

# PointCompress3D - A Point Cloud Compression Framework for Roadside LiDARs in Intelligent Transportation Systems

Walter Zimmer<sup>★</sup> 

Ramandika Pranamulia<sup>★</sup> 

Xingcheng Zhou 

Mingyu Liu 

Alois C. Knoll 

<https://pointcompress3d.github.io>

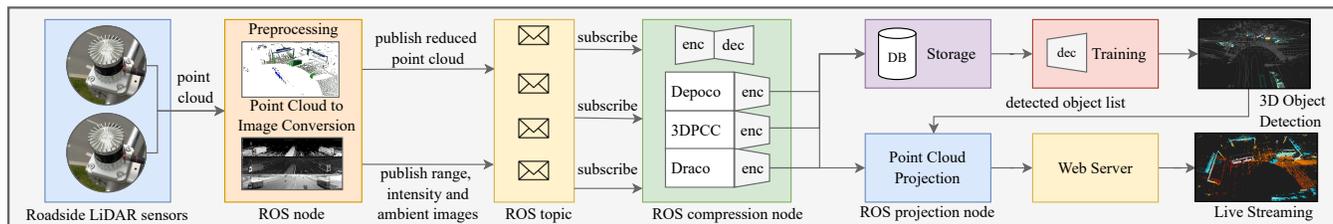


Fig. 1: Visualization of our point cloud compression and streaming pipeline, showcasing raw LiDAR point clouds on the left, progressing through preprocessing, compression, storage, training, and 3D object detection, merging in the streaming of compressed point clouds with 3D detection results.

**Abstract**—In the context of Intelligent Transportation Systems (ITS), efficient data compression is crucial for managing large-scale point cloud data acquired by roadside LiDAR sensors. The demand for efficient storage, streaming, and real-time object detection capabilities for point cloud data is substantial. This work introduces *PointCompress3D*, a novel point cloud compression framework tailored specifically for roadside LiDARs. Our framework addresses the challenges of compressing high-resolution point clouds while maintaining accuracy and compatibility with roadside LiDAR sensors. We adapt, extend, integrate, and evaluate three cutting-edge compression methods using our real-world-based *TUMTraf* dataset family. We achieve a frame rate of 10 FPS while keeping compression sizes below 105 Kb, a reduction of 50 times, and maintaining object detection performance on par with the original data. In extensive experiments and ablation studies, we finally achieved a *PSNR d2* of 94.46 and a *BPP* of 6.54 on our dataset. Future work includes the deployment on the live system. The code is available on our project website.

## I. INTRODUCTION

A project named *AUTOtech.agil* is currently being undertaken by the German government in collaboration with research institutions and universities in Germany, as well as the industry, to achieve a reliable mobility system. This project is expected to contribute significantly to various aspects of mobility, including creating a live digital twin of the traffic and achieving level 4 autonomous shuttles. One concept for developing ITS systems to create digital twins of traffic involves the installation of cameras and LiDAR sensors at the sides of roads and junctions. These sensors capture data

from the roads, which is then processed, and disseminated to all relevant parties. A LiDAR employs laser light for distance measurement, creates precise three-dimensional maps of the surroundings, and translates the environment into a collection of 3D points. LiDAR sensors operate at speeds ranging from 10-30 Hz and emit up to 2.6 Mio. points per second (pps). This results in a significant amount of data, ranging from 5-25 MB per frame. There are many use cases for roadside infrastructure LiDAR sensors. One of the most crucial tasks is to process point cloud data to detect traffic participants. Numerous scene understanding [14] and object detection methods [15]–[18] have been developed for that purpose, using e.g. *PointPillars* [19] as their baseline. In other cases, they can be incorporated with multiple input sources like images to perform multi-modal object detection [20]–[23]. However, the large amount of data poses challenges in ease of processing, storing, and transmitting.

### Our contribution is the following:

- We propose a point cloud compression framework shown in Fig. 1 for roadside infrastructure LiDARs.
- We provide an in-depth comparison and analysis of state-of-the-art compression methods on *SemanticKITTI* [24], *Ford* [25] and the *TUMTraf* datasets [26]–[29].
- We extend existing compression methods to make them compatible with our roadside *Ouster* LiDAR sensors.
- We perform extensive experiments and ablation studies on our real *TUMTraf Intersection* and *TUMTraf V2X Coop. Perception* dataset and achieve a *PSNR d2* of 94.46 and a *BPP* of 6.54 using our dev-kit<sup>1</sup>.
- We open-source our framework, which contains the point cloud projection and compression module and provide a project website with video results.

This research was supported by the Federal Ministry of Education and Research in Germany within the *AUTOtech.agil* project (01IS22088U).

The authors are with the School of Computation, Information and Technology, Technical University of Munich, 85748 Garching, Germany

Corresponding author: [walter.zimmer@tum.de](mailto:walter.zimmer@tum.de)

<sup>★</sup>These authors contributed equally.

<sup>1</sup><https://github.com/tum-traffic-dataset/tum-traffic-dataset-dev-kit>

TABLE I: Performance evaluation of state-of-the-art methods. We report the *BPP* metric for *PSNR* d2. and highlight the supported methods in gray. Best results are marked in bold, second best are underlined.

Method	Year	BPP↓	PSNR (d1) [db]↑	PSNR (d2) [db]↑	Enc. Speed [ms]↓	Dec. Speed [ms]↓	Code	Dataset
MPEG-GPCC [1]	2018	1.00 [2]	67.00 [2]	73.00 [2]	160.20 [3]	112.80 [3]	✓	SemanticKITTI
Draco [4]	2018	5.98 [5]	–	44.19 [5]	25.80 [3]	15.30 [3]	✓	-
OctSqueeze [6]	2020	1.00 [2]	69.00 [2]	74.00 [2]	35.75	92.96	✓	SemanticKITTI
MuSCLE [3]	2020	5.00	–	<u>81.00</u>	20.80	49.00	-	SemanticKITTI
VoxelContext-Net [7]	2021	1.00 [2]	<u>72.00</u> [2]	76.00 [2]	388.90 <sup>†</sup> [2]	374.10 <sup>†</sup> [2]	-	SemanticKITTI
Depoco [8]	2021	4.98 [5]	–	40.01 [5]	32.00 [5]	<b>2.00</b> [5]	✓	SemanticKITTI
OctAttention [9]	2022	1.00 [2]	<b>73.00</b> [2]	77.00 [2]	<b>2.30</b> <sup>‡</sup> [2]	2060.2 <sup>‡</sup> [2]	✓	SemanticKITTI
GrASPNet [10]	2022	1.00	53.00	57.00	193.80	345.92	✓	Ford
RIDDLE [11]	2022	5.00	–	<b>95.00</b>	532.51	966.30	–	SemanticKITTI
D-PCC [5]	2022	4.23	–	47.98	80.00	24.00	✓	SemanticKITTI
3DPCC [12]	2022	10.00	–	–	35.00*	5.00*	✓	–
OctFormer [2]	2023	1.00	<u>72.00</u>	76.00	5.20	7.30	✓	SemanticKITTI
InterFrame [13]	2024	<b>0.06</b>	<u>70.10</u>	73.00	364.00	714.00	–	8iVFB v2

<sup>†</sup>Using a voxel size of 5.

