

# Q-LiDAR: Efficient and Accurate Training-Free Quantization for Point Cloud 3D Object Detection Models

Anonymous CVPR submission

Paper ID \*\*\*\*

## Abstract

001 *3D object detection with point clouds is crucial for appli-*  
002 *cations in autonomous driving, robotics, and augmented re-*  
003 *ality. As these applications advance towards real-time pro-*  
004 *cessing on edge devices, they demand models that enable*  
005 *efficient and flexible inference. While recent post-training*  
006 *quantization methods address some of these challenges by*  
007 *reducing model size and computational load without re-*  
008 *training, they are limited by calibration data biases and*  
009 *suffer from sub-optimal accuracy without precise calibra-*  
010 *tion. Moreover, existing methods often overlook the distinct*  
011 *characteristics of various components in 3D LiDAR models,*  
012 *where high variance in value distributions poses additional*  
013 *quantization challenge. To overcome these issues, we pro-*  
014 *pose Q-LiDAR, a novel quantization approach that incor-*  
015 *porates techniques, including SmoothQConv, fine-grained*  
016 *quantization for sparse operators, and Hessian-guided bit-*  
017 *width allocation. Our approach achieves W4A8 mixed-*  
018 *precision quantization on state-of-the-art 3D LiDAR mod-*  
019 *els while retaining XX% model accuracy, without requiring*  
020 *a calibration dataset or retraining.*

## 021 1. Introduction

022 3D object detection with point clouds is a critical task in  
023 various applications such as autonomous driving, robotics,  
024 and augmented reality. These applications rely on accurate  
025 and efficient detection of objects in 3D space to navigate  
026 and interact with their environment safely and effectively.  
027 As they move toward real-time processing on edge devices,  
028 the demand for efficient models has grown even higher.

029 Model quantization has proven to be an effective com-  
030 pression method. By compressing high bit-width floating-  
031 point (FP) data into lower bit-width integers, the computa-  
032 tional and memory costs of the model can be significantly  
033 reduced. Prior works have studied quantization methods  
034 for 2D object detection models and achieved promising re-  
035 sults [20, 22, 34]. However, directly applying these quan-

tization methods to 3D point cloud object detection models  
leads to sub-optimal accuracy [10]. Moreover, prior work  
often relies on quantization-aware training (QAT) [10, 37],  
which requires extensive fine-tuning, limiting its flexibility  
for rapid deployment in resource-constrained environments.

Recently, LiDAR-PTQ introduces a post-training quan-  
tization (PTQ) approach that reduces model size and com-  
putational demands for 3D object detection models without  
retraining [38]. While achieving promising results, LiDAR-  
PTQ faces two main challenges:

- **Calibration data bias:** Despite eliminating the need for retraining, LiDAR-PTQ relies on calibration data during the quantization process. The quality of the quantization can be negatively affected depending on the calibration data provided. For example, if the calibration data is not representative of the full dataset, the quantization might not generalize well.
- **Poor accuracy without calibration:** Empirical results show that quantizing LiDAR models for 3D object detection with point cloud is challenging due to the complex mix of diverse layers specifically designed for point cloud processing. We observe high variance across four different components in LiDAR models: (1) 2D/1D convolution (Conv2D/1D), (2) sparse convolution (SPConv3D), (3) submanifold convolution (SubMConv3D), and (4) multi-layer perceptron (MLP). The unique distributions of activations and weights across these components result in significant variations, making it very challenging to apply a uniform set of quantization parameters, such as, W8A8, to values with high variance.

Based on our observations of the data distribution, we draw ideas from the large language model (LLM) compression literature, where training-free quantization methods have been shown to be effective [8, 27, 33]. To avoid quantization errors from outliers in convolution operations, we extend a smoothing-based quantization technique [33] to transform the quantization difficulties in activations to convolution weights or vice versa. To compress SPConv3D and SubMConv3D layers, we adopt channel-wise quantization. However, when assigning proper bit-width for each compo-

076 ent, we find that simple Mean Square Error with the un-  
077 compressed model is inadequate to distinguish quantization  
078 effects on various components. To overcome this, we pro-  
079 pose to leverage Hessian information to estimate the quan-  
080 tization sensitivity of each component and guide mixed-  
081 precision bit-width allocation across components.

082 The contributions of the paper are four-fold:

- 083 1. **Comprehensive Component Analysis:** We conduct an  
084 in-depth analysis of the data distribution across different  
085 components in 3D LiDAR object detection models. This  
086 investigation reveals the unique quantization challenges  
087 associated with the diverse layers in these models.
- 088 2. **Development of Calibration-Free Quantization**  
089 **Method:** We propose a calibration-free quantization  
090 method, Q-LiDAR, that employs component-specific  
091 quantization strategies, including SmoothQConv  
092 for Conv2D/1D and MLP layers, and channel-wise  
093 quantization for sparse operators.
- 094 3. **Sensitivity-Based Mixed-Precision Quantization:** To  
095 address the bit-width allocation challenges in vari-  
096 ous components, such as SPConv3D, SubMConv3D,  
097 Conv2/1d, and MLP layers, Q-LiDAR incorporates  
098 a sensitivity-based bit-width allocation policy based  
099 on Hessian information, tailored to each component’s  
100 unique characteristics to mitigate the accuracy loss.
- 101 4. **Extensive Experimental Validation:** We validate the  
102 effectiveness of Q-LiDAR across a range of state-of-the-  
103 art LiDAR models for 3D object detection. Experimen-  
104 tal results demonstrate that Q-LiDAR achieves XX com-  
105 pression ratio while obtaining very comparable accuracy  
106 (XX%) as the uncompressed model including [Hongbo:  
107 will add the finalized models over here]. [Explicit accu-  
108 racy here]. We empirically select recent advances that  
109 have been adopted in many industries to construct a uni-  
110 fied baseline. The follow-up experiments show that we  
111 achieved [n%], [n%] and [n%] respectively on KITTI,  
112 nuScenes and Waymo datasets as compared to direct  
113 quantization and achieve mAP and NDS of [n%], [n%]  
114 and [n%]. [Hongbo: evaluation results here]

