiting input sparsity throughout their entire pipeline to reduce runtime and resource usage while preserving detection performance. All models, weights, experimental configurations, and datasets used are publicly available¹.

I. INTRODUCTION

In modern robot perception systems, fitting state-of-the-art machine learning models within the power and compute constraints of on-board computational platforms is a major challenge. In the autonomous vehicle space, high-end desktop GPUs and CPUs are often available on-board, but this hardware still faces power and cost limits and must be shared with other components of the control stack. This challenge is even more pronounced for intelligent mobile robots; for example, even a larger, high-end robot like the Fetch Freight [1] is not able to power several desktop-grade GPUs and high-end CPUs in order to run its control stack. Instead, roboticists are forced to settle for smaller embedded systems like NVidia’s Jetson [2], while smaller platforms such as quadcopters often struggle to handle the weight or power requirements of even these embedded systems [3]. This motivates the problem of developing machine learning models that have significantly reduced resource usage compared to existing models while preserving their performance — models need to be shrunk not just to fit on smaller devices, but to fit while *sharing* these resources with other components.

In this paper we address this problem of reducing resource usage while preserving performance within PointPillars [4], a point cloud-based 3D object detector that is *very* popular

among autonomous vehicle makers due to its speed and performance. PointPillars segregates the raw points into pillars and vectorizes these point collections into N dimensional vectors, resulting in a sparse Bird’s Eye View (BEV) pseudoimage of the scene. PointPillars then processes the pseudoimage with a dense 2D convolutional Backbone and predicts bounding boxes using Single-Stage Detector [5]. The PointPillars Backbone is the most expensive component of the pipeline, yet the pseudoimage it processes is highly sparse. As a result, large, empty sections of the image are unnecessarily convolved by the Backbone and runtimes directly scale with the pseudoimage area, presenting a trade-off.

Our method, Sparse PointPillars, avoids these inefficiencies and tradeoffs with a novel processing pipeline that maintains *and* exploits end-to-end sparsity to reduce runtime and resource usage on embedded systems while maintaining reasonable performance. Our main contributions include:

- 1) A **new pipeline** that **maintains and exploits representational sparsity** and **avoids the trade-off between pseudoimage area and runtime**.
- 2) **$2\times$ or more speedup on embedded systems**, allowing practitioners to gain large runtime speedups for the **same power budget** or **modest runtime speedups for a significantly smaller power budget**, all in exchange for a modest decrease in detection quality.
- 3) A general **design approach centered around representational sparsity** for efficient embedded system pipelines.

An overview of the original PointPillars’ pipeline and our new pipeline (Sparse PointPillars) is shown in Fig. 1.

II. RELATED WORK

Solutions to the problem of reducing resource usage of machine learning models in order to deploy on embedded hardware typically fall into three categories. One general-purpose solution to this problem, model quantization [6]–[10], first trains models in a standard fashion using floating point weights and then, after training, converts some [6] or all [7], [8] weights into integer [9] or binary [10] quantized values that are faster to multiply than floating point values. The quantized network is then finetuned, resulting in similar performance while running faster on accelerators (e.g. GPUs [11]), low-end compute devices (e.g. mobile phones [12]), or specialized hardware (e.g. FPGAs [13]).

A second model-agnostic approach, often paired with quantization, is model *weight* sparsification, also known as model pruning. As per the Lottery Ticket Hypothesis [14], most weights contribute little to performance [15]–[17] and can be pruned to improve runtime [18]–[24], either with

*Kyle Vedder is with the Dept. of Computer and Information Science, University of Pennsylvania, Philadelphia, PA 19104, USA kvedder@seas.upenn.edu

†Eric Eaton is with the Dept. of Computer and Information Science, University of Pennsylvania, Philadelphia, PA 19104, USA eeaton@seas.upenn.edu