<sup>‡</sup>Using a context window size of N=1024.

\*Using JPEG2000 compression instead of RNN.

## II. RELATED WORK

**Point cloud data** consists of  $(x, y, z)$  geometric coordinates alongside additional attributes such as color, normal, intensity, and reflectivity. However, most point cloud compression methods still focus on geometry compression. Methods doing attribute compression in addition to geometry compression have done both compressions separately, indicating the importance of compressing the geometrical part exclusively. Geometry compression is highly prioritized because it is still difficult to get a good balance of fast compression speed, low bit rate, and low reconstruction error. Thus, achieving superior results for attribute compression is irrelevant if the performance in geometric compression is still lacking.

**Geometry compression** focuses on compressing the 3D coordinates, representing geometrical shapes, with various techniques. The first step usually involves converting 3D coordinates - which is very hard to compress because of its sparsity - into a dense representation, such as tree structures [1], [4], voxels [7], [30], or 2D images [12], [31]. After obtaining this dense representation, various methods can be applied depending on the previously chosen structure. In general, the processing methods will exploit the statistical properties of the data using hand-crafted or learning-based methods.

**Octree** - We can encode the *octree*'s occupancy nodes using predictive and arithmetic coding for the *octree* representation. Predictive coding is used to predict the occupancy symbol of the *octree* nodes based on their neighbors. This can be achieved by hand-crafted entropy models [1] or by using deep learning-based methods to automatically learn the entropy model [6], [9], [32]. Arithmetic coding is mostly used to encode the residual error from comparing the original data to the reconstructed data [33]. Unfortunately, this approach is highly complex when dealing with large point cloud data since constructing the *octree* and performing traversal operations can be time-consuming. A precision loss might also be a part of the *octree*-based codec because of the quantization process that took place. We might minimize *octree* complexity and the precision loss by limiting *octree*

depth and using a less aggressive quantization level respectively, but it results in a lower compression ratio.

### Algorithm 1 Voxelization Process

- 1: **Input:** Point cloud data  $\mathcal{P}$ , Voxel size  $v$
- 2: **Output:** Voxel grid representation  $\mathcal{V}$
- 3:
- 4: **Procedure:**
- 5: Define the resolution of the voxel grid by specifying the voxel size  $v$ .
- 6: Divide the 3D space occupied by the point cloud  $\mathcal{P}$  into a regular grid of voxels.
- 7: **for** each voxel  $voxel$  in  $\mathcal{V}$  **do**
- 8: Initialize voxel values based on assignment method:
- 9:     **Binary occupancy:** Set voxel value to 1 if it contains one or more points, 0 otherwise.
- 10:     **Averaged properties:** Calculate average properties (e.g., color, normal) of points within voxel.
- 11:     **Density:** Count the number of points within the voxel and store them as a value.
- 12: **end for**

**Voxelization** - Voxelizing the point cloud into 3D structural grids with specified grid sizes as shown in Table 1. is another common representation. This extends the concept of 2D images constructed by pixels, while in this case, the 3D space is made of *voxels*. This aims to bring the concepts of 2D processing into 3D space, one of which is a *3D CNN autoencoder* that we can employ to compress those *voxels* through the encoder part and decompress them through the decoder part [7]. Another method we can apply to those *voxel* representations is masked convolutions to predict *voxel* occupancy probabilities sequentially, feeding previously predicted *voxels* as context [34], which is inspired by combining *voxels* and *octree* representation. This voxelization method has a similar problem to that of *octree* methods, where the size of the *voxels* will determine a tradeoff between codec speed and compression ratio.

**3D Transformation** - As for projection transformation, where the transformation results are images, compression technology for images has advanced significantly. Hence,

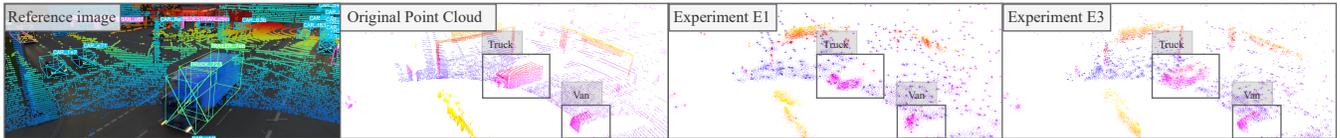


Fig. 2: Qualitative results on our *TUMTraf V2X Coop. Perception* [29] dataset. From left to right: a) Reference image b) View of the original point cloud. c) Experiment E1: Compressed point cloud with a subsampling distance of 3.0 and a minimum kernel radius of 1.5. d) Experiment E3: Compressed point cloud with a subsampling distance of 1.0 and a minimum kernel radius of 1.2. Both compressed point clouds are generated with max. 30,000 points and a grid size of [40, 40, 15] m.

we can utilize available technologies such as *JPEG* or deep learning-based image compression [35] methods to get a good result like in [12]. The advantage of this method is that the transformation directly reduces the size without the compression process. Each point cloud represented as  $(x, y, z)$  gets reduced to a single variable representing the range of the point from the observation source. It is also stated in [12] that they achieve a lossless compression, which means precisely reconstructing the original point is possible. Nevertheless, this approach is limited to a specific scenario. There should be no point occluded by another point. It also relies on the projection parameters, such as the position of the sensor and the azimuth angles, which are specific to the LiDAR type and manufacturer. Moving the sensor will cause the transformation to break.

