

# SYNMOTOR: A BENCHMARK SUITE FOR OBJECT ATTRIBUTE REGRESSION AND MULTI-TASK LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

In this paper, we develop a novel benchmark suite including both a 2D synthetic image dataset and a 3D synthetic point cloud dataset. Our work is a sub-task in the framework of a remanufacturing project, in which small electric motors are used as fundamental objects. Apart from the given detection, classification, and segmentation annotations, the key objects also have multiple learnable attributes with ground truth provided. This benchmark can be used for computer vision tasks including 2D/3D detection, classification, segmentation, and multi-attribute learning. It is worth mentioning that most attributes of the motors are quantified as continuously variable rather than binary, which makes our benchmark well suited for the less explored regression tasks. In addition, appropriate evaluation metrics are adopted or developed for each task and promising baseline results are provided. We hope this benchmark can stimulate more research efforts on the sub-domain of object attribute learning and multi-task learning in the future.

## 1 INTRODUCTION

Machine learning researchers have developed tremendous inventive network models and algorithms during the past decade. In parallel, a relatively small number of benchmarks have been developed for evaluating and comparing the performance of various models. Datasets and benchmarks play an important role in the development of neural networks. A good dataset can boost the development of a certain computer vision domain, *e.g.*, ImageNet Deng et al. (2009) to image classification, PASCAL VOC Everingham et al. (2009) and COCO Lin et al. (2014) to image detection and segmentation, KITTI Geiger et al. (2012) to point cloud detection, or ShapeNet Chang et al. (2015) and S3IDS Armeni et al. (2016) to point cloud segmentation. However, among all the computer vision tasks, the task of attribute regression is less explored due to the scarcity of suitable benchmarks. Current attribute learning methods mostly focus on outdoor pedestrians Li et al. (2015; 2018) or human facials Sarafianos et al. (2018); Kalayeh et al. (2017). The attributes in those dataset are mostly binary (*e.g.* gender, with glasses or not), only few attributes are continuous variables (*e.g.* age). On the other hand, regression models with neural networks are mostly used for non-vision data. We think it would be interesting if we could contribute in bridging the gap between attribute learning and regression for the computer vision community.

In this paper, we propose SynMotor, a benchmark that gives the possibility to perform object attribute regression. Multi-task learning and multi-modal learning are also possible with the provided dataset. The benchmark originates from a sub-task within our manufacturing project which aims at the automatic disassembly of small electric motors. A mesh model dataset is firstly generated with a carefully developed Blender addon, in which the motor specifications are saved as object attributes. Subsequently, synthetic motor datasets of both 2D image dataset and 3D point cloud dataset are created for deep learning purposes. Apart from the object attributes ground truth, the ground truth label of common computer vision tasks including detection, classification and segmentation are also provided. Those labels are generated automatically along with the image or point cloud generation, no manual annotation required. On the other hand, the metrics for common computer vision tasks are already well-developed, while developing metrics for object attribute regression is more difficult since the metrics need to be attribute-oriented. In our benchmark, we have designed several metrics for attribute regression given the motor attributes. Baseline results are also provided.

The remainder of this paper is structured as follows: Section 2 summarizes the state-of-the-art of 3D dataset creation, multi-task learning, and attribute learning. Section 3 shows a pipeline of creating our synthetic dataset. Section 4 describes the tasks and the developed metrics. Section 5 gives baseline results with some widely recognized network models. Finally, Section 6 summarizes presented results and discusses future work.