¹https://vedder.io/sparse_point_pillars

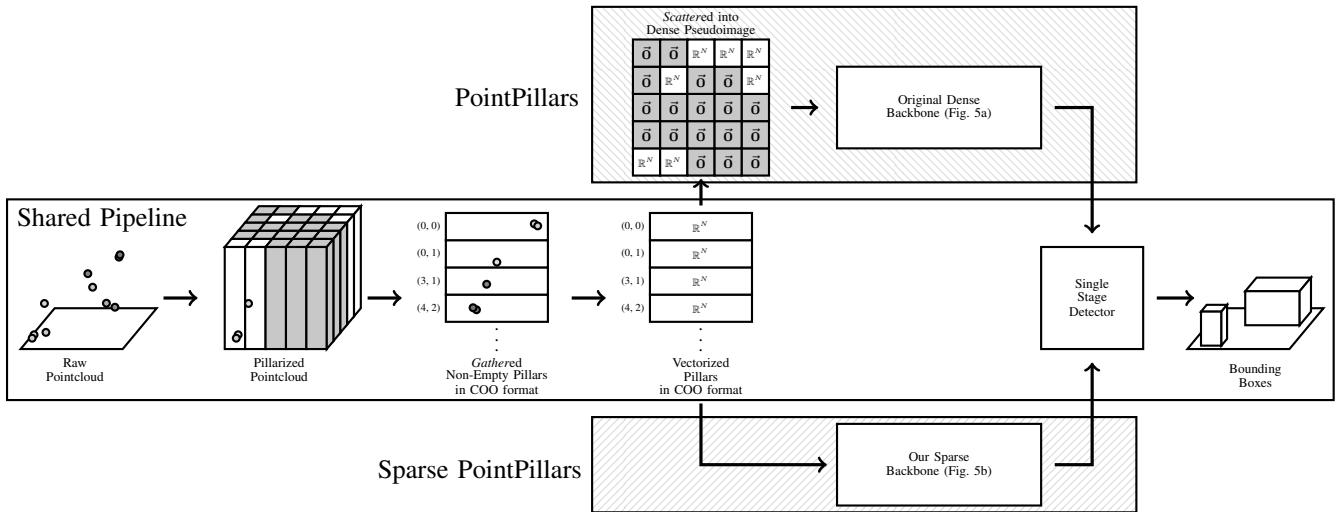


Fig. 1: Overview of PointPillars’ [4] pipeline and Sparse PointPillars’ pipeline. Both use a shared pipeline which consumes a raw point cloud, pillarizes the scene (empty pillars are depicted in gray), *gathers* the non-empty pillars into a sparse COO matrix, and vectorizes the pillars. PointPillars’ Pillar Feature Net then *scatters* this sparse matrix into a dense matrix, passing this to its dense Backbone for processing. Sparse PointPillars’ Pillar Feature Net preserves the sparse COO representation, passing this to our sparse Backbone. Both PointPillars and Sparse PointPillars use Single Stage Detector [5] to regress output bounding boxes. PointPillars spends the majority of its runtime running its dense Backbone; Sparse PointPillars reduces this runtime by maintaining and exploiting the natural sparsity of the COO matrix via its new Pillar Feature Net and Backbone.

regular structure to exploit hardware properties [19]–[21], or without structure [22]–[24], without much performance loss.

A third approach is to exploit data representations to reduce the computational burden, such as input representational sparsity. For example, a common approach for 3D object detection in point clouds is to voxelize the 3D space and perform 3D convolutions through a pipeline similar to 2D object detection [25]–[27]. 3D convolution of these voxels dominates network runtime [25], but the voxels tend to be highly sparse; point clouds from KITTI contain over 100,000 points, but when voxelized into 16cm cubes, over 95% are empty [4], [26]. This motivates sparsity-aware 3D convolutional methods such as SECOND [26], which performs sparse 3D convolutions to produce mathematically identical results with significantly less computation and runtime. SBNet [28], a BEV object detector, takes a spiritually similar approach by doing dense 2D convolutions only in coarsely masked regions of relevance to reduce the area convolved, decreasing runtime without a significant performance hit.

Our method, Sparse PointPillars, falls into this third category, using sparse matrix representations and multiple types of specialized 2D sparse convolutions to decrease runtime without a significant impact on performance. Importantly, our input sparsity-aware contributions are *orthogonal* to model quantization and weight sparsification, thus providing practitioners another useful tool in improving model runtime.

III. MAINTAINING AND EXPLOITING END-TO-END SPARSITY WITHIN SPARSE POINTPILLARS

This section describes Sparse PointPillars’s new processing pipeline. We then theoretically analyze (Section IV) and empirically validate (Section V) our modifications.

A. New Pillar Feature Net

PointPillars performs pillarization (Fig. 1) via its Pillar Feature Net which *gathers* the non-zero entries of the full scene into a coordinate (COO) format sparse tensor [29] with a fixed number of pillars and a fixed number of points per pillar (set a priori), along with the coordinate location of each pillar. The Feature Net then vectorizes each pillar using a PointNet-like vectorizer [30]. In the original PointPillars’ Feature Net, these resulting vectors are then *scattered* back into a standard dense tensor in the shape of the full scene.

In our modified Pillar Feature Net, we elide this *scatter* step, thus maintaining the sparse COO matrix format of the vectorized pillars. This significantly reduces GPU requirements by avoiding an additional allocation of a $\text{scene_width} \times \text{scene_height} \times N$ matrix plus complex matrix masking to insert the appropriate values, *and* allows the Pillar Feature Net to emit a sparse pseudoimage output, enabling its further exploitation by our new Backbone.

B. New Backbone

The original PointPillars Backbone (Fig. 5a) takes in the dense tensor format pseudoimage from the *scatter* step and processes it with a convolutional feature pyramid network [31] style Backbone. This Backbone emits a single large pseudoimage of half the width and height of the input pseudoimage, composed of intermediary pseudoimages from the three layers of the Backbone concatenated along the channel axis. Due to the heavy use of standard 3×3 stride-1 convolutions throughout the Backbone, there is significant smearing of non-zero entries across the pseudoimage; Fig. 2 visualizes this for an input from the KITTI dataset: Fig. 2a shows the input pseudoimage, and Figs. 2b–2d show the smearing of

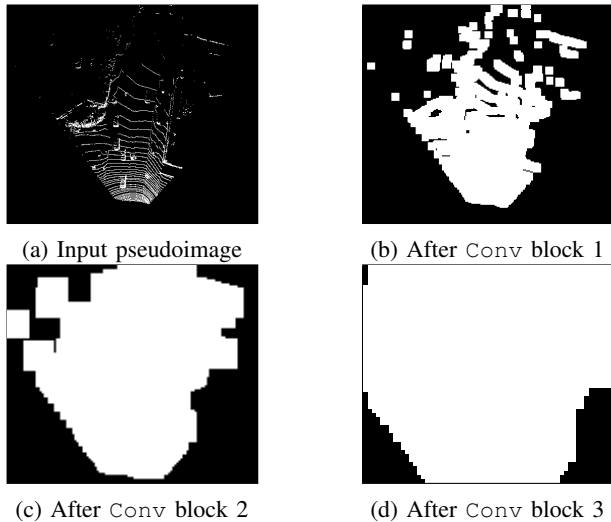


Fig. 2: Pseudoimages from original PointPillars with BatchNorm removed for visualization; black represents zero entries on all channels and white represents at least one non-zero channel entry. With BatchNorm retained, sparsity is entirely destroyed as zero entries are modified.

non-zero entries across zero entries in the pseudoimage as it travels through the Backbone.