### III. METHODOLOGY

We extend three available state-of-the-art methods (*Depoco*, *3DPCC* and *Draco*) to make them compatible with the *TUMTraf* datasets [36] and the Ouster LiDAR sensor. Then, we compare these methods and analyze their performance based on available metrics. Afterward, we choose the best methods, adapt them to roadside Ouster LiDAR sensors, and improve the compression. Finally, we integrate them into our *PointCompress3D* compression and streaming framework and connect them to a LiDAR 3D object detector to stream compressed point clouds with detected 3D objects.

#### A. Method Selection

To choose the most promising method from the state of the art, we first evaluate and inspect various aspects of each method, such as compression efficiency, compression effectiveness, compression speed, and code availability (see Table I).

The bitrate indicates the number of bits required per data point in a point cloud. It is influenced by factors like spatial information and additional data stored. Each bitrate corresponds to a unique *Peak Signal-to-Noise Ratio (PSNR)* value. Still, variations in the *PSNR* formula, particularly the *MAX* parameter, can affect comparisons between methods. Due to machine specifications and implementation details, speed is not directly comparable across methods. Additionally, some metrics mentioned may not be consistent or available in other sources.

The strategies employed to ensure comparability between methods include:

- We only focus on geometry compression, i.e., spatial attributes
- As *SemanticKITTI* is the most widely used dataset for point cloud compression, we focus on gathering metrics evaluated on this dataset.
- We only consider methods that use 1 or 5 *BPP* for each *PSNR* value evaluation.
- When we can't find the performance of a given method regarding a common metric, i.e., speed, we refer to other papers that provide the desired metric to provide cross-sourced metric findings.
- We only consider *bits per point* required to encode spatial attributes.

Although all the strategies above do not rule out the possibility of drawing incorrect conclusions about which methods are the best, at least the strategies mentioned can reduce the likelihood of drawing incorrect conclusions.

Looking at Table I, with priority given to speed and code availability before reconstruction error, we find four promising methods that run under 10 FPS. Although *MuSCLE* [3] fulfills the speed requirement, the code is not publicly available. Hence, we rule it out. Of all four candidates, *Depoco* is the fastest one, followed by *Draco*, with their total encoding speed being 34 ms and 40 ms, respectively. *Depoco* and *Draco* reconstruction errors are also comparable and not significantly different. Hence, we evaluate both *Draco* and *Depoco* alongside another additional method, *3DPCC-RNN* [12], on our dataset. *3DPCC-RNN* provides insights into converting 3D point clouds to range images and vice versa, enabling us to utilize image compression algorithms, which are actually more advanced and promising than directly compressing the 3D point cloud.

#### B. Extensions

To gain insight into the performance of a compression method with respect to our data, we perform several experiments on the *TUMTraf Intersection* and the *TUMTraf V2X Cooperative Perception* dataset. Some modifications and extensions to the original methods were implemented.

1) *Depoco*: *Depoco* is a learning-based compression method using autoencoders [8]. Configuration files and pre-trained weights of the *Depoco* model are provided.

We tune the parameters listed in Table II and start from the configuration file with the best qualitative result. Figure 2 shows how the initial configuration performs on our *TUMTraf Intersection* dataset. We can see that the configuration

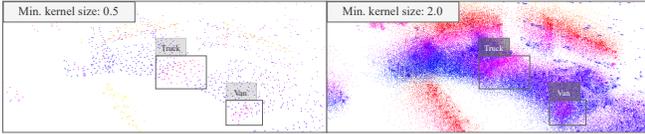


Fig. 3: Adjusting the *minimum kernel radius* parameters by 0.5 times and 2 times the original value, respectively.

with a subsampling distance of 1.0 and a minimum kernel radius of 1.2 performs best. Both compressed images are generated with max. 30,000 points and a grid size of [40, 40, 15] m. The truck and the car look almost identical as in the original point cloud frame. The mean encoding time for experiment E1 is 64 - 65 ms, and for experiment E3 is 110 - 120 ms, while the mean decoding time for both configurations is 90 ms. Therefore, we opted to start tuning from configuration E3, prioritizing encoding speed within the 0.1 s limit while ensuring superior quality performance.

We begin by tuning the following parameters: *subsampling distance*, *min. kernel radius*, and *kernel radius*. To maintain compatibility with the provided pre-trained weights, we retain the original *subsampling distance*. Adjustments to this parameter were found to cause block dimension mismatches, thus, it remains unchanged to utilize the existing weights effectively. Next, we adjust the *min. kernel radius* parameter. The result is depicted in Figure 3. We observe that enlarging the variable results in denser points while reducing it decreases the number of points. However, our current value appears optimal, as the object boundaries are well-defined and perceivable. Moving on, we tune the *kernel radius* in encoder and decoder blocks, with the default settings being 1.0 for the encoder and 0.05 for the decoder. Once again, examining the results in Figure 4, we find that the original parameter yields the best result. Subsequently, we adjust the *max. number of points*, similar to the *min. kernel radius* and *kernel radius*. Increasing the *max. number of points* adds more points, while decreasing it reduces the number of points. Notably, for *max. number of points*, the points added to the scene are around the dynamic objects, where the object remains perceptible but becomes densely populated with points, lacking inner details shown in Figure 5. We do not inspect the *grid feature dimension* because the grid size is the tile unit where compression is applied. This implies that smaller grid sizes should result in better compression quality. Moreover, the *grid feature dimension* does not affect the reconstruction quality. It is only modified if additional attributes need to be encoded, such as the intensity.

2) *3DPCC*: In [12], Beemelmans *et al.* presents a novel approach to convert 3D point cloud data acquired from *Velodyne* LiDARs to range images. The authors then employ *JPEG2000*, *PNG*, and *RNN*-based compression methods [35] to compress the images. We use *Ouster* LiDAR sensors, which require specific code modifications due to differences in settings compared to *Velodyne* LiDARs.