## 2 RELATED WORK

**3D dataset.** For 3D object dataset, ModelNet Wu et al. (2015) builds a huge dataset of 3D CAD models and provides the ModelNet40 benchmark for tasks including 3D shape classification and retrieval. Another similar work of ShapeNet Chang et al. (2015) provides more detailed semantic annotations, its subsequent work of PartNet Mo et al. (2019) additionally offers fine-grained semantic segmentation information for a subset of the models. Tremblay et al. (2018) creates a dataset for object pose estimation. A large dataset of 3D-printing models is provided in Thingi10K Zhou & Jacobson (2016), while a more recent ABC dataset Koch et al. (2019) collects over 1 million CAD models including many mechanical components. Regarding 3D scenes, KITTI Geiger et al. (2012) uses a vehicle-mounted platform outfitted with a variety of sensors to capture and record road information. Data for autonomous driving tasks like 3D object detection and 3D tracking are collected. While Ros et al. (2016) and Khan et al. (2019) generate synthetic datasets for the segmentation and detection of objects in virtual urban scenes, Le Hoang-An et al. (2021) generates images from virtual garden scenes. SynthCity Griffiths & Boehm (2019) generates point clouds of urban scenes using Blender. Pierdicca et al. (2019) also uses Blender but for the generation of point clouds of historical objects. For indoor scenes, SUN-RGBD Song et al. (2015), S3DIS Armeni et al. (2016), and ScanNet Dai et al. (2017) use different cameras to scan rooms to get 3D indoor point cloud dataset.

**Multi-task learning.** Overfeat Sermanet et al. (2014) trains on classification, localization, and detection tasks simultaneously using a single shared convolutional network. Eigen & Fergus (2015) implements depth prediction, normal estimation and semantic labeling using a single multiscale convolutional network. Kendall et al. (2018) optimizes the weight of loss during multi-task training according to the uncertainty of each task, thereby realizing the simultaneous learning of classification and regression tasks of different orders of magnitude. MultiNet Teichmann et al. (2018) and UberNet Kokkinos (2017) offer methods for jointly performing classification, detection and semantic segmentation using a unified architecture and achieve very efficient results. The above approaches work on the fundamental principle of a global feature extractor comprised of convolutional layers shared by all tasks and a different output branch for each task. In contrast to previous approaches, the strategy used in Misra et al. (2016) and Gao et al. (2019) is to have a separate network for each task while exchanging parameters between their parallel layers. Ruder Ruder et al. (2019) divides each layer into task blocks and information sharing blocks of previous layers, and the input of each layer is a linear combination of these two blocks, so that when learning, you can choose to focus more on the relevance of the task or the task itself. PAD-Net Xu et al. (2018) uses multi-task learning making preliminary predictions for depth, scene, surface normal estimation, then combining these predictions to obtain the refined depth and scene parsing results.

**Attribute learning.** Attribute learning refers to the process of discovering relevant attributes based on known logical principles. Visual attribute recognition has become an important research area due to its high-level semantic information. Previous work like Russakovsky & Fei-Fei (2010) establishes transfer learning by learning visual attributes such as colors, shapes or textures of images, and makes connections between semantically unrelated categories. Attribute learning has a wide range applications in the recognition of pedestrians. Li et al. (2015) suggests two deep learning models for surveillance scenarios to recognize based on single attribute and multiple attributes respectively. PGDM Li et al. (2018) analyses the salient features of the pedestrian’s body structure, it is helpful for the recognition of pedestrian attributes. DIAA Sarafianos et al. (2018) utilizes multi-scale visual attention and weighted focal loss for person attribute recognition. With the widespread popularity of posting selfies in social software, the practice of recognizing facial attributes is growing, Liu et al. (2015) provide a cascaded deep learning framework for jointly predicting attributes and localizing faces. Kalayeh et al. (2017) leverages the information obtained from semantic segmentation to improve facial feature prediction, Liu et al. (2018) obtains equivalent discriminative ability for face attribute recognition on CelebA and LFWA Huang et al. (2008) datasets by learning disentangled but complementing face features with minimal supervision.