Due to this smearing, naïvely modifying the original PointPillars Backbone to perform sparse convolutions [26] on the sparse COO format pseudoimage from our new Pillar Feature Net would fail to *preserve* input sparsity. A sparse convolution (Fig. 4b) is mathematically equal to a standard convolution (zero entries are simply skipped to save computation), and thus later pseudoimages would require convolution of almost all entries, e.g. Fig. 2d, hurting runtime.

Thus, Sparse PointPillars’ new Backbone (Fig. 5b) operates on the sparse COO format pseudoimage with a combination of carefully placed sparse convolutions and submanifold (SubM) convolutions [32], a type of convolution that only convolves filters centered on existing non-zero entries in order to maintain *and* exploit the pseudoimage’s sparsity throughout the Backbone. Fig. 4 shows the results of these different types of convolutions, and Fig. 3 visually demonstrates their impact on the pseudoimage’s sparsity as it flows through the Sparse PointPillars Backbone. Compared to PointPillars’ Backbone, Sparse PointPillars’ Backbone maintains and exploits sparsity in three key ways: 1) its BatchNorm is only applied to the recorded non-zero entries, leaving the zero entries unaffected (this is trivial for COO tensors as only non-zero values are recorded), 2) it replaces the 3×3 stride-2 convolutions with 2×2 stride-2 convolutions when shrinking the pseudoimage, ensuring that non-zero entries in the higher resolution pseudoimage appear only once in the lower resolution pseudoimage, and 3) it replaces the 3×3 stride-1 convolutions with 3×3 stride-1 SubM convolutions. As we will show, these changes require the new Backbone to perform fewer computations (Section IV) resulting in significantly improved runtime in practice (Section V).

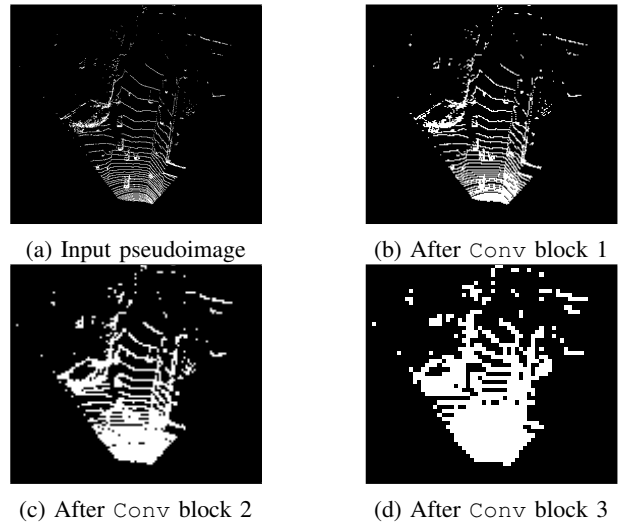


Fig. 3: Pseudoimages from our Sparse PointPillars run on a sample from KITTI. Black represents zero entries on all channels and white represents at least one non-zero channel entry. Due to the use of SubM convs and BatchNorm only operating over non-zero entries, sparsity is maintained.

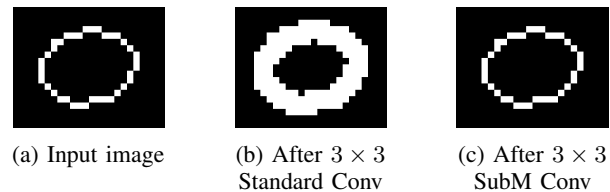


Fig. 4: 3×3 stride-1 Standard Convolution versus 3×3 Submanifold (SubM) Convolution. Black represents zero entries on all channels and white represents at least one non-zero channel entry. Standard convolutions can be centered on zero entries next to non-zero entries, resulting in a new non-zero entry, causing smearing and destroying sparsity. SubM convolutions are only centered on non-zero entries, preventing smearing and maintaining sparsity.

IV. THEORETICAL ANALYSIS

Our new Backbone maintains and exploits input sparsity to perform fewer operations, reduce memory usage, and avoid the pseudoimage area vs runtime tradeoff in order to achieve faster runtime compared to PointPillars’ Backbone. For PointPillars, the memory consumed and convolution count of its Backbone (Fig. 5a) is a function of the area and number of channels of the input pseudoimage; the *values* of the input pseudoimage are irrelevant. By comparison, the memory consumed and convolution count of Sparse PointPillars’ Backbone (Fig. 5b) is a function of the area, number of channels, *and* pseudoimage density (the fraction of non-zero values of the input pseudoimage). Due to the strategic use of sparse 2×2 and SubM stride-2 convolutions in each Conv block, Sparse PointPillars’ Backbone increases the pseudoimage density by *at most* $4 \times$ per block; however, in practice this density increase is far below $4 \times$ as non-zero entries tend to cluster near one another (e.g. Fig. 3).