It’s important to note that *Velodyne* LiDARs point cloud data do not come with range images, necessitating a conver-

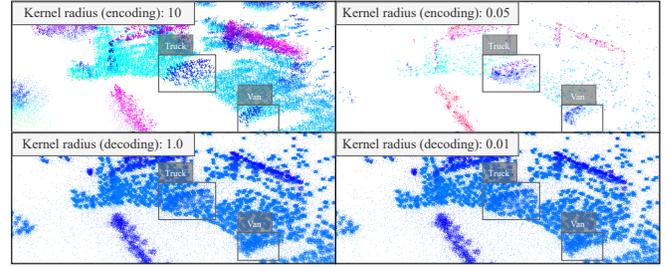


Fig. 4: We set the *kernel radius* in encoding block to 10 and 0.05 (top row) while keeping the decoding block values constant at 0.05. Subsequently, we set *kernel radius* in the decoding block to 1 and 0.01 (bottom row) while keeping the encoding block values constant at 1.0.

sion algorithm. In contrast, *Ouster* LiDARs provide range images alongside point cloud data. However, without azimuth maps, we are restricted to using the built-in range images directly. Some experiments that are going to be done involve replacing the point cloud-transformed range image with the built-in range image for the reconstruction, as can be seen in Tables III and IV.

The *3DPCC* compression framework is integrated into a *ROS* system. Communication occurs through *ROS* topics and messages. To utilize point clouds from *Ouster* LiDARs, we need to convert them into the correct point cloud data format. Later, instead of sending the images to another *ROS* topic, we store them in the filesystem. Three images representing one point cloud frame are produced: the range, intensity, and azimuth images. Converting back to the point cloud frame is accomplished by reading these images from the file system, performing a 2D to 3D transformation, and storing them back to the file system as *PCD* point cloud files.

Next, we need to adjust the hardcoded azimuth and elevation angles in the configuration file due to differences in sensor positioning and beam configuration<sup>2</sup>. In *Velodyne* sensors used in [12], the number of LiDAR beams is 32, where 16 are above the horizon, and the other 16 are below. In our *Ouster* sensor, there are 64 beams, all set up to point out under the horizon. We also need to order the elevation angle from lower to higher, and the azimuth angle should correspond to the right elevation angle representing a configuration of a laser beam.

In the last step, we modify the input point cloud. The initial order of points in our original point cloud files is sorted based on width, which is 2048 in our case. Therefore, the first 2048 rows correspond to one LiDAR beam closest to the horizon. The next 2048 rows correspond to the second LiDAR beam closest to the horizon. This is equivalent to all 64 beams. Hence, we have  $2048 * 64 = 131,072$  points in a single point cloud scan. We change the order of the points to columns: The first column represents the different 64 LiDAR scans, and then the second column represents another 64 LiDAR scans, up to 2048 columns scanned.

<sup>2</sup>The *Ouster* LiDAR mounted on our roadside infrastructure has a below horizon beam configuration.

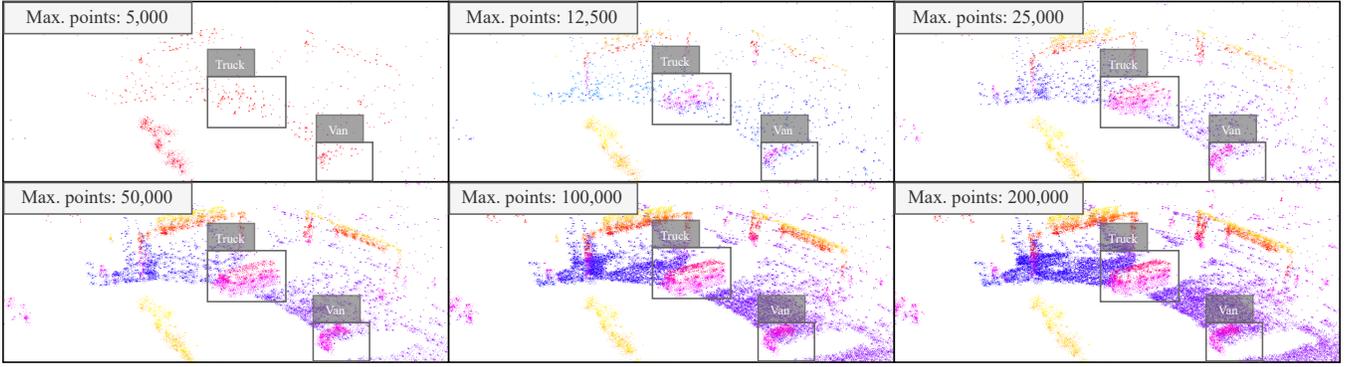


Fig. 5: Visualization of ablation study for the max. number of points. The point cloud shows a side view of a truck and a car with 5,000, 12,500, 25,000, 50,000, 100,000, and 200,000 points.

3) *Draco*: To use the *Draco* point cloud compression method in our framework, we integrate the *DracoPy* Python bindings [37], which wrap the *Draco* executables. We optimize the quantization bits and compression level parameters specifically for *Ouster* LiDAR data. By fine-tuning these parameters, we ensured an efficient and effective compression of point clouds generated from *Ouster* LiDARs to decrease the file storage size and improve transmission bandwidth for LiDAR-based applications.

#### IV. EVALUATION

##### A. Performance Metrics

In compression, we encounter diverse metrics assessing performance and quality. We balance speed, quality, and performance, yet lack a unified metric for comparison, unlike object detection. We typically assess reconstruction quality at specific compression rates, comparing metric graphs to determine method superiority.

1) *PSNR*: Following [38] based on [1], we use *PSNR* (Peak Signal-to-Noise Ratio) with the  $d1$  and  $d2$  metric, as reconstruction metric. *PSNR* is a widely used metric in image and video compression to evaluate the quality of the compressed image or video compared to the original ones. It measures the reconstruction quality by comparing the *mean squared error* (MSE) between the original and compressed signals. The *PSNR* metric is calculated based on normalized point cloud data, in a bidirectional way  $PSNR_{AB}$  and  $PSNR_{BA}$  and using the *max pooling* of these as our final *PSNR* value.

$$PSNR = 10 \cdot \log_{10} \left( \frac{MAX^2}{MSE} \right) \quad (1)$$

where  $MAX$  is the maximum possible pixel value of the image (255 for an 8-bit grayscale image) and  $MSE$  is the mean squared error between the original and compressed image.

*PSNR* is measured in decibels (dB), and a higher *PSNR* value indicates better quality because the  $MSE$  is smaller. In the case of point clouds, the  $MAX$  value can be set in various ways, depending on the compression type, such as voxelized or not [38]. The most common method is to set it

to the value of the maximum diagonal of the box bounding the point cloud.