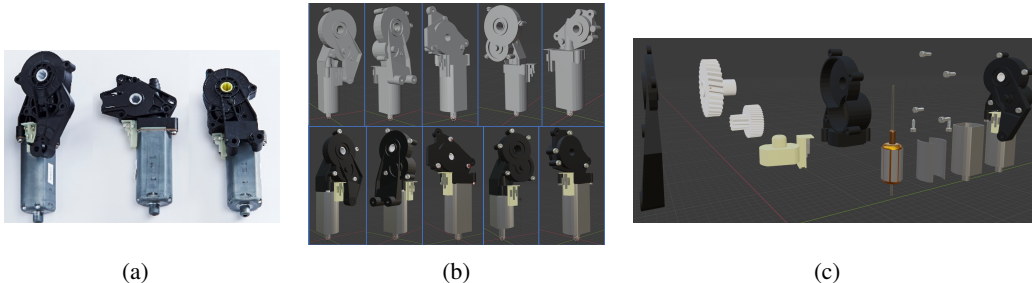


Figure 1: Real-world motors and synthetic motors. (b) Upper row: no textures added; bottom row: textures added and rendered. (c) An explosion figure of a generated motor. The original assembled model is also shown at the right most.

### 3 SYNTHETIC DATASET GENERATION

#### 3.1 MESH MODEL GENERATION

A Blender addon is created for the easy generation of synthetic motor mesh models. As an open source software, Blender Bonatti et al. (2016) is a proven tool that performs well in modeling shapes and creating highly customizable addons. Our addon is able to generate motor mesh models with various specifications and save them in desired file formats. Each component of a generated motor can also be saved separately. The generated models contain following components: (i) Pole Pot; (ii) Electric Connection; (iii) Gear Container; (iv) Cover; (v) Side Screws and (vi) Cover screws. Those are the six main categories we need perform segmentation on for the first step of disassembly. Additionally, following inner components have also been generated : (vii) Magnets; (viii) Armature; (ix) Lower Gear and (x) Upper Gear, as presented in Figure 1(c). To generate motors with various specifications, we provide lots of parameter options that control the type, size, position and rotation of different parts of motor, *e.g.* screw position, gear size, or pole pot length. Figure 1(b) shows ten generated demo motors with different parameters and an exploded view of a demo motor. All the individual components mentioned above are modeled separately as illustrated.

#### 3.2 IMAGE AND POINT CLOUD GENERATION

The generated mesh models are further used to create synthetic image and point cloud datasets. A simulated scene is built in Blender for it. Apart from the lights and cameras, to make the scene more realistic, a model of the real-world clamping system and a background panel have been added additionally. The camera rotates randomly on top of the scene within a certain view range yet always towards the motor. To create image dataset, apart from the scene images rendered by Blender directly, BlenderProc Denninger et al. (2019) is used to generate corresponding depth images, normal images, and segmentation ground truth images. Detection ground truth of bounding boxes of motor and screws are also provided. On the 3D synthetic data side, BlenSor Gschwandtner et al. (2011) is used to simulate the sensors to create point cloud data as well as to generate their segmentation ground truth. 3D bounding boxes are also given. Moreover, for the better learning of key objects, we additionally provide corresponding sub-point clouds for each scene by cropping random cuboid regions around the motors, but make sure to include them. The current version of our benchmark focuses more on the undismantled motors, hence inner components will not be investigated. However, any researcher interested in this part is free to use the provided addon and scripts to generate corresponding dataset according to their own needs.

#### 3.3 DATASET DETAILS

We have created a synthetic motor mesh dataset of 1000 motor mesh models with different specifications. They are placed in the same simulated scene but with random camera settings, random light conditions and random mild translations and rotations. With these 1000 scenes, 1000 sets of images and 1000 point clouds are generated respectively. 80% of the data are randomly selected as the training data, while the other 20% are used as the test data. Note that to ensure the correspon-

Table 1: Object attributes and their notations and ranges. The x/y/z values of a key point location are treated as three separate attributes. Some of the attributes are only valid under certain conditions. The attributes are used in different designed metrics. SRE stands for size relative error, GLE stands for gear location error, MRE stands for motor rotation error, SLE stands for screw location error. See a detailed explanation regarding these metrics in Section 4.2.