## 115 2. Related Work

116 **3D object detection.** 3D object detection (3DOD) is a  
117 pivotal area of research for autonomous driving, robotics,  
118 and augmented reality. This process heavily relies on so-  
119 phisticated sensor technologies such as LiDAR (Light De-  
120 tection and Ranging), radar, and stereo vision cameras that  
121 capture detailed three-dimensional information about the  
122 environment. Among them, LiDAR has become one of the  
123 most widely used sensors for its real-time feedback and high  
124 accuracy, and since the data collected are separated points  
125 with different properties, they are also called point cloud.

126 Several notable methods are introduced to capture pre-

127 cise 3D spatial information, including PointNet [24] and  
128 its variants PointNet++ [25] and PointNeXt [26], which di-  
129 rectly process point clouds, and voxel-based methods like  
130 VoxelNet [39], Voxel Transformer [21], and VoxelNeXt [4],  
131 which convert point clouds into structured grids for easier  
132 processing.

133 To efficiently manage sparse point cloud data, which is  
134 inherently memory-intensive, prior work introduce custom  
135 layers, such as SparseConv and SubMConv layers to han-  
136 dle sparse point cloud data [5]. These layers leverages the  
137 inherit sparsity of the input data by performing convolu-  
138 tions exclusively on non-zero elements, which drastically  
139 reduces both memory consumption and computational over-  
140 head, making it particularly suitable for large-scale 3D data  
141 processing.

**Training-free quantization.** While early model com-  
142 pression techniques focus on improving model accuracy  
143 through retraining or fine-tuning [6, 7, 15, 19, 23], they  
144 face challenges in flexibility, which hinders the widespread  
145 adoption of those methods across diverse deployment en-  
146 vironments. Recent advancements in training-free com-  
147 pression have significantly improved the efficiency and  
148 deployment of vision transformers [13, 14, 17, 20, 36].  
149 In NLP, techniques such as GPTQ [9], AWQ [16],  
150 SmoothQuant [33] have also demonstrated their success  
151 in quantizing large language models. These advancement  
152 highlight the ongoing efforts to compress DNN models.  
153 While demonstrating promising results, few studies have  
154 looked into training-free compression for 3D LiDAR object  
155 detection models.  
156

## 157 3. Methodology

158 In this section, we first introduce Q-LiDAR, a novel  
159 training-free quantization method to compress 3D object  
160 detectors while retaining accuracy. And then we develop  
161 a Hessian-guided method for bit-width allocation to reduce  
162 quantization errors.

163 Giving the hybrid architecture of 3D LiDAR mod-  
164 els, we investigate the impact of quantization on specific  
165 components within 3D object detection models, especially  
166 six components commonly used in 3DOD models: Spar-  
167 seConv3d, SubMConv3d, SparseConv2d, SubMConv2d,  
168 Conv2d/1d, MLP.

### 169 3.1. Improving 3D LiDAR Model Compression via 170 SmoothQConv

171 We start by directly applying W8A8 post-training quan-  
172 tization to 3D LiDAR models. However, we find that  
173 W8A8 leads to large accuracy drop. Table 1 shows the  
174 results of various combination of quantization bit-width  
175 with both static and dynamic round-to-nearest (RTN) post-  
176 training quantization over CenterPoint-Voxel [35] and the

177 autonomous driving dataset NuScenes *val*. The quantiza-  
178 tion operation is formulated as:

$$179 \quad \mathbf{X}^{\text{INT}} = \text{clamp} \left( \left\lfloor \frac{x}{s} \right\rfloor + z, q_{\min}, q_{\max} \right) \quad (1)$$

180 where  $\lfloor \cdot \rfloor$  is the rounding-to-nearest operator,  $s$  is the scal-  
181 ing factor, and  $z$  is the zero-point. As shown, while W8A8  
182 quantization has been considered quite robust to leads to  
183 traditional 2D convolution tasks [], it causes around 20  
184 mAP and NDS loss, which is quite significant. In contrast,  
185 W4A16 quantization leads to relatively lower accuracy loss.

Table 1. Quantization Results with Performance Gaps

Method	Bits(W/A)	Metrics	
		mAP	NDS
Full Prec.	32/32	59.22	66.48
Dynamic	8/8	39.56 (-19.66)	47.63 (-18.85)
Static	8/8	38.46 (-20.76)	46.13 (-20.35)
Dynamic	4/8	34.36 (-24.86)	44.87 (-21.61)
Static	4/8	33.80 (-25.42)	44.16 (-22.32)
Dynamic	4/16	51.24 (-7.98)	59.38 (-7.10)
Static	4/16	51.24 (-7.98)	59.38 (-7.10)
Dynamic	16/4	xx.xx (-7.98)	xx.xx (-7.10)
Static	16/4	xx.xx (-7.98)	xx.xx (-7.10)

186 To investigate why activation quantization leads to more  
187 significant accuracy drop, we further collect 1) the top-8  
188 maximum weight values; and 2) the average of the top-8  
189 activation values. As illustrated in Figure 1 and 2, the re-  
190 sults show that the activations, especially those convolution  
191 layers, contain outliers, whereas the weight values contain a  
192 much smaller dynamic range. Notably, the majority of out-  
193 liers are found in the activation values with the maximum  
194 value being up to 120, and the scaling of the activation is  
195 way greater than that of the weight(magnitude of 30 com-  
196 pared to 2). This confirms the imbalanced scaling and mag-  
197 nitude of the activation and weight value within the model.