2) *Chamfer Distance*: It is similar to *PSNR* in terms of a measurement system for reconstruction error. The *Chamfer* distance between two point clouds  $P_1 = \{\mathbf{x}_i \in \mathbb{R}^3\}_{i=1}^n$  and  $P_2 = \{\mathbf{x}_j \in \mathbb{R}^3\}_{j=1}^m$  is defined as the average distance between pairs of nearest neighbors between  $P_1$  and  $P_2$  i.e.

$$C(P_1, P_2) = \frac{1}{2n} \sum_{i=1}^n |\mathbf{x}_i - \text{NN}(\mathbf{x}_i, P_2)| + \frac{1}{2m} \sum_{j=1}^m |\mathbf{x}_j - \text{NN}(\mathbf{x}_j, P_1)| \quad (2)$$

where  $\text{NN}(\mathbf{x}, P) = \text{argmin}_{\mathbf{x}' \in P} \|\mathbf{x} - \mathbf{x}'\|$  is the nearest neighbor function and  $n, m$  are the number of points in  $P_1, P_2$  respectively. The *Chamfer* distance value has a minimum of 0, meaning that the two sets are spatially identical. We can also adjust the distance metric  $\|\cdot\|$  with various distance metrics such as *Manhattan* or *Euclidian* distance metric.

3) *BPP (bits per point)*: The compression rate is measured using *bits per point* (*BPP*), indicating the bits needed to store a single data point. For instance, a binary *pcd* file with 16-bit precision for  $(x, y, z)$  coordinates requires  $3 \times 16 = 48$  bits, which equals 6 bytes. In contrast, *ASCII* format uses 3 to 36 bytes due to its readable text storage (1 byte per character). Compression formats achieve significantly smaller sizes (refer to Table I).

4) *Metric Graph*: As we can see in Table I, for each *bits per point* value representing compression method efficiency, there will be a *PSNR* value indicating the reconstruction quality of that specific *BPP* from a certain compression method. We can control *BPP* indirectly in several ways, including quantization level, *octree* depth, or entropy coding efficiency. This will also indirectly affect the *PSNR* value. It is important to note that the *BPP* value is proportional to *PSNR*. When the *BPP* value is higher, *PSNR* will also be higher and vice versa. Therefore, we usually use a 2D graph to map the efficiency and effectiveness of compression methods and compare their performance using those graphs.

5) *Compression Speed*: We also measure the speed of compression and decompression individually to run each

TABLE II: Parameter tuning for *Depoco* on the max. number of points and the grid size. We use the *TUMTraf Intersection* dataset to find the best parameters and evaluate the compression method on the *TUMTraf V2X Cooperative Perception* dataset.

Grid Size	Max. Points	Enc. (ms)↓	Dec. (ms)↓	Enc. VRAM (GB)↓	Dec. VRAM (GB)↓	PSNR d1↑	PSNR (d2)↑	BPP↓	mAP <sub>3D</sub> ↑
8x8x3	50,000	<b>220</b>	2.7	<b>3.9</b>	<b>3.9</b>	15.32	23.68	<b>7.48</b>	13.32
	100,000	410	2.7	6.0	4.3	21.81	30.13	11.72	17.29
	200,000	520	2.8	5.5	7.1	<b>24.18</b>	<b>32.88</b>	13.57	<b>19.39</b>
16x16x6	50,000	230	2.6	<b>3.0</b>	<b>2.7</b>	-7.81	-0.19	<b>5.17</b>	12.21
	100,000	350	<b>2.5</b>	4.3	3.7	-3.59	3.87	7.03	19.50
	200,000	510	2.7	6.8	4.0	<b>-1.46</b>	<b>6.31</b>	7.73	<b>20.91</b>
24x24x9	50,000	240	<b>2.5</b>	<b>2.2</b>	<b>2.7</b>	-5.57	0.99	<b>3.90</b>	13.59
	100,000	410	2.6	5.5	3.1	-1.87	5.50	4.94	<b>19.75</b>
	200,000	530	2.7	6.6	3.2	<b>-0.18</b>	<b>7.47</b>	5.31	19.32

operation. The speed is calculated only for compression and decompression processes, neglecting the speed required to write or read data from storage.

6) *Memory Consumption*: Finally, we report the allocated memory on the GPU (VRAM) for each compression algorithm. We measure the memory consumption separately for the encoding and decoding steps.

7) *Mean Average Precision*: The *mean average precision* (mAP) metric quantifies the average precision scores across multiple object categories, evaluating a detection model’s performance across different classes. For our original uncompressed point cloud data we get an *mAP* of 20.11

### B. Experiment Setup

**Datasets** We use the *TUMTraf Intersection* dataset [27] (release R02, sequence S02) and the *TUMTraf V2X Cooperative Perception* [29] dataset. Both were labeled with the 3D bounding box annotation tool (3D BAT) [29], [39]. The latter one contains ten 10 s long sequences. Each sequence has 100 LiDAR point cloud scans and their corresponding 2D projection image (range, reflectivity, signal, near-infrared).

**Hardware Setup** Our framework runs inside a docker container with Docker version 20.10.20 under Ubuntu 20.04. We use 3x NVIDIA GeForce RTX 3090 GPUs, each having 24576 MiB of VRAM.

### C. Quantitative Results

In this section, we evaluate different methods on our *TUMTraf datasets* and provide quantitative results for *bits per point*, *PSNR d1* and *d2*, and *codec*<sup>3</sup> speed.

Referring to Sec. III-B.1, we find that the *max. number of points* and *grid size* are the most impactful parameters for *Depoco*. Tuning them further, we present the results

<sup>3</sup>The codec speed is given for the encoding and decoding process separately.

TABLE III: Point cloud image compression size (in KB).