Notation	Attribute	Range	Validity	Involved metric
$T$	Motor type	{0, 1, 2, 3, 4}	always	Cls. accuracy
$N_s$	Number of cover screws	{3, 4, 5}	always	Cls. accuracy
$L_b$	Bottom length	6.2 ~ 8.0	always	SRE
$L_{sb}$	Sub-bottom length	0.6 ~ 2.0	always	SRE
$D_{lg}$	Lower gear region diameter	3.5 ~ 4.5	always	SRE
$D_{ug}$	Upper gear region diameter	5.0 ~ 6.5	if $T = 0/1/2$	SRE
$X_{lg}, Y_{lg}, Z_{lg}$	Lower gear center location (xyz)	(1.7 ~ 2.3, 1.0, 10.6 ~ 14.2)	always	GLE <sub>xyz</sub> , GLE <sub>xz</sub>
$X_{ug}, Y_{ug}, Z_{ug}$	Upper gear center location (xyz)	(1.7 ~ 2.5, 0.3 ~ 0.5, 13.5 ~ 17.3)	if $T = 0/1/2$	GLE <sub>xyz</sub> , GLE <sub>xz</sub>
$R_x, R_y, R_z$	Motor rotation (rx/ry/rz)	( $\pm 15^\circ, \pm 5^\circ, \pm 5^\circ$ )	always	MRE
$X_{s1}, Y_{s1}, Z_{s1}$	Cover screw 1 location (xyz)		always	SLE <sub>xyz</sub> , SLE <sub>xz</sub>
$X_{s2}, Y_{s2}, Z_{s2}$	Cover screw 2 location (xyz)	(-4.9 ~ 5.0,	always	SLE <sub>xyz</sub> , SLE <sub>xz</sub>
$X_{s3}, Y_{s3}, Z_{s3}$	Cover screw 3 location (xyz)	-3 ~ -1.4,	always	SLE <sub>xyz</sub> , SLE <sub>xz</sub>
$X_{s4}, Y_{s4}, Z_{s4}$	Cover screw 4 location (xyz)	8.6 ~ 20.7)	if $N_s = 4/5$	SLE <sub>xyz</sub> , SLE <sub>xz</sub>
$X_{s5}, Y_{s5}, Z_{s5}$	Cover screw 5 location (xyz)		if $N_s = 5$	SLE <sub>xyz</sub> , SLE <sub>xz</sub>

dence between images and point clouds for each scene, the camera information have been saved and shared between BlenderProc and Blesor. The dataset are organized as follows.

**Mesh model dataset.** Part obj files and assembled full shape obj file are both saved, with corresponding material files. Apart from that, all attributes are saved in a csv file. We select 30 key attributes as the main learning targets for the benchmark and have pre-processed them for a more convenient usage. The 30 attributes are given in Table 1. Attribute *Type* and *Number of cover screws* are suitable for classification tasks, while other 28 attributes are suitable for attribute learning. Note that some attributes are invalid in some cases, e.g., the xyz coordinates of 5th cover screw when the motor only has 4 cover screws. While the attribute values are set to zero when they are invalid, an additional binary mask csv file is provided for all 28 attributes of all 1000 motors. In our benchmark, all 28 attributes are continuous variables hence are suitable for attribute regression.

**Image dataset.** 1000 sets of images are provided. In each set, there are one rgb image, one depth image, one normal image and one segmentation ground truth image. Detection ground truth of 2D bounding boxes are provided in the COCO fashion, with a visualized detection result supplemented. A demo of generated images are given in Figure 2.

**Point cloud dataset.** 1000 sets of point clouds are provided. In each set, there are one scene point cloud and one cropped point cloud which focuses on the motor region. 3D bounding boxes are provided. Segmentation ground truth labels are also provided. Figure 3 gives a demo of generated point clouds colored in their segmentation ground truth.