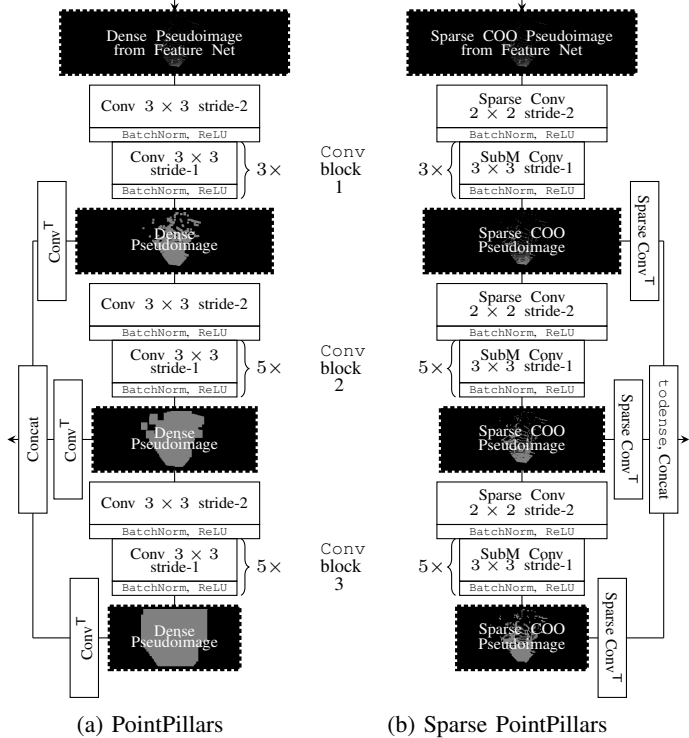


Fig. 5: PointPillars vs Sparse PointPillars Backbone. The Sparse PointPillars Backbone maintains and exploits pseudoimage sparsity by using SubM convs and 2×2 stride-2 convs to avoid the smearing effect of 3×3 stride-1 convs.

Table I outlines the type and number of convolutions performed by each Backbone, providing an exact count for PointPillars and an *upper bound* for Sparse PointPillars. When the input pseudoimage density D approaches 1, Sparse PointPillars’ Backbone converges to PointPillars’ Backbone with the 3×3 stride-2 convolutions replaced with 2×2 stride-2 convolutions; when D approaches 0, Sparse PointPillars’ Backbone performs significantly fewer convolutions.

In practice, input density D is very small for realistic data; for KITTI’s test data, the median D is 0.02459 (min 0.013321, max 0.03899) and for Matterport-Chair’s test data, the median D is 0.00750 (min 0.00029, max 0.01679). Translating this to convolution operations count, for KITTI’s median density, our Backbone performs *at least 50%* fewer convolutions, and for Matterport-Chair’s median density, our Backbone performs *at least 79%* fewer convolutions.

As Sparse PointPillars’ runtime is dependent on input density D and PointPillars’ is not, Sparse PointPillars avoids the tradeoff between pseudoimage area and runtime, providing practitioners more flexibility in tuning their detector without significant runtime increases. For example, practitioners want to decrease the pillar size to allow for the capture of finer-grained objects or increase in the maximum sensor range to detect of objects further away (both of which decrease D as much as increase $W \times H$), and Sparse PointPillars enables both of these without the same runtime hit as PointPillars.

TABLE I: Number of convolutions performed by PointPillars’ Backbones and an *upper bound* on number of convolutions performed by Sparse PointPillars’ Backbone, for an input pseudoimage of size $W \times H$ with C channels and D density.

Operation	Baseline Count	Sparse PointPillars’ Upper Bound
3×3 Conv	$\frac{15}{4} C^2 HW$	$C^2 HW (\min(\frac{3}{4}, 3D) + \min(\frac{5}{4}, 20D) + \min(\frac{5}{4}, 80D))$
2×2 Conv	0	$C^2 HW (\frac{D}{4} + \min(\frac{1}{8}, \frac{1}{2}D) + \min(\frac{1}{8}, 2D))$
1×1 Conv ^T	$\frac{1}{2} C^2 HW$	$C^2 HW \min(\frac{1}{2}, 2D)$
2×2 Conv ^T	$\frac{1}{4} C^2 HW$	$C^2 HW \min(\frac{1}{4}, 4D)$
4×4 Conv ^T	$\frac{1}{8} C^2 HW$	$C^2 HW \min(\frac{1}{8}, 8D)$

V. EMPIRICAL EVALUATION

To validate the design of Sparse PointPillars, we implemented it in Open3D-ML [33], a high-quality third party implementation of PointPillars, using the Minkowski Engine [34] for sparse convolutions. To demonstrate Sparse PointPillars’ value on embedded systems via a realistic task, we compare it against PointPillars on *Matterport-Chair*, a custom chair detection task derived from Matterport3D [35], an indoor 3D scan dataset designed to simulate a task required of real service robots (Section V-A). We evaluate the two trained Matterport-Chair models across three compute platforms: a desktop with a high-end GPU, an embedded ML accelerator configured for minimal and maximal power modes, and the CPU of a high-end commercial robot.

To contextualize Sparse PointPillars’ performance, we evaluate it against PointPillars on the KITTI [36] dataset (Section V-B), a well understood 3D self-driving dataset used for evaluation in the original PointPillars paper, using a high-end GPU desktop setup common in these evaluations, and compare these results to other baseline models. We also perform several ablative studies on our Backbone to demonstrate that it produces a reasonable trade-off between runtime and detection performance.

A. Matterport-Chair Evaluation with Embedded Performance

To simulate a realistic detection task faced by a service robot or other embodied platform using embedded compute systems, we constructed a labeled chair detection dataset *Matterport-Chair* using point clouds and their object labels sampled from houses in Matterport3D [35]. Matterport3D is a dataset of multiple building-scale indoor 3D meshes constructed using many high-resolution panoramic RGBD views taken inside real houses and labeled with 3D bounding boxes and semantic labels for over 20 different object classes. To generate our training and test dataset, we sampled point clouds of random views from the perspective of a robot sitting one meter off the ground across four different high quality house meshes, producing a train/test split of 7,500 point clouds each (the same size as the KITTI splits). We post-processed the bounding boxes, aligning them vertically and rejecting boxes that were highly occluded, associated with too few points, or caused by dataset noise such as holes

TABLE II: Per instance model component runtime and standard deviation in milliseconds, run on Matterport-Chair’s test set, averaged over ten trials for Desktop and Xavier High and averaged over three trials for Xavier Low and Robot. Models differ only in their Feature Net and Backbone; all other components are identical. Lower is better.