Input Image	Size (gen.)	Size (orig.)	TinyJPG	PNG	ImgMagic
Range image	113.70	43.50	13.40	111.70	92.10
Azimuth image	56.30	-	10.20	65.30	56.30
Intensity image	<b>20.72</b>	-	<b>6.20</b>	<b>21.90</b>	<b>20.72</b>
<b>Total Size (KB)</b>	190.70	120.50	29.80	204.90	169.10

in Table II. The *decoding speed* remains consistent across configurations, while the *encoding speed*, *GPU usage*, and *BPP* increase with higher *max. number of points*. We calculate the *BPP* by dividing the encoded binary file size by the number of points in the original point cloud. We proceed by analyzing the conversion of 3D point clouds to 2D images. Each 3D point cloud frame yields three images (AIR): azimuth, intensity, and range. An example result is shown in Fig. 9. The transformation speed from 3D to 2D points is 300 ms, and from 2D to 3D points is 19 ms, while the *ImageMagick* PNG to JPG conversion takes approximately 15 ms. The average size of the range, azimuth, and intensity images from conversion compared to the original ones can be seen in Tab. III. Further, we compress the images using *PNG* compression using the *PIL* python library, online compression *TinyJPG* [40], and *ImageMagick* library with compression level 100 for the range image only. As we can see from the table, the *PNG* compression by *PIL* hardly compresses the images followed by the *ImageMagick*, while the *TinyJPG* performs really well. At last, we evaluate the *Draco* method, achieving encoding and decoding times of 60 ms and 80 ms, a *PSNR* of 52 dB and 144 dB, and an *mAP* of 20.01 and 20.11, respectively. Finally, we compare *Depoco*’s performance to *3DPCC* and *Draco*, visualizing it in a 2D graph (see Fig. 8).

### D. Qualitative Results

Here, we show the reconstruction result for different parameters of *Depoco* as shown in Tab. II. Furthermore, we provide qualitative results for the reconstruction from range images as depicted in Tab. IV. We illustrate how the *PSNR* value corresponds to the qualitative result. As shown in Fig. 6, increasing the *max. number of points* parameter will add

TABLE IV: Evaluation of the compression and the detection performance based on different input representations (3D to 2D transformation).

Input Representation	PSNR d1	PSNR d2	mAP@50
Generated RIA <sup>†</sup>	25.93	<b>36.23</b>	<b>7.49</b>
Generated IA + Ori. R <sup>†</sup>	21.57	24.87	0.00
Generated IA + Comp. R <sup>†</sup> (ImageMagick)	<b>26.38</b>	33.72	2.48
Comp. RIA <sup>†</sup> (TinyJPEG)	26.10	26.92	0.04

<sup>†</sup>A=azimuth, I=intensity, R= range image.

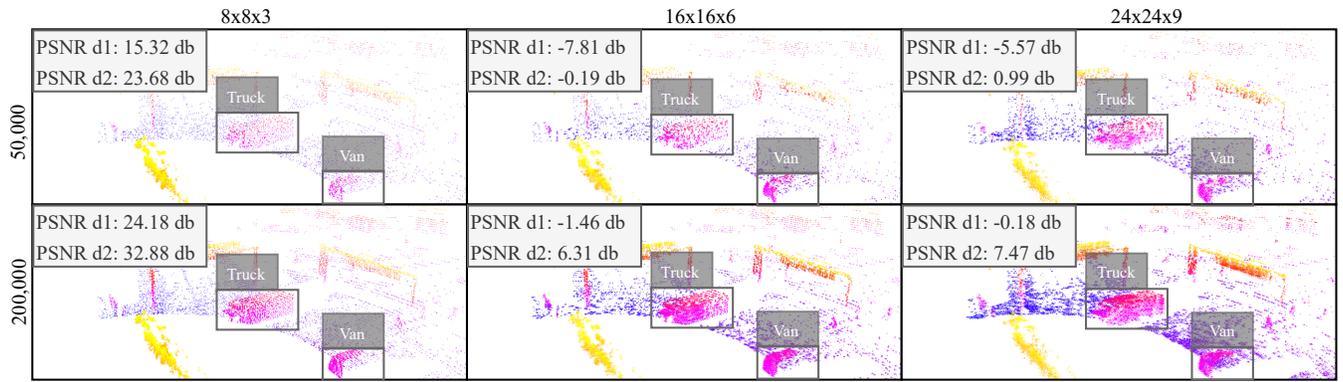


Fig. 6: We show the ablation study results for *Depoco* of reconstructed point clouds for different point and grid sizes on the *TUMTraF V2X Coop. Perception* [29] dataset.

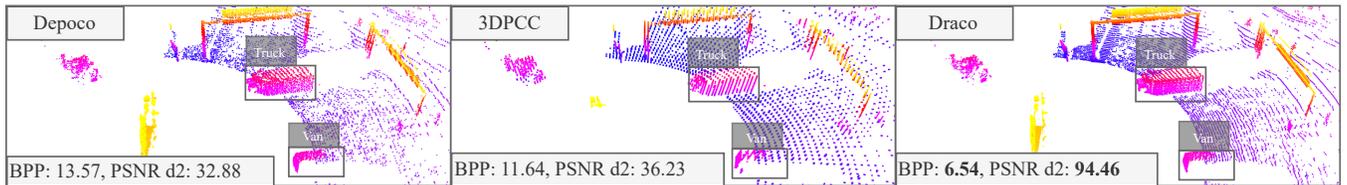


Fig. 7: Qualitative reconstruction results of *Depoco*, *3DPCC* and *Draco* on the *TUMTraF V2X Coop. Perception* [29] dataset.

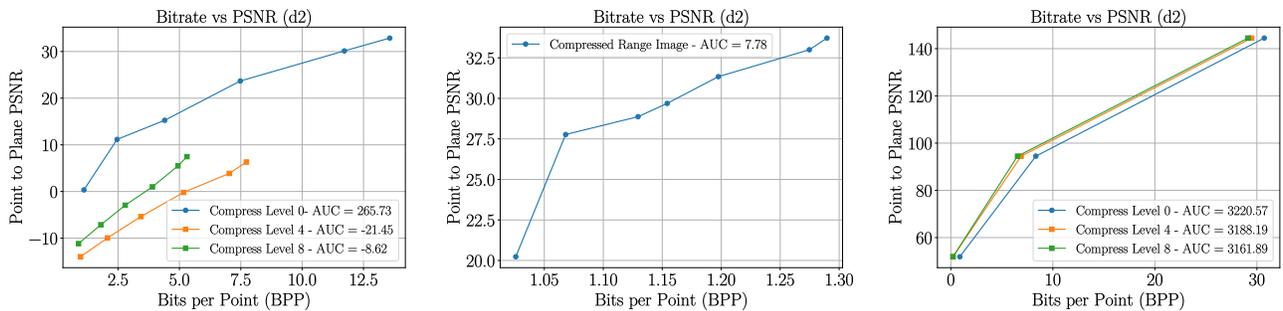


Fig. 8: From left to right: Visualization of *BPP* and *PSNR* metrics for *Depoco*, *3DPCC* and *Draco*.