## 4 DESIGNED TASKS AND METRICS

### 4.1 COMMON COMPUTER VISION TASKS

**Detection.** The detection task is to detect the positions of the key objects in the scene and to draw 2D or 3D bounding boxes around each object of interest in RGB images or point clouds. For 2D detection, the widely used mean average precision (mAP) was originally proposed in the VOC challenge Everingham et al. (2009). We use an advanced version in COCO Lin et al. (2014) which further considers different IoU thresholds. The two metrics are (i) mAP with IoU threshold of 0.5; (ii) average mAP with IoU threshold of 0.5, 0.55, 0.6, . . . , 0.95. For 3D detection, we require a 3D bounding box overlap of 70% for computing the precision-recall curve and the mAP.

**Classification.** The task of classification is to classify the key object in the scene into the prior-defined categories. There are 5 types of motors in our dataset, the main differences between them are the number of gears and the shape of covers. The classification task can be performed on both 2D dataset and 3D dataset. The well-known classification accuracy is used as the metric for the

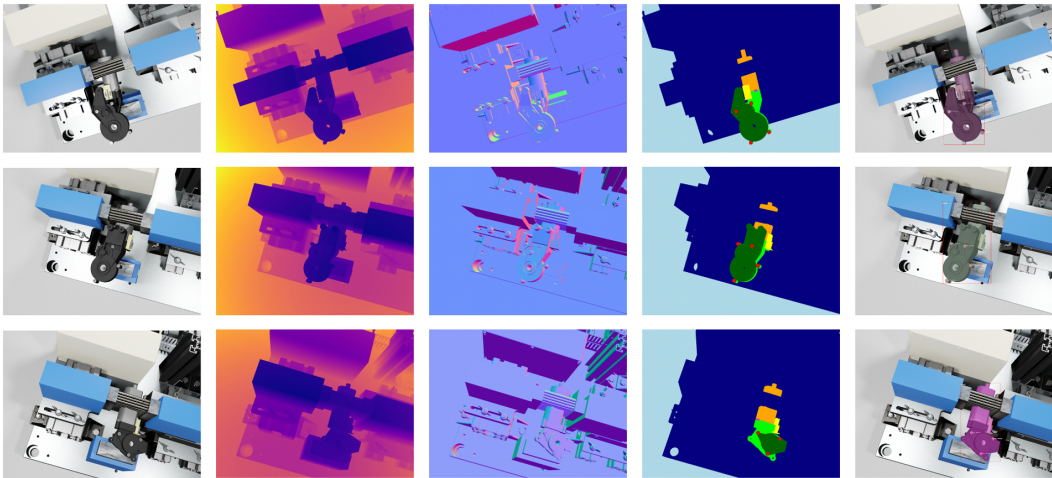


Figure 2: Demos of synthetic image data.

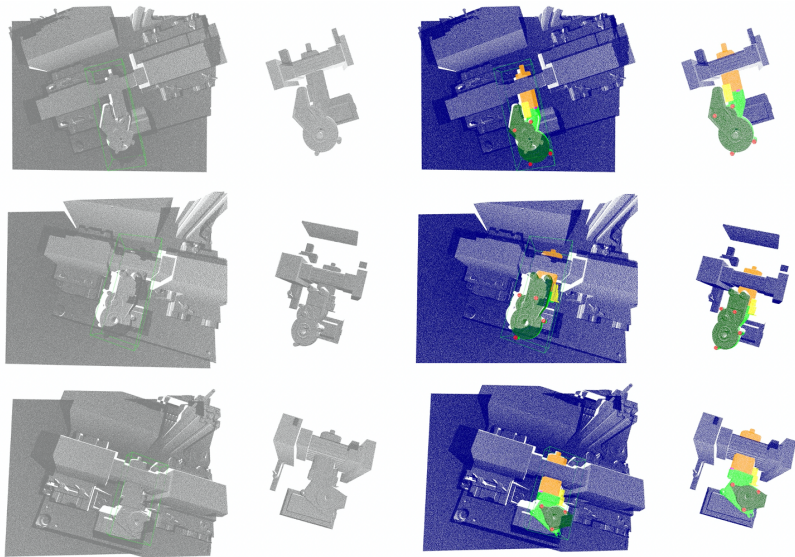


Figure 3: Demos of synthetic point cloud data.

classification task. This is a relatively easy task in our setting since we only have 5 categories for classification and none of them is a tail category.