198 [Minjia: TODO: Add the figure that show the dynamic  
199 range results of weights and activations across layers here.  
200 I remember asking Banghao to collect these results before,  
201 so we should have them. Also, it would be better to show  
202 the results of conv2D, e.g., the one that correspond to the  
203 120.]

204 To mitigate the errors introduced by extreme outliers  
205 and the imbalanced quantization difficulty between activa-  
206 tion and weight, a technique called SmoothQuant [33] can  
207 be implemented. SmoothQuant redistributes quantization  
208 complexity from one tensor to another (e.g., from activation  
209 to weight). This approach is especially effective to reduce  
210 outlier impact, particularly in the context of linear operators  
211 in large language models (LLMs), where outliers are often  
212 found in per-token areas. However, applying SmoothQuant

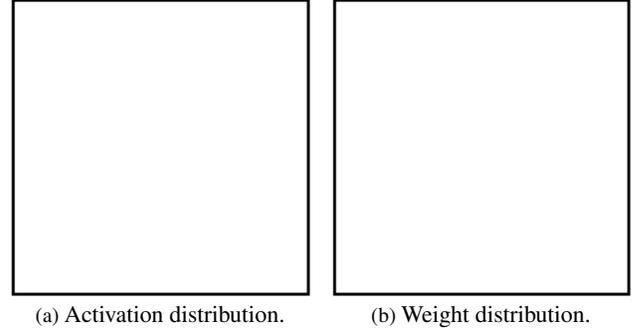


Figure 1. The left figure shows the dynamic range of activations across different convolution layers. The right figure shows the dynamic range of weights across layers.

213 to 3D LiDAR models presents a unique challenge, as there  
214 is no direct mathematical mechanism for shifting quanti-  
215 zation complexity between activations and convolutional  
216 weights.

217 To overcome this, we introduce SmoothQConv, an ex-  
218 tension of SmoothQuant tailored to convolutional operators.  
219 Our primary insight is that convolution can be reformulated  
220 as a matrix multiplication by transforming the input data  
221 through the *im2col* (image-to-column) operation [3]. The  
222 *im2col* operation rearranges the input activation (feature  
223 map) into a matrix where each column represents a local  
224 region (receptive field) of the input that the convolutional  
225 filter will slide over. This process effectively “unfolds” the  
226 input data into a 2D matrix. allowing the convolution to be  
227 treated as standard matrix multiplication. The resulting ma-  
228 trix from the multiplication is then reshaped back (“fold”)   
229 into the original spatial dimensions of the output feature  
230 map.

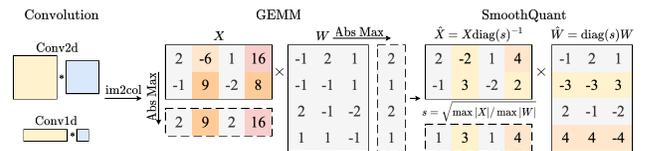


Figure 2. Overview of SmoothQConv operation when  $\alpha = 0.5$ . “\*” and  $\times$  indicate convolution and matrix multiplication operations respectively.

231 The original General Matrix Matrix Multiplication  
232 (GEMM) floating-point operation of Conv2d after unfold-  
233 ing is:

$$234 \quad \mathbf{Y} = \mathbf{X}^{\text{FP32}} \mathbf{W}^{\text{FP32}} \quad (2)$$

235 where  $\mathbf{X} \in \mathbb{R}^{bn \times ihw}$  and  $\mathbf{W} \in \mathbb{R}^{ihw \times c}$ .

236 To leverage INT8 GEMM acceleration on general hard-  
237 ware, we implement weight-per-channel and activation-per-  
238 tensor quantization. The output  $\mathbf{Y}$  is approximated using

Table 2. Notations

in channels	$i$
out channels	$c$
kernel height	$h$
kernel width	$w$
kernel depth	$d$
batch size	$b$
number of sliding	$n$
number of active voxels	$v$

239 quantized INT8 operands as:

$$240 \mathbf{Y} \approx (\hat{\mathbf{X}}^{\text{INT8}} \odot \Delta_{\mathbf{X}}^{\text{FP32}})(\hat{\mathbf{W}}^{\text{INT8}} \text{diag}(\Delta_{\mathbf{W}}^{\text{FP32}})) \quad (3)$$

$$241 \approx \text{diag}(\Delta_{\mathbf{X}}^{\text{FP32}})(\hat{\mathbf{X}}^{\text{INT8}} \hat{\mathbf{W}}^{\text{INT8}}) \text{diag}(\Delta_{\mathbf{W}}^{\text{FP32}}) \quad (4)$$

242 where  $\hat{\mathbf{X}}^{\text{INT8}}$  and  $\hat{\mathbf{W}}^{\text{INT8}}$  are the quantized activation and  
243 weight matrices, and  $\Delta_{\mathbf{X}} \in \mathbb{R}$  and  $\Delta_{\mathbf{W}} \in \mathbb{R}^c$  denote  
244 the scaling factors for activation-per-tensor and weight-per-  
245 channel quantization respectively.