Fig. 9: Result of the 3D to 2D transformation. From top to bottom: range, azimuth, and intensity image. We post-process the images and show the front part of the LiDAR.

more points to the scene after decoding. Note that we keep the grid size low, either  $8 \times 8 \times 3$  or  $16 \times 16 \times 6$ , to get a defined object. Otherwise, the object will be blurry, as we can see in the  $200k | 24 \times 24 \times 9$  setting, and it will get even more blurry as the grid size gets larger. Keeping the grid size small also keeps the *PSNR* higher, meaning that the point cloud reconstructed looks more similar to the original one. Fig. 7 illustrates the optimal reconstruction quality achieved at minimal bits per pixel for *Depoco*, *3DPCC*, and *Draco*.

## V. CONCLUSION & FUTURE WORKS

After evaluating the three methods, *Draco* emerges as the most promising method, with encoding and decoding speeds exceeding 10 FPS alongside good *PSNR* and *mAP* values. *Draco*'s *PSNR* and *mAP* values consistently correlate, with higher *PSNR* values leading to better *mAP* values. Quantization bits are revealed as the most influential parameters in *Draco*, with value of 16 yielding *mAP* values of 20.01 and 20.11 when set to the maximum of 30 bits. On the other hand, *Depoco* faces speed issues despite configurations surpassing the baseline *mAP* of 20.11. The surpassing result is due to *Depoco*'s introduction of additional points, enhancing object recognition but leading to inconsistent *PSNR*-*mAP* relationships. Lastly, *3DPCC* exhibits a good *PSNR* value without compression (36.23 vs. 32.88 for *Depoco*) but a very poor *mAP* of only 7.49, leaving no option for improvement except fixing the transformation to prevent loss of point clouds as what we have found. Future work includes deploying our framework on the live system, and the extension to other LiDAR types and manufacturers.

## ACKNOWLEDGMENT

This research was supported by the Federal Ministry of Education and Research in Germany within the project *AUTOtech.agil*, Grant Number: 01IS22088U.