**Segmentation.** Segmentation in 2D image is the process of assigning a label to every pixel in the image such that pixels with the same label share certain characteristics. 3D point cloud segmentation is the process of classifying point clouds into different regions, so that the points in the same isolated region have similar properties. Common segmentation metrics are category-wise IoU and overall mIoU. In our case, the overall mIoU metric is used. Moreover, since the screw categories are the key categories in real-world applications, an additional metric of screw mIoU is performed.

#### 4.2 OBJECT ATTRIBUTE REGRESSION

When generating the motor mesh models, their detailed geometric attributes are also saved, including the length of the pole pot, the diameter of the gear region, the positions of the screws, etc., which provides the possibility for regression learning. Evaluation metrics for assessing these learned attributes are proposed with consideration on multi-perspectives, including object size, object orientation, and object key point positions. It is worth noting that since the attributes are not always valid in all scenes, a binary mask is used for the metrics to get more accurate error information. In the

following definition of each metric, for a more clearly description, we use  $A$  to denote the union set of involved attributes, with a corresponding binary mask  $M_a$ . Note that the following metrics are only used for evaluation, the loss used for the training is computed with a masked mean square error (MSE) between the predicted results and the ground truth.

**Size Relative Error (SRE).** This metric evaluates the predicted overall motor size using four key attributes which represent the main body of motors: the lengths of the bottom and the sub-bottom part, *i.e.* pole pot, and the diameters of the gear regions. Denote  $A$  as the union set of those attributes:  $A \in \{L_b, L_{sb}, D_{lg}, D_{ug}\}$  with a corresponding binary mask  $M_a$ , for a batch of  $N$  motor point clouds with ground truth  $A^*$ , the SRE metric is given as

$$\text{SRE} = \frac{1}{N} \sum_{i=1}^N \left( \frac{1}{\sum M_{a_i}} \sum_{A_i} \frac{|M_{a_i} A_i - A_i^*|}{1 - M_{a_i} + A_i^*} \right) \quad (1)$$

where  $\sum M_{a_i}$  counts the valid attribute number in  $i$ th motor and  $1 - M_{a_i}$  is served as a safety term to ensure the denominator is greater than 0 When  $A_i^* = 0$ . For a certain attribute of a certain motor that is invalid, its ground truth value and mask value are both 0, *i.e.*  $A_i^* = M_{a_i} = 0$ , its size relative error is computed as zero and this attribute will not be counted in the denominator.

Note that we compute a mean error over the instance level for SRE. We believe this can represent the metric information better. The following three metrics are performed over the batch level.

**Gear Location Error (GLE).** This metric evaluates the distance error between the predicted location of gear region center and its ground truth. Involved attributes are the center point coordinate values. Denote  $A_l$  and  $A_u$  as the union sets of those attributes:  $A_l \in \{X_{lg}, Y_{lg}, Z_{lg}\}$  and  $A_u \in \{X_{ug}, Y_{ug}, Z_{ug}\}$  with masks  $M_{al}$  and  $M_{au}$ , for a batch of  $N$  motor point clouds, the  $\text{GLE}_{xyz}$  metric is given as

$$\text{GLE}_{xyz} = \frac{\sum_{i=1}^N (\sqrt{\sum_{A_{li}} (M_{ali} A_{li} - A_{li}^*)^2} + \sqrt{\sum_{A_{ui}} (M_{aui} A_{ui} - A_{ui}^*)^2})}{\sum_{i=1}^N (\sum M_{ali}/3 + \sum M_{aui}/3)} \quad (2)$$