	To Device	Feature Extract	Feature Net	Backbone	Head	BBox Extract	Total time
Desktop Dense	0.072 ± 0.001	1.783 ± 0.010	0.390 ± 0.005	2.512 ± 0.017	0.199 ± 0.001	15.198 ± 0.145	20.154 ± 0.143
Desktop Sparse	0.073 ± 0.001	1.747 ± 0.015	0.130 ± 0.002	6.633 ± 0.038	0.218 ± 0.002	4.854 ± 0.010	13.655 ± 0.060
Xavier High Dense	0.657 ± 0.056	13.526 ± 0.195	4.509 ± 0.074	17.335 ± 0.163	1.586 ± 0.016	126.374 ± 0.227	163.987 ± 0.485
Xavier High Sparse	0.602 ± 0.007	13.451 ± 0.076	1.341 ± 0.030	43.935 ± 0.176	2.199 ± 0.012	27.054 ± 0.139	88.584 ± 0.349
Xavier Low Dense	2.128 ± 0.422	28.340 ± 0.098	6.907 ± 0.094	14.557 ± 0.292	1.511 ± 0.009	407.499 ± 0.302	460.941 ± 0.246
Xavier Low Sparse	2.163 ± 0.106	28.813 ± 0.047	1.728 ± 0.007	60.233 ± 0.076	2.385 ± 0.004	62.169 ± 0.073	157.492 ± 0.199
Robot Dense	1.531 ± 0.216	43.073 ± 0.482	29.237 ± 0.319	879.225 ± 6.116	115.363 ± 1.065	13.706 ± 0.064	1,082.135 ± 6.692
Robot Sparse	1.383 ± 0.111	41.073 ± 0.994	0.313 ± 0.004	66.045 ± 0.409	114.911 ± 1.171	13.012 ± 0.171	236.737 ± 2.491

TABLE III: Performance of PointPillars as % AP and performance of Sparse PointPillars and its ablations as the relative % AP difference (Δ) to PointPillars on KITTI with 16cm×16cm pillars. Higher is better.

	PointPillars			Sparse PointPillars			Sparse1+Dense23			Sparse12+Dense3			Sparse+WideConv		
	Easy	Med.	Hard	Easy	Med.	Hard	Easy	Med.	Hard	Easy	Med.	Hard	Easy	Med.	Hard
BEV AP	90.75	89.79	89.48	-5.71 Δ	-6.53 Δ	-4.98 Δ	-5.25 Δ	-5.44 Δ	-3.73 Δ	-5.54 Δ	-5.89 Δ	-5.60 Δ	-6.89 Δ	-7.24 Δ	-9.25 Δ
3D AP	82.30	80.34	79.11	-8.82 Δ	-7.83 Δ	-9.16 Δ	-7.62 Δ	-7.36 Δ	-8.85 Δ	-8.88 Δ	-7.71 Δ	-9.27 Δ	-11.1 Δ	-11.57 Δ	-13.69 Δ

TABLE IV: Per instance model component runtime and standard deviation in milliseconds, run on KITTI’s test set, averaged over ten trials. Models differ only in their Feature Net and Backbone; all other components are identical. Lower is better.

	To Device	Feature Extract	Feature Net	Backbone	Head	BBox Extract	Total time
PointPillars	0.064 ± 0.000	2.375 ± 0.010	0.303 ± 0.004	2.358 ± 0.027	0.204 ± 0.002	9.105 ± 0.030	14.410 ± 0.033
Sparse PointPillars	0.064 ± 0.000	2.330 ± 0.009	0.133 ± 0.002	7.578 ± 0.038	0.244 ± 0.004	3.877 ± 0.010	14.226 ± 0.056
Sparse1+Dense23	0.064 ± 0.000	2.341 ± 0.011	0.134 ± 0.002	7.394 ± 0.049	0.231 ± 0.003	4.437 ± 0.013	14.602 ± 0.068
Sparse12+Dense3	0.064 ± 0.000	2.359 ± 0.011	0.134 ± 0.002	7.803 ± 0.050	0.236 ± 0.002	4.406 ± 0.018	15.001 ± 0.080
Sparse+WideConv	0.066 ± 0.001	2.356 ± 0.015	0.137 ± 0.002	17.286 ± 0.071	0.242 ± 0.007	6.184 ± 0.038	26.270 ± 0.124

in the house mesh. The training and test splits along with the generation code are available on the project webpage.

Both PointPillars and Sparse PointPillars are trained with 5cm×5cm pillars and 768×512 pillar pseudoimage (set using the max absolute X and Y-axis values of the training data point clouds), 66%/33% train/validation splits, and standard hyperparameters, with the exception that Sparse PointPillars performs 50 more epochs in order to converge. Despite performing 25% more epochs, Sparse PointPillars trains 1.8× faster. Our evaluation follows the KITTI protocol of measuring the average precision (AP) at a detection threshold of 50% Intersection over Union (IoU) of the bounding box relative to ground truth on two key benchmarks: the bounding boxes from BEV (*BEV AP*) and in full 3D (*3D AP*).