246 The quantization of activations and weights is defined as:

$$247 \hat{\mathbf{X}}^{\text{INT8}} = \left\lfloor \frac{\mathbf{X}^{\text{FP32}}}{\Delta_{\mathbf{X}}^{\text{FP32}}} \right\rfloor \quad \hat{\mathbf{W}}^{\text{INT8}} = \left\lfloor \frac{\mathbf{W}^{\text{FP32}}}{\Delta_{\mathbf{W}}^{\text{FP32}}} \right\rfloor \quad (5)$$

248 where  $\lfloor \cdot \rfloor$  denotes rounding to the nearest integer.

249 The scaling factors are computed to map the floating-  
250 point values to the INT8 quantization range:

$$251 \Delta_{\mathbf{X}} = \frac{\max(\mathbf{X}^{\text{fp32}}) - \min(\mathbf{X}^{\text{fp32}})}{2^b - 1} \quad (6)$$

252 and

$$253 \Delta_{\mathbf{W},j} = \frac{\max_{i=1,\dots,ihw}(\mathbf{W}_{ij}^{\text{fp32}}) - \min_{i=1,\dots,ihw}(\mathbf{W}_{ij}^{\text{fp32}})}{2^b - 1} \quad (7)$$

254 To extend the SmoothQuant technique to convolution op-  
255 erators, we introduce a dedicated hyperparameter  $\alpha$ , which  
256 controls the degree to which quantization difficulty is redis-  
257 tributed between tensors. The scaling value  $s_k$  is calculated  
258 as:

$$259 s_k = \max(|\mathbf{X}_k^{\text{FP32}}|)^\alpha / \max(|\mathbf{W}_k^{\text{FP32}}|)^{1-\alpha} \quad (8)$$

260 Utilizing  $s_k$ , we can transfer the quantization difficulty  
261 from one to the other by applying this scaling value to our  
262 activation and weight prior to the actual quantization stage.  
263 This yields the new quantized representations:

$$264 \hat{\mathbf{X}}^{\text{INT8}} = \left\lfloor \frac{\mathbf{X}^{\text{FP32}}}{\Delta_{\mathbf{X}}^{\text{FP32}}} \text{diag}(s_k)^{-1} \right\rfloor \quad (9)$$

$$266 \hat{\mathbf{W}}^{\text{INT8}} = \left\lfloor \frac{\mathbf{W}^{\text{FP32}}}{\Delta_{\mathbf{W}}^{\text{FP32}}} \text{diag}(s_k) \right\rfloor \quad (10)$$

The final matrix multiplication is then approximated using these quantized INT8 operands,  $\hat{\mathbf{X}}$  from equation 9 and  $\hat{\mathbf{W}}$  from equation 10, as per the approximation in equation 4.

As seen in Figure 1 and 2, the scaling of activation is way greater than that of weight, meaning  $s_k$  is mostly greater than 1. Therefore, after scaling up the weight and scaling down the activation, we manage to reduce the rounding error by reducing the value of  $\max(|\mathbf{X}|)$  in every tensor. Hence, we successfully reduce the quantization error and achieve more accurate results of the convolution operation even though the parameters of the convolution is INT8 quantized.

### 3.2. Fine-grained Quantization for Sparse Convolutions

3D LiDAR models incorporate sparse operators, such as submanifold convolution [12] and sparse convolution [18], to reduce the computation load. Specifically, these operators selectively process only the active voxels, bypassing non-active regions. Let  $\mathbf{x}_u$  represent an input feature vector of an active voxel located at 3-dimensional coordinates  $u \in \mathbb{R}^3$ . The submanifold operator  $F_0$  by a kernel for  $\mathbf{X}_u$  is formulated as:

$$F_0(\mathbf{W}, \mathbf{X}_u) = \sum_{i \in N(u)} \mathbf{W}_i \mathbf{X}_{u+i} \quad (11)$$

where  $N(u)$  denotes the set of offsets in the 3-dimensional cube centered at origin relative to  $u$ . Each offset is associated with a specific kernel weight parameterized by  $\mathbf{W}_i$ .

Since the sparse operators perform convolution only in active regions of the feature map, we adopt a channel-wise quantization approach for both weights and activations in SPConv and SubMConv layers. The weights of these sparse convolutions,  $\mathbf{W} \in \mathbb{R}^{c \times i \times h \times w \times d}$ , extend Conv2d weights with an additional depth dimension. To quantize these weights, we first reshape  $\mathbf{W}$  into a 2D matrix  $\mathbf{W} \in \mathbb{R}^{c \times ihwd}$  and apply channel-wise quantization along the output channel dimension  $c$ .

Activations in sparse convolution layers, represented as  $\mathbf{X} \in \mathbb{R}^{v \times 3}$ , where  $v$  is the number of active voxels in the feature map and 3 indicates the coordinates  $(x, y, z)$  of each active voxel, are similarly quantized. We apply channel-wise quantization along each spatial axis  $(x, y, z)$  of the coordinates to ensure independent quantization for each spatial dimension.

### 3.3. Searching to Allocate Bit-Width/Hessian-Guided Bit-width Allocation

[Hongbo: start fixing here] To consider layer sensitivity, we choose to automatically search for an optimal bit-width allocation policy that minimizes the output difference (e.g., L1 loss) after the quantization for a certain layer.

316 [Minjia: TODO: Depending on the Hessian results, we  
317 may consider the Hessian-guide bit-width allocation or the  
318 original sensitivity analysis based quantization. @Banghao,  
319 please share the Hessian results as soon as you get them.]