## REFERENCES

- [1] S. Schwarz, M. Preda, V. Baroncini, M. Budagavi, P. Cesar, P. A. Chou, R. A. Cohen, M. Krivokuća, S. Lasserre, Z. Li, *et al.*, “Emerging mpeg standards for point cloud compression,” *IEEE Journal Emerging and Selected Topics in Circuits and Systems*, vol. 9, pp. 133–148, 2018.
- [2] M. Cui, J. Long, M. Feng, B. Li, and H. Kai, “Octformer: Efficient octree-based transformer for point cloud compression with local enhancement,” in *Proc. of AAAI Conf. on Artificial Intellig.*, vol. 37, no. 1, 2023, pp. 470–478.
- [3] S. Biswas, J. Liu, K. Wong, S. Wang, and R. Urtasun, “Muscle: Multi sweep compression of lidar using deep entropy models,” *Adv. in Neural Information Processing Systems*, vol. 33, pp. 22 170–22 181, 2020.
- [4] F. Galligan, M. Hemmer, O. Stava, F. Zhang, and J. Brettle, “Google/draco: a library for compressing and decompressing 3d geometric meshes and point clouds,” <https://github.com/google/draco>, 2018.
- [5] Y. He, X. Ren, D. Tang, Y. Zhang, X. Xue, and Y. Fu, “Density-preserving deep point cloud compression,” in *Proc. of IEEE/CVF Conf. on Computer Vision and Pattern Recognition*, 2022, pp. 2333–2342.
- [6] L. Huang, S. Wang, K. Wong, J. Liu, and R. Urtasun, “Octsqueeze: Octree-structured entropy model for lidar compression,” in *Proc. of IEEE/CVF conf. on computer vision and pattern recognition*, 2020.
- [7] Z. Que, G. Lu, and D. Xu, “Voxelcontext-net: An octree based framework for point cloud compression,” in *Proc. of IEEE/CVF Conf. on Computer Vision and Pattern Recognition*, 2021, pp. 6042–6051.
- [8] L. Wiesmann, A. Milioto, X. Chen, C. Stachniss, and J. Behley, “Deep compression for dense point cloud maps,” *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 2060–2067, 2021.
- [9] C. Fu, G. Li, R. Song, W. Gao, and S. Liu, “Octattention: Octree-based large-scale contexts model for point cloud compression,” in *Proc. of AAAI conf. on artificial intellig.*, vol. 36, no. 1, 2022, pp. 625–633.
- [10] J. Pang, M. A. Lodhi, and D. Tian, “Grasp-net: Geometric residual analysis and synthesis for point cloud compression,” in *Proc. of Workshop on Advances in Point Cloud Compression, Processing and Analysis*, 2022, pp. 11–19.
- [11] X. Zhou, C. R. Qi, Y. Zhou, and D. Anguelov, “Riddle: Lidar data compression with range image deep delta encoding,” in *Proc. of IEEE/CVF Conf. on Computer Vision and Pattern Recognition*, 2022, pp. 17 212–17 221.
- [12] T. Beemelmans, Y. Tao, B. Lampe, L. Reiher, R. van Kempen, T. Woopen, and L. Eckstein, “3d point cloud compression with recurrent neural network and image compression methods,” in *2022 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2022, pp. 345–351.
- [13] A. Akhtar, Z. Li, and G. Van der Auwera, “Inter-frame compression for dynamic point cloud geometry coding,” *IEEE Trans. on Image Processing*, 2024.
- [14] X. Zhou, M. Liu, B. L. Zagar, E. Yurtsever, and A. C. Knoll, “Vision language models in autonomous driving and intelligent transportation systems,” *arXiv preprint arXiv:2310.14414*, 2023.
- [15] W. Zimmer, J. Wu, X. Zhou, and A. C. Knoll, “Real-time and robust 3d object detection with roadside lidars,” in *Proc. of Int. Scientific Conf. on Mobility and Transport: Mobility Innovations for Growing Megacities*. Springer, 2023, pp. 199–219.
- [16] W. Zimmer, J. Birkner, M. Brucker, H. T. Nguyen, S. Petrovski, B. Wang, and A. C. Knoll, “Infradet3d: Multi-modal 3d object detection based on roadside infrastructure camera and lidar sensors,” in *2023 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2023.
- [17] W. Zimmer, M. Grabler, and A. Knoll, “Real-time and robust 3d object detection within road-side lidars using domain adaptation,” *arXiv preprint arXiv:2204.00132*, 2022.
- [18] W. Zimmer, E. Ercecik, X. Zhou, X. J. D. Ortiz, and A. Knoll, “A survey of robust 3d object detection methods in point clouds,” *arXiv preprint arXiv:2204.00106*, 2022.
- [19] A. H. Lang, S. Vora, H. Caesar, L. Zhou, J. Yang, and O. Beijbom, “Pointpillars: Fast encoders for object detection from point clouds,” in *Proc. of IEEE/CVF conf. on computer vision and pattern recognition*, 2019, pp. 12 697–12 705.
- [20] Y. Li, A. W. Yu, T. Meng, B. Caine, J. Ngiam, D. Peng, J. Shen, Y. Lu, D. Zhou, Q. V. Le, *et al.*, “Deepfusion: Lidar-camera deep fusion for multi-modal 3d object detection,” in *Proc. of IEEE/CVF Conf. on Computer Vision and Pattern Recognition*, 2022, pp. 17 182–17 191.
- [21] G. Melotti, C. Premebida, and N. Gonçalves, “Multimodal deep-learning for object recognition combining camera and lidar data,” in *2020 IEEE Int. Conf. on Autonomous Robot Systems and Competitions (ICARSC)*. IEEE, 2020, pp. 177–182.
- [22] A. Ghita, B. Antoniuussen, W. Zimmer, R. Greer, C. Creß, A. Møgelmoose, M. M. Trivedi, and A. C. Knoll, “Activeanno3d—an active learning framework for multi-modal 3d object detection,” in *2024 IEEE Intelligent Vehicles Symposium, IV 2024, Korea, June 2-6, 2024*. IEEE, 2024, p. 8.
- [23] A. Hekimoglu, P. Friedrich, W. Zimmer, M. Schmidt, A. Marcos-Ramiro, and A. Knoll, “Multi-task consistency for active learning,” in *Proc. of IEEE/CVF Int. Conf. on Comp. Vision*, 2023.
- [24] J. Behley, M. Garbade, A. Milioto, J. Quenzel, S. Behnke, C. Stachniss, and J. Gall, “Semantickitti: A dataset for semantic scene understanding of lidar sequences,” in *Proc. of IEEE/CVF int. conf. on computer vision*, 2019, pp. 9297–9307.
- [25] S. Agarwal, A. Vora, G. Pandey, W. Williams, H. Kourous, and J. McBride, “Ford multi-av seasonal dataset,” *The Int. Journal of Robotics Research*, vol. 39, no. 12, pp. 1367–1376, 2020.
- [26] C. Creß, W. Zimmer, L. Strand, M. Fortkord, S. Dai, V. Lakshminarasimhan, and A. Knoll, “A9-dataset: Multi-sensor infrastructure-based dataset for mobility research,” in *2022 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2022, pp. 965–970.
- [27] W. Zimmer, C. Creß, H. T. Nguyen, and A. C. Knoll, “Tumtraf intersection dataset: All you need for urban 3d camera-lidar roadside perception,” in *2023 IEEE 26th Int. Conf. on Intelligent Transportation Systems (ITS)*. IEEE, 2023, pp. 1030–1037.
- [28] C. Creß, W. Zimmer, N. Purschke, B. N. Doan, S. Kirchner, V. Lakshminarasimhan, L. Strand, and A. C. Knoll, “Tumtraf event: Calibration and fusion resulting in a dataset for roadside event-based and rgb cameras,” *IEEE Transactions on Intelligent Vehicles*, pp. 1–19, 2024.
- [29] W. Zimmer, G. A. Wardana, S. Sriharan, X. Zhou, R. Song, and A. C. Knoll, “Tumtraf v2x cooperative perception dataset,” in *Proc. of IEEE/CVF Conf. on Comp. Vision and Pattern Recog.*, 2024, p. 10.
- [30] J. Pang, K. Bui, and D. Tian, “Pivot-net: Heterogeneous point-voxel-tree-based framework for point cloud compression,” *arXiv preprint arXiv:2402.07243*, 2024.
- [31] H. Houshiar and A. Nüchter, “3d point cloud compression using conventional image compression for efficient data transmission,” in *2015 XXV int. conf. on information, communication and automation technologies (ICAT)*. IEEE, 2015, pp. 1–8.
- [32] M. Tatarchenko, A. Dosovitskiy, and T. Brox, “Octree generating networks: Efficient convolutional architectures for high-resolution 3d outputs,” in *Proc. of IEEE int. conf. on computer vision*, 2017.
- [33] M. Quach, J. Pang, D. Tian, G. Valenzise, and F. Dufaux, “Survey on deep learning-based point cloud compression,” *Frontiers in Signal Processing*, vol. 2, p. 846972, 2022.
- [34] D. T. Nguyen, M. Quach, G. Valenzise, and P. Duhamel, “Learning-based lossless compression of 3d point cloud geometry,” in *ICASSP 2021-2021 IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 4220–4224.
- [35] G. Toderici, D. Vincent, N. Johnston, S. Jin Hwang, D. Minnen, J. Shor, and M. Covell, “Full resolution image compression with recurrent neural networks,” in *Proc. of IEEE conf. on Computer Vision and Pattern Recognition*, 2017, pp. 5306–5314.
- [36] M. Liu, E. Yurtsever, J. Fossaert, X. Zhou, W. Zimmer, Y. Cui, B. L. Zagar, and A. C. Knoll, “A survey on autonomous driving datasets: Statistics, annotation quality, and a future outlook,” *IEEE Transactions on Intelligent Vehicles*, pp. 1–29, 2024.
- [37] S. Lab, “DracoPy: Python bindings for Draco,” <https://github.com/seung-lab/DracoPy>, Accessed: 2024-05-01.
- [38] A. Javaheri, C. Brites, F. Pereira, and J. Ascenso, “Improving psnr-based quality metrics performance for point cloud geometry,” in *Int. Conf. on Image Processing (ICIP)*. IEEE, 2020, pp. 3438–3442.
- [39] W. Zimmer, A. Rangesh, and M. Trivedi, “3d bat: A semi-automatic, web-based 3d annotation toolbox for full-surround, multi-modal data streams,” in *2019 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2019, pp. 1816–1821.
- [40] Tinyjpg - compress jpeg images intelligently. [Online]. Available: <https://tinyjpg.com/>