Sparse PointPillars lags behind PointPillars in performance by 6.04% AP on *BEV* and by 4.61% AP on *3D* (PointPillars achieved 84.09% AP on *BEV* and 80.66% AP on *3D*). However, as shown in Table II, due to Matterport-Chair’s low density, Sparse PointPillars is significantly faster than PointPillars across the full range of compute platforms available to embodied agents: a desktop with an AMD Ryzen 7 3700X CPU and an NVidia 2080ti GPU (denoted *Desktop*), an NVidia Jetson Xavier embedded ML accelerator configured for the highest and lowest power settings (30 Watt, 8 core mode denoted *Xavier High* and 10 Watt, 2 core mode denoted *Xavier Low*), and a Fetch Freight [1] robot’s built-in four core Intel i5-4590S CPU (denoted *Robot*). Sparse PointPillars is more than 1.5× as fast as PointPillars on *Desktop* (fast enough for 60Hz inference), more than 2× as fast on *Xavier High* (fast enough for 10Hz inference), almost 3× as fast on *Xavier Low*

(fast enough for 6Hz inference), and more than 4× as fast on *Robot* (fast enough for 4Hz inference).

Of note, the recorded Backbone runtimes for Sparse PointPillars on GPU accelerated platforms is *slower* than PointPillars, but the BBox Extract stage is *faster* despite using *identical* code. The runtime difference comes from the time taken to allocate the memory for the anchor boxes—both allocate the same size GPU array, but due to pipelining and earlier memory cleanup that inflated the Sparse PointPillars’ Backbone’s runtime, it is able to allocate the final anchor boxes faster. The *Robot* evaluations demonstrate that when GPU pipelining is not a factor, Sparse PointPillars’ Backbone is far faster than PointPillars’ Backbone, and the BBox Extract stage runs at roughly the same speed.

B. KITTI Evaluation with Ablative Studies

KITTI [36], a self-driving car dataset of LiDAR point clouds with human-annotated 3D bounding boxes, is a common benchmark dataset in 3D object detection and is used as the evaluation dataset in the original PointPillars paper. Both Sparse PointPillars and PointPillars are trained on the KITTI Car detection task, configured with the default 16cm×16cm pillars, 504×440 pillar pseudoimage, 50%/50% train/validation split, and hyperparameter configurations outlined in the PointPillars paper, with the exception that Sparse PointPillars performs 50 more epochs in order to converge. Despite performing 25% more epochs, Sparse PointPillars trains in roughly the same amount of time. Our evaluation follows the prescribed KITTI evaluation protocol of measuring the average precision (AP) at a detection threshold of 70%

Intersection over Union (IoU) of the bounding box relative to ground truth on two key benchmarks: the bounding boxes from a BEV (*BEV AP*) and the full 3D bounding boxes (*3D AP*). KITTI does not have public labels for its test set, so in keeping with the literature [4], [25], [26] we report results on the validation set. Results are separated for the three KITTI difficulty levels (Easy, Medium, Hard), and runtimes are recorded on a dedicated desktop with an AMD Ryzen 7 3700X CPU and an NVidia 2080ti GPU.

Additionally, to better understand our contributions, we perform two ablative studies to answer these questions:

- 1) Would making the new Backbone only partially sparse provide a better performance-runtime trade-off?
- 2) Would approximating the flow of smeared information with wider convolutions improve performance?

To answer 1), we replace the later sections of our new Backbone with their dense counterparts from the original Backbone to construct two variants. Using Fig. 5’s `Conv` block definitions, the ablated variant *Sparse1+Dense23* uses the sparse `Conv` block 1 and dense `Conv` blocks 2 and 3, and the variant *Sparse12+Dense3* uses sparse `Conv` blocks 1 and 2 with a dense `Conv` block 3. To answer 2), we modify the filter size of the first SubM convolution of each `Conv` block to be 9×9 in order to simulate the information transfer caused by pseudoimage smearing in the original Backbone. We refer to this variant as *Sparse+WideConv*.

The absolute percentage of Average Precision (% AP) for PointPillars on each benchmark and the relative performance of Sparse PointPillars and its ablations are shown in Table III. Relative to PointPillars, Sparse PointPillars performs roughly 5% AP worse on *BEV* and roughly 8.5% AP worse on *3D*, and roughly equally to the ablative models, with *Sparse1+Dense23* performing slightly better and with *Sparse12+Dense3* and *Sparse+WideConv* performing worse. Together, these results indicate that SubM convolutional blocks in the Backbone are more difficult to train, even if the block has access to the same information as the dense model.

The runtime for each component of each method is reported in Table IV, with our Sparse PointPillars running 0.18ms faster than PointPillars. Our Feature Net runs 0.18ms faster as it avoids the *scatter* step, but like with Matterport-Chair, the recorded runtime for our Backbone is actually 5.22ms *slower* than PointPillars Backbone and the BBox Extract stage with its *identical* code is 5.22ms *faster* due to memory pipelining. Unsurprisingly, *Sparse1+Dense23* and *Sparse12+Dense3* are both slower than PointPillars and Sparse PointPillars due to the Backbone pipelining interruption when converting from a sparse to a dense tensor, and *Sparse+WideConv* is significantly slower due to its very large convolutions.

To contextualize Sparse PointPillars’ 14.2ms runtime vs PointPillars’ 14.4ms runtime in Table III, we refer to the reported KITTI runtime numbers for other sparsity based 3D detection approaches discussed in Section II. SECOND [26], a sparse 3D convolution-based object detector, reports a runtime of 50ms for its large variant (which the authors use to evaluate performance) and 25ms for its small variant using $20cm \times 20cm \times 40cm$ voxels. SBNNet [28], a 2D BEV object detector

that performs dense convolution inside coarse pseudolabel masks, reports 17.9ms runtime using $10cm \times 10cm$ pillars. Direct head-to-head performance and runtime comparisons against prior art can be found in the PointPillars [4] paper.