---

**Algorithm 1** Auto-Sensitive Analysis
 

---

**Require:** Pretrained FP model with  $N$  layers; Calibration dataset  $D^c$ ; Standard quantization module replacement map  $M^q$ ; **XXX** module replacement map  $M^{sq}$ ; Calibration number  $T$

**Ensure:** quantization method assigned to each type of layer using corresponding map  $M_q$  or  $M_{sq}$ .

- 1: Input  $T$  samples of  $D^c$  to FP network to get averaged FP output of each layer  $O_{fp}$ ;
  - 2: **for**  $L_i = \{L_i | i = 1, 2, \dots, N\}$  **do**
  - 3: Find quantized layer  $L_i^q$  with map  $M^q$ ;
  - 4: Replace the original layer within the model  $L_i$  with quantized layer  $L_i^q$ ;
  - 5: **end for**
  - 6: Input  $T$  samples of  $D^c$  to standard-quantized network to get averaged standard-quantized output of each layer  $O^{qint}$ ;
  - 7: **for**  $L_i = \{L_i | i = 1, 2, \dots, N\}$  **do**
  - 8: Find quantized layer  $L_i^{sq}$  with map  $M^{sq}$ ;
  - 9: Replace the original layer within the model  $L_i$  with quantized layer  $L_i^{sq}$ ;
  - 10: **end for**
  - 11: Input  $T$  samples of  $D^c$  to **XXX** network to get averaged **XXX** output of each layer  $O^{sqint}$ ;
  - 12: Check standard-quantized network output  $O^{qint}$  and FP final output  $O_{fp}$  to calculate  $L1_{qint}$ ;
  - 13: Check **XXX** network output  $O^{sqint}$  and FP final output  $O_{fp}$  to calculate  $L1_{sqint}$ ;
  - 14: Check  $L1_{qint}$  and  $L1_{sqint}$  to get the list of layers that can be quantized under **XXX**, with others being standard-quantized;
- 

## 320 4. Experiments

321 We conduct experiments to evaluate the effectiveness of Q-  
322 LiDAR in terms of accuracy preserving and compression  
323 ratio. Our evaluation aims to answer the following ques-  
324 tions:

- 325 • Can Q-LiDAR enable high compression ratio for 3D Li-  
326 DAR models without compromising accuracy?
- 327 • Does Q-LiDAR effectively generalize across diverse  
328 model architectures?
- 329 • How does Q-LiDAR compare to existing 3D LiDAR  
330 model compression methods in terms of the trade-off be-  
331 tween compression ratio and accuracy?

## 4.1. Evaluation Methodology

**Models** Our experiments include both transformer-based and convolution-based state-of-the-art 3D LiDAR models. The transformer-based models, DSVT-Voxel [32] and TransFusion-L [1], feature a voxel transformer backbone, a 2D convolution backbone, and a dense convolutional head. In contrast, the convolution-based models, PV-RCNN++[30], PV-RCNN[28], Part- $A^2$ -Anchor [29], and CenterPoint-Voxel [35], are equipped with a sparse 3D convolution backbone, a 2D convolution backbone, and a convolutional dense head.

**Datasets** We use Waymo Open Dataset (WOD) [31], nuScenes [2], and KITTI [11] for evaluation.

**Baselines** We evaluate the performance of Q-LiDAR by comparing it against two primary baselines: (1) the full-precision model (FP32) to establish an upper-bound reference, and (2) the standard Max-min quantized model (W8A8), commonly used in edge deployments. Since LiDAR-PTQ does not have their code released, we skip it for quantitative comparison.

## 4.2. Performance Comparison on Datasets

**Waymo Dataset.** To evaluate the performance of Q-LiDAR, several experiments are performed with DSVT-Voxel and PV-RCNN++ models on the Waymo dataset.

Table 3. Waymo Results for Different Detectors.

Models	Methods	Bits(W/A)	Vehicle	Pedestrian	Cyclist
DSVT-Voxel	Full Prec.	32/32	x.x	x.x	
	Max-min	8/8	x.x	x.x	
	QL-0.XX	8/8	x.x	x.x	
	Max-min	4/8	x.x	x.x	
	QL-X.XX	4/8	x.x	x.x	
PV-RCNN++	Full Prec.	32/32	67.68	60.17	72.55
	Max-min	8/8	x.x	x.x	
	QL-0.40	8/8	67.01	59.57	72.17
	Max-min	4/8	x.x	x.x	
	QL-0.xx	4/8	x.x	x.x	

**nuScenes Dataset.** To evaluate the performance of Q-LiDAR, several experiments are performed with BEVFusion and TransFusion-L models on the nuScenes dataset.

**KITTI Dataset.** To evaluate the performance of our method, we conduct experiments on 2 models, PV-RCNN and Part- $A^2$ , on KITTI dataset.

As shown in Table 5, achieves superior performance compared to max-min quantization method. It manages to

Table 4. nuScenes Results for Different Detectors.

Models	Methods	Bits(W/A)	mAP	NDS
TransFusion-L	Full Prec.	32/32	x.x	x.x
	Max-min	8/8	x.x	x.x
	QL-X.XX	8/8	x.x	x.x
	Max-min	4/8	x.x	x.x
	QL-X.XX	4/8	x.x	x.x
CP-Voxel	Full Prec.	32/32	59.22	66.48
	Max-min	8/8	x.x	x.x
	QL-0.80	8/8	59.16	66.40
	Max-min	4/8	x.x	x.x
	QL-X.XX	4/8	x.x	x.x

Table 5. KITTI Results for Different Detectors

Models	Methods	Bits(W/A)	Car	Pedestrian	Cyclist
PV-RCNN	Full Prec.	32/32	83.69	54.84	68.92
	Max-min	8/8	79.28	54.65	69.45
	QL-0.30	8/8	82.98	54.83	69.65
	Max-min	4/8	78.01	56.54	62.88
	QL-0.40	4/8	78.74	54.74	67.41
Part-A <sup>2</sup> -Anchor	Full Prec.	32/32	79.40	60.11	69.92
	Max-min	8/8	79.40	60.90	70.67
	QL-0.40	8/8	79.41	60.28	69.95
	Max-min	4/8	78.17	53.28	67.18
	QL-0.35	4/8	79.25	55.05	68.77

364 minimize the accuracy loss within less than 1% for both  
365 W8A8 and W4A8.

366 [Minjia: TODO: Add more in-depth description of the  
367 results. 1. Describe how to interpret the results in the table.  
368 2. Main observations. 3. Explanation of why we see these  
369 results.]

### 370 4.3. Ablation Study

371 We conducted an ablation study to evaluate the effects of  
372 the three key components of our framework, using the XXX  
373 model on the XXX dataset. As illustrated in Table 6, the ap-  
374 plication of channel-wise quantization to the 3D backbone  
375 network yielded a modest improvement in performance.  
376 Building on this, the introduction of SmoothQuant to the  
377 model’s 2D backbone resulted in a substantial performance  
378 leap from xx.x to xx.x. Finally, by employing the auto-  
379 sensitive-analysis algorithm (Algorithm 1) to identify and  
380 exclude layers particularly susceptible to quantization, we  
381 achieved a peak accuracy of xx.x.

382 [Minjia: TODO: Suggested ablation studies: 1. Q-  
383 LiDAR, 2. Q-LiDAR- layer sensitivity, 3. Q-LiDAR- layer  
384 sensitivity - channelwise quantization for sparse ops, 4.  
385 Q-LiDAR- layer sensitivity - channelwise quantization for

sparse ops - SmoothQConv.]

## 5. Conclusion

In this work, we introduce a training-free approach for effi-  
cient and accurate 3D object detection. Our approach em-  
ploys tailored optimizations against different components in  
3D LiDAR models, including SmoothQConv, subchannel-  
wise grouped quantization for SPConv and SubMConv. Ad-  
ditionally, we introduce a Hessian-guided method for bit-  
width allocation. Together, Q-LiDAR achieves state-of-the-  
art compression ratio for LiDAR models over 3D object de-  
tection.

## Acknowledgments

Use unnumbered third level headings for the acknowledg-  
ments. All acknowledgments, including those to funding  
agencies, go at the end of the paper.

## References

- [1] Xuyang Bai, Zeyu Hu, Xinge Zhu, Qingqiu Huang, Yilun Chen, Hongbo Fu, and Chiew-Lan Tai. Transfusion: Robust lidar-camera fusion for 3d object detection with transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1090–1099, 2022. 5
- [2] 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*, pages 11621–11631, 2020. 5
- [3] Kumar Chellapilla, Sidd Puri, and Patrice Simard. High performance convolutional neural networks for document processing. In *Tenth international workshop on frontiers in handwriting recognition*. Suvisoft, 2006. 3
- [4] Yukang Chen, Jianhui Liu, Xiangyu Zhang, Xiaojuan Qi, and Jiaya Jia. Voxelnex: Fully sparse voxelnet for 3d object detection and tracking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 21674–21683, 2023. 2
- [5] Spconv Contributors. Spconv: Spatially sparse convolution library. <https://github.com/traveller59/spconv>, 2022. 2
- [6] Peiyan Dong, LEI LU, Chao Wu, Cheng Lyu, Geng Yuan, Hao Tang, and Yanzhi Wang. Packqvit: Faster sub-8-bit vision transformers via full and packed quantization on the mobile. In *Advances in Neural Information Processing Systems*, pages 9015–9028. Curran Associates, Inc., 2023. 2
- [7] Zhen Dong, Zhewei Yao, Amir Gholami, Michael W. Mahoney, and Kurt Keutzer. Hawq: Hessian aware quantization of neural networks with mixed-precision. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019. 2
- [8] Zhen Dong, Zhewei Yao, Amir Gholami, Michael W Mahoney, and Kurt Keutzer. Hawq: Hessian aware quantization

Table 6. Ablation study of different components of XXX on XXX dataset.

Models	Methods	Bits(W/A)	mAP	NDS
	Full Prec.	32/32	x.x	x.x
DSVT-Voxel	+ASA	8/8	x.x	x.x
	+CWQ	8/8	x.x	x.x
	+SQConv	8/8	x.x	x.x