VI. CONCLUSION AND FUTURE WORK

This work demonstrates that Sparse PointPillars allows practitioners to trade small amounts of model performance for significant decreases in runtime and resource usage on embedded systems. For example, on our Matterport-Chair dataset, **Sparse PointPillars runs faster on the Jetson Xavier in low power mode than PointPillars does in high power mode**, allowing a practitioner to save power *and* get reduced runtimes at the cost of a few % AP. Alternatively, PointPillars runs at less than 1Hz on the robot’s CPU; **with Sparse PointPillars, practitioners can reliably run at 1Hz and have more than 75% of the CPU budget left to run other components of the robot control stack**. By providing faster runtimes via our architectural design, Sparse PointPillars provides practitioners new tools in their toolbox to build and optimize their full control stack.

This work can be extended by exploring model quantization and weight pruning in tandem with our new pipeline. Prior art has shown significant quantization of PointPillars results in only minor drops in performance [13]. When combined with Sparse PointPillars, this may result in significant further reductions in runtime for a modest drop in performance, or enable inference on more exotic hardware, e.g. FPGAs.

Additionally, this work would benefit from further performance evaluation using a Streaming AP [37] style measure extended to 3D detectors. In this work, we evaluated detection quality with AP, a standard metric in the vision literature that matches output detections to the *input* point cloud. However, this evaluation protocol does not represent the problem practitioners face: in dynamic environments, the detection is most useful if it matches the state of the world *at the time it is emitted*. The world changes while the detector is performing inference and so a quick, lower quality detection is potentially better representative of the world upon emission than a slow, higher quality detection. A streaming measure would directly consider the dramatic latency reductions of Sparse PointPillars in the evaluation of its accuracy, better reflecting the problem formulation that practitioners face.

Finally, sparse perception pipelines, such as this work, are (and will continue to be) **directly impactful for autonomous vehicle makers**. New LiDAR sensors offer increased detection ranges [38]–[41], causing increased pseudoimage area and reduced density, making Sparse PointPillars’ sparsity exploitation even more important in maintaining low runtimes.

ACKNOWLEDGMENTS

We gratefully acknowledge the useful feedback on this work from Ishan Khatri and Jorge Mendez. This research was partially supported by the DARPA Lifelong Learning Machines program under grant FA8750-18-2-0117, the DARPA SAIL-ON program under contract HR001120C0040, and the Army Research Office under MURI grant W911NF20-1-0080.