- 438 of neural networks with mixed-precision. In *Proceedings of*  
439 *the IEEE/CVF international conference on computer vision*,  
440 pages 293–302, 2019. 1
- 441 [9] Elias Frantar, Saleh Ashkboos, Torsten Hoefer, and Dan  
442 Alistarh. Gptq: Accurate post-training quantization  
443 for generative pre-trained transformers. *arXiv preprint*  
444 *arXiv:2210.17323*, 2022. 2
- 445 [10] Huan-ang Gao, Beiwen Tian, Pengfei Li, Hao Zhao, and  
446 Guyue Zhou. Dqs3d: Densely-matched quantization-  
447 aware semi-supervised 3d detection. In *Proceedings of the*  
448 *IEEE/CVF International Conference on Computer Vision*,  
449 pages 21905–21915, 2023. 1
- 450 [11] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel  
451 Urtasun. Vision meets robotics: The kitti dataset. *The Inter-*  
452 *national Journal of Robotics Research*, 32(11):1231–1237,  
453 2013. 5
- 454 [12] Benjamin Graham and Laurens Van der Maaten. Sub-  
455 manifold sparse convolutional networks. *arXiv preprint*  
456 *arXiv:1706.01307*, 2017. 4
- 457 [13] Jung Hwan Heo, Arash Fayyazi, Mahdi Nazemi, and Mas-  
458 soud Pedram. A fast training-free compression framework  
459 for vision transformers. *CoRR*, abs/2303.02331, 2023. 2
- 460 [14] Woosuk Kwon, Sehoon Kim, Michael W. Mahoney, Joseph  
461 Hassoun, Kurt Keutzer, and Amir Gholami. A fast post-  
462 training pruning framework for transformers. In *Advances in*  
463 *Neural Information Processing Systems 35: Annual Confer-*  
464 *ence on Neural Information Processing Systems 2022*, 2022.  
465 2
- 466 [15] Yanjing Li, Sheng Xu, Baochang Zhang, Xianbin Cao, Peng  
467 Gao, and Guodong Guo. Q-vit: Accurate and fully quan-  
468 tized low-bit vision transformer. In *Advances in Neural In-*  
469 *formation Processing Systems*, pages 34451–34463. Curran  
470 Associates, Inc., 2022. 2
- 471 [16] Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Wei-Ming  
472 Chen, Wei-Chen Wang, Guangxuan Xiao, Xingyu Dang,  
473 Chuang Gan, and Song Han. Awq: Activation-aware weight  
474 quantization for on-device llm compression and acceleration.  
475 *Proceedings of Machine Learning and Systems*, 6:87–100,  
476 2024. 2
- 477 [17] Yang Lin, Tianyu Zhang, Peiqin Sun, Zheng Li, and  
478 Shuchang Zhou. Fq-vit: Post-training quantization  
479 for fully quantized vision transformer. *arXiv preprint*  
480 *arXiv:2111.13824*, 2021. 2
- 481 [18] Baoyuan Liu, Min Wang, Hassan Foroosh, Marshall Tappen,  
482 and Marianna Pinsky. Sparse convolutional neural networks.  
483 In *Proceedings of the IEEE conference on computer vision*  
484 *and pattern recognition*, pages 806–814, 2015. 4
- [19] Shih-Yang Liu, Zechun Liu, and Kwang-Ting Cheng. 485  
Oscillation-free quantization for low-bit vision transformers. 486  
In *Proceedings of the 40th International Conference on Ma-* 487  
*chine Learning*, pages 21813–21824. PMLR, 2023. 2 488
- [20] Zhenhua Liu, Yunhe Wang, Kai Han, Wei Zhang, Siwei Ma, 489  
and Wen Gao. Post-training quantization for vision trans- 490  
former. *Advances in Neural Information Processing Systems*, 491  
34:28092–28103, 2021. 1, 2 492
- [21] Jiageng Mao, Yujing Xue, Minzhe Niu, Haoyue Bai, Jiashi 493  
Feng, Xiaodan Liang, Hang Xu, and Chunjing Xu. Voxel 494  
transformer for 3d object detection. In *Proceedings of the* 495  
*IEEE/CVF International Conference on Computer Vision* 496  
*(ICCV)*, pages 3164–3173, 2021. 2 497
- [22] Markus Nagel, Rana Ali Amjad, Mart Van Baalen, Chris- 498  
tos Louizos, and Tijmen Blankevoort. Up or down? adap- 499  
tive rounding for post-training quantization. In *International* 500  
*Conference on Machine Learning*, pages 7197–7206. PMLR, 501  
2020. 1 502
- [23] Eunhyeok Park, Sungjoo Yoo, and Peter Vajda. Value-aware 503  
quantization for training and inference of neural networks. 504  
In *Proceedings of the European Conference on Computer Vi-* 505  
*sion (ECCV)*, pages 580–595, 2018. 2 506
- [24] Charles R. Qi, Hao Su, Kaichun Mo, and Leonidas J. Guibas. 507  
Pointnet: Deep learning on point sets for 3d classification 508  
and segmentation. In *Proceedings of the IEEE Conference* 509  
*on Computer Vision and Pattern Recognition (CVPR)*, 2017. 510  
2 511
- [25] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J 512  
Guibas. Pointnet++: Deep hierarchical feature learning on 513  
point sets in a metric space. In *Advances in Neural Infor-* 514  
*mation Processing Systems*. Curran Associates, Inc., 2017. 515  
2 516
- [26] Guocheng Qian, Yuchen Li, Houwen Peng, Jinjie Mai, 517  
Hasan Hammoud, Mohamed Elhoseiny, and Bernard 518  
Ghanem. Pointnext: Revisiting pointnet++ with improved 519  
training and scaling strategies. In *Advances in Neural In-* 520  
*formation Processing Systems*, pages 23192–23204. Curran 521  
Associates, Inc., 2022. 2 522
- [27] Sheng Shen, Zhen Dong, Jiayu Ye, Linjian Ma, Zhewei Yao, 523  
Amir Gholami, Michael W Mahoney, and Kurt Keutzer. Q- 524  
bert: Hessian based ultra low precision quantization of bert. 525  
In *Proceedings of the AAAI Conference on Artificial Intelli-* 526  
*gence*, pages 8815–8821, 2020. 1 527
- [28] Shaoshuai Shi, Chaoxu Guo, Li Jiang, Zhe Wang, Jianping 528  
Shi, Xiaogang Wang, and Hongsheng Li. Pv-rcnn: Point- 529  
voxel feature set abstraction for 3d object detection. In *Pro-* 530  
*ceedings of the IEEE/CVF conference on computer vision* 531  
*and pattern recognition*, pages 10529–10538, 2020. 5 532