REFERENCES

- [1] M. Wise, M. Ferguson, D. King, E. Diehr, and D. Dymesich, "Fetch & Freight : Standard Platforms for Service Robot Applications," in *Proceedings of the 2016 International Joint Conference on Artificial Intelligence Workshop on Autonomous Mobile Service Robots*, 2016.
- [2] D. Franklin, V. Nguyen, and R. Love, "Jetson AGX Xavier and the New Era of Autonomous Machines," 2019.
- [3] X. Liu, G. V. Nardari, F. C. Ojeda, Y. Tao, A. Zhou, T. Donnelly, C. Qu, S. W. Chen, R. A. F. Romero, C. J. Taylor, and V. Kumar, "Large-Scale Autonomous Flight With Real-Time Semantic SLAM Under Dense Forest Canopy," *IEEE Robotics and Automation Letters*, 2022.
- [4] A. Lang, S. Vora, H. Caesar, L. Zhou, J. Yang, and O. Beijbom, "PointPillars: Fast Encoders for Object Detection From Point Clouds," in *Proceedings of the 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 12 689–12 697.
- [5] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "SSD: Single Shot MultiBox Detector," in *Proceedings of the European Conference on Computer Vision (ECCV) 2016*, 2016, pp. 21–37.
- [6] Z. Cai and N. Vasconcelos, "Rethinking Differentiable Search for Mixed-Precision Neural Networks," in *Proceedings of the 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [7] S. Shin, Y. Boo, and W. Sung, "Fixed-point optimization of deep neural networks with adaptive step size retraining," in *Proceedings of the 2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2017, pp. 1203–1207.
- [8] K. Wang, Z. Liu, Y. Lin, J. Lin, and S. Han, "HAQ: Hardware-Aware Automated Quantization With Mixed Precision," in *Proceedings of the 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [9] K. Zhao, S. Huang, P. Pan, Y. Li, Y. Zhang, Z. Gu, and Y. Xu, "Distribution Adaptive INT8 Quantization for Training CNNs," in *Proceedings of the Thirty-Fifth AAAI Conference on Artificial Intelligence (AAAI)*, 2021.
- [10] M. Kim and P. Smaragdis, "Bitwise neural networks," in *ICML Workshop on Resource-Efficient Machine Learning*, 2016.
- [11] B. Jacob, S. Kligys, B. Chen, M. Zhu, M. Tang, A. Howard, H. Adam, and D. Kalenichenko, "Quantization and Training of Neural Networks for Efficient Integer-Arithmetic-Only Inference," in *Proceedings of the 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [12] J. Wu, C. Leng, Y. Wang, Q. Hu, and J. Cheng, "Quantized Convolutional Neural Networks for Mobile Devices," in *Proceedings of the 2016 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [13] J. Stanisiz, K. Lis, T. Kryjak, and M. Gorgon, "Optimisation of the PointPillars network for 3D object detection in point clouds," in *2020 Signal Processing: Algorithms, Architectures, Arrangements, and Applications (SPA)*, 2020, pp. 122–127.
- [14] J. Frankle and M. Carbin, "The Lottery Ticket Hypothesis: Finding Sparse, Trainable Neural Networks," in *Proceedings of the Conference on Learning Representations, (ICLR)*, 2019.
- [15] Y. Lecun, J. Denker, and S. Solla, "Optimal Brain Damage," in *Advances in Neural Information Processing Systems (NeurIPS)*, 1989.
- [16] B. Hassibi, D. Stork, and G. Wolff, "Optimal Brain Surgeon and general network pruning," in *Proceedings of the IEEE International Conference on Neural Networks*, 1993.
- [17] R. Reed, "Pruning Algorithms - A Survey," *IEEE Transactions on Neural Networks*, vol. 4, no. 5, 1993.
- [18] D. Blalock, J. J. Gonzalez Ortiz, J. Frankle, and J. Guttag, "What is the State of Neural Network Pruning?" in *Proceedings of Machine Learning and Systems*, I. Dhillon, D. Papailiopoulos, and V. Sze, Eds., 2020.
- [19] J. Pool, A. Sawarkar, and J. Rodge, "Accelerating Inference with Sparsity Using the NVIDIA Ampere Architecture and NVIDIA TensorRT," 2021.
- [20] W. Wen, C. Wu, Y. Wang, Y. Chen, and H. Li, "Learning Structured Sparsity in Deep Neural Networks," 2016.
- [21] Y. He, X. Zhang, and J. Sun, "Channel Pruning for Accelerating Very Deep Neural Networks," in *Proceedings of the 2017 IEEE/CVF International Conference on Computer Vision (ICCV)*, 10 2017.
- [22] K. Ullrich, E. Meeds, and M. Welling, "Soft Weight-Sharing for Neural Network Compression," in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2017.
- [23] J. Frankle, G. Dziugaite, D. M. Roy, and M. Carbin, "Linear Mode Connectivity and the Lottery Ticket Hypothesis," in *Proceedings of the International Conference on Machine Learning (ICML)*, 2020.
- [24] S. Han, J. Pool, J. Tran, and W. J. Dally, "Learning Both Weights and Connections for Efficient Neural Networks," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2015.
- [25] Y. Zhou and O. Tuzel, "VoxelNet: End-to-End Learning for Point Cloud Based 3D Object Detection," in *Proceedings of the 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018.
- [26] Y. Yan, Y. Mao, and B. Li, "SECOND: Sparsely Embedded Convolutional Detection," *Sensors*, vol. 18, 2018.
- [27] A. Bewley, P. Sun, T. Mensink, D. Anguelov, and C. Sminchisescu, "Range Conditioned Dilated Convolutions for Scale Invariant 3D Object Detection," in *Conference on Robot Learning*, 2020.
- [28] M. Ren, A. Pokrovsky, B. Yang, and R. Urtasun, "SBNNet: Sparse Blocks Network for Fast Inference," in *Proceedings of the 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [29] P. A. Tew, "An Investigation of Sparse Tensor Formats for Tensor Libraries," Master's thesis, Massachusetts Institute of Technology, Cambridge, MA, September 2016.
- [30] R. Charles, H. Su, K. Mo, and L. Guibas, "PointNet: Deep Learning on Point Sets for 3D Classification and Segmentation," in *Proceedings of the 2017 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [31] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature Pyramid Networks for Object Detection," in *Proceedings of the 2017 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [32] B. Graham and L. van der Maaten, "Submanifold Sparse Convolutional Networks," *arXiv preprint arXiv:1706.01307*, 2017.
- [33] Q.-Y. Zhou, J. Park, and V. Koltun, "Open3D: A modern library for 3D data processing," *arXiv:1801.09847*, 2018.
- [34] C. Choy, J. Gwak, and S. Savarese, "4D Spatio-Temporal ConvNets: Minkowski Convolutional Neural Networks," in *Proceedings of the 2019 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [35] I. Armeni, O. Sener, A. R. Zamir, H. Jiang, I. Brilakis, M. Fischer, and S. Savarese, "3D Semantic Parsing of Large-Scale Indoor Spaces," in *Proceedings of the 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [36] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for Autonomous Driving? The KITTI Vision Benchmark Suite," in *Proceedings of the 2012 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012.
- [37] M. Li, Y. Wang, and D. Ramanan, "Towards Streaming Perception," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020.
- [38] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom, "nuScenes: A multimodal dataset for autonomous driving," *arXiv preprint arXiv:1903.11027*, 2019.
- [39] M.-F. Chang, J. W. Lambert, P. Sangkloy, J. Singh, S. Bak, A. Hartnett, D. Wang, P. Carr, S. Lucey, D. Ramanan, and J. Hays, "Argoverse: 3D Tracking and Forecasting with Rich Maps," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [40] B. Wilson, W. Qi, T. Agarwal, J. Lambert, J. Singh, S. Khandelwal, B. Pan, R. Kumar, A. Hartnett, J. K. Pontes, D. Ramanan, P. Carr, and J. Hays, "Argoverse 2: Next Generation Datasets for Self-Driving Perception and Forecasting," in *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021.
- [41] P. Sun, H. Kretschmar, X. Dotiwalla, A. Chouard, V. Patnaik, P. Tsui, J. Guo, Y. Zhou, Y. Chai, B. Caine, V. Vasudevan, W. Han, J. Ngiam, H. Zhao, A. Timofeev, S. Ettinger, M. Krivokon, A. Gao, A. Joshi, and D. Anguelov, "Scalability in Perception for Autonomous Driving: Waymo Open Dataset," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.