- 533 [29] Shaoshuai Shi, Zhe Wang, Jianping Shi, Xiaogang Wang,  
534 and Hongsheng Li. From points to parts: 3d object detection  
535 from point cloud with part-aware and part-aggregation net-  
536 work. *IEEE transactions on pattern analysis and machine*  
537 *intelligence*, 43(8):2647–2664, 2020. 5
- 538 [30] Shaoshuai Shi, Li Jiang, Jiajun Deng, Zhe Wang, Chaoxu  
539 Guo, Jianping Shi, Xiaogang Wang, and Hongsheng Li. Pv-  
540 rnn++: Point-voxel feature set abstraction with local vector  
541 representation for 3d object detection. *International Journal*  
542 *of Computer Vision*, 131(2):531–551, 2023. 5
- 543 [31] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien  
544 Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou,  
545 Yuning Chai, Benjamin Caine, Vijay Vasudevan, Wei Han,  
546 Jiquan Ngiam, Hang Zhao, Aleksei Timofeev, Scott Et-  
547 tinger, Maxim Krivokon, Amy Gao, Aditya Joshi, Yu Zhang,  
548 Jonathon Shlens, Zhifeng Chen, and Dragomir Anguelov.  
549 Scalability in perception for autonomous driving: Waymo  
550 open dataset. In *Proceedings of the IEEE/CVF Conference*  
551 *on Computer Vision and Pattern Recognition (CVPR)*, 2020.  
552 5
- 553 [32] Haiyang Wang, Chen Shi, Shaoshuai Shi, Meng Lei, Sen  
554 Wang, Di He, Bernt Schiele, and Liwei Wang. Dsvt: Dy-  
555 namic sparse voxel transformer with rotated sets. In *Pro-  
556 ceedings of the IEEE/CVF Conference on Computer Vision*  
557 *and Pattern Recognition*, pages 13520–13529, 2023. 5
- 558 [33] Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien  
559 Demouth, and Song Han. Smoothquant: Accurate and effi-  
560 cient post-training quantization for large language models.  
561 In *International Conference on Machine Learning*, pages  
562 38087–38099. PMLR, 2023. 1, 2, 3
- 563 [34] Hongyi Yao, Pu Li, Jian Cao, Xiangcheng Liu, Chen-  
564 ying Xie, and Bingzhang Wang. Rapq: Rescuing accuracy  
565 for power-of-two low-bit post-training quantization. *arXiv*  
566 *preprint arXiv:2204.12322*, 2022. 1
- 567 [35] Tianwei Yin, Xingyi Zhou, and Philipp Krahenbuhl. Center-  
568 based 3d object detection and tracking. In *Proceedings of*  
569 *the IEEE/CVF conference on computer vision and pattern*  
570 *recognition*, pages 11784–11793, 2021. 2, 5
- 571 [36] Zhihang Yuan, Chenhao Xue, Yiqi Chen, Qiang Wu, and  
572 Guangyu Sun. Ptq4vit: Post-training quantization for vision  
573 transformers with twin uniform quantization. In *European*  
574 *conference on computer vision*, pages 191–207. Springer,  
575 2022. 2
- 576 [37] Yifan Zhang, Zhen Dong, Huanrui Yang, Ming Lu, Cheng-  
577 Ching Tseng, Yuan Du, Kurt Keutzer, Li Du, and Shanghang  
578 Zhang. Qd-bev : Quantization-aware view-guided distilla-  
579 tion for multi-view 3d object detection. In *Proceedings of*  
580 *the IEEE/CVF International Conference on Computer Vision*  
581 *(ICCV)*, pages 3825–3835, 2023. 1
- 582 [38] Sifan Zhou, Liang Li, Xinyu Zhang, Bo Zhang, Shipeng  
583 Bai, Miao Sun, Ziyu Zhao, Xiaobo Lu, and Xiangxiang Chu.  
584 Lidar-ptq: Post-training quantization for point cloud 3d ob-  
585 ject detection. *arXiv preprint arXiv:2401.15865*, 2024. 1
- 586 [39] Yin Zhou and Oncel Tuzel. Voxelnet: End-to-end learning  
587 for point cloud based 3d object detection. In *Proceedings of*  
588 *the IEEE conference on computer vision and pattern recog-*  
589 *niton*, pages 4490–4499, 2018. 2

## A. Appendix

590

You may include other additional sections here.

591

Table 7. nuScenes Results for Different Detectors

Models	Methods	Bits(W/A)	mAP	NDS	Car	Truck	CV	Bus	Trail
	Full Prec.	32/32	59.22	66.48	84.86	57.38	16.85	70.75	38.1
CP-Voxel	SQ-0.80	8/8	59.16	66.40	84.69	57.31	16.89	70.70	38.1
	SQ-X.XX	4/8	x.x						

Table 8. KITTI Results for Different Detectors

Models	Methods	Bits(W/A)	Car	Pedestrian	Cyclist
	Full Prec.	32/32	78.62	52.97	67.14
	Max-min	8/8	78.23	52.96	62.01
SECOND	SQ-0.60	8/8	78.69	52.97	67.03
	Max-min	4/8	69.41	42.81	52.99
	SQ-0.65	4/8	78.29	54.72	64.03
	Full Prec.	32/32	77.28	52.30	62.71
	Max-min	8/8	74.68	50.83	60.44
PointPillar	SQ-0.70	8/8	76.79	51.96	62.84
	Max-min	4/8	63.74	44.15	55.50
	SQ-0.35	4/8	75.11	49.79	60.02