where  $M_{aui}/3$  means for the 3D coordinate of each center point, the mask should only be counted for one time. In real-world applications of disassembly, the small error along the motor normal direction (in our case, the Y axis) is sometimes irrelevant. We further provide another metric that only considers the distance error on the XZ plane with projected points. By redefining  $A_l \in \{X_{lg}, Z_{lg}\}$  and  $A_u \in \{X_{ug}, Z_{ug}\}$ , a similar metric of  $\text{GLE}_{xz}$  is given as

$$\text{GLE}_{xz} = \frac{\sum_{i=1}^N (\sqrt{\sum_{A_{li}} (M_{ali} A_{li} - A_{li}^*)^2} + \sqrt{\sum_{A_{ui}} (M_{aui} A_{ui} - A_{ui}^*)^2})}{\sum_{i=1}^N (\sum M_{ali}/2 + \sum M_{aui}/2)} \quad (3)$$

**Motor Rotation Error (MRE).** This metric evaluates the absolute motor rotation error. Involved attributes are the motor rotations along three axes. Denote  $A$  as the union set of those involved attributes:  $A \in \{R_x, R_y, R_z\}$ . Since the rotation attributes are always valid, the mask  $M_a$  is unnecessary for computation. The metric MRE is defined as

$$\text{MRE} = \frac{1}{3N} \sum_{i=1}^N \sum_{A_i} |A_i - A_i^*| \quad (4)$$

**Screw Location Error (SLE).** This metric evaluates the distance error between the predicted cover screw positions and their ground truth. Involved attributes are the screw position coordinate values. Denote  $A_j$  as the union set of those attributes:  $A_j \in \{X_{sj}, Y_{sj}, Z_{sj}\}$  with masks  $M_{aj}$  where  $j = 1, 2, 3, 4, 5$ . For a batch of  $N$  motor point clouds, the  $\text{SLE}_{xyz}$  metric is given as

$$\text{SLE}_{xyz} = \frac{\sum_{i=1}^N \sum_{j=1}^5 \sqrt{\sum_{A_{ji}} (M_{aji} A_{ji} - A_{ji}^*)^2}}{\sum_{i=1}^N \sum_{j=1}^5 \sum M_{aji}/3} \quad (5)$$

where  $M_{aji}/3$  means for the 3D coordinate of each screw, the mask should only be counted for one time. Same as GLE, a metric  $\text{SLE}_{xz}$  only considers the XZ plane is defined with  $A_j \in \{X_{sj}, Z_{sj}\}$ :

$$\text{SLE}_{xz} = \frac{\sum_{i=1}^N \sum_{j=1}^5 \sqrt{\sum_{A_{ji}} (M_{aji} A_{ji} - A_{ji}^*)^2}}{\sum_{i=1}^N \sum_{j=1}^5 \sum M_{aji}/2} \quad (6)$$

### 4.3 MULTI-TASK LEARNING

Multi-task learning Ruder (2017); Zhang & Yang (2017) means to solve multiple learning tasks simultaneously, while exploiting the commonalities and differences between tasks. This can improve the learning efficiency and prediction accuracy of task-specific models compared to training the models individually. While acceptable performance can usually be obtained by focusing on a single task, the information that might help on getting better performance are possibly ignored. Specifically, the information come from training signals for related tasks. By sharing representations between related tasks, the generalization ability of the model can be improved on the original task.

In our case, above tasks can be performed simultaneously with a same backbone network. For example, classification and segmentation can be performed at the same time. Or as mentioned in the last subsection, the classification, detection and segmentation results may be used as additional input for the regression task. Multi-modal learning is also possible if both 2D and 3D data are used.

## 5 BASELINE RESULTS

### 5.1 2D DETECTION

Since there is no such network architectures for 3D point cloud detection as widely recognized as YOLO series for 2D image detection, for the detection task, we give baseline results on the 2D image dataset. YOLO models of different sizes are used. All the experiments are performed with same parameter settings. The input image resolution is of 640, the batch size is 16. The optimizer SGD is used with a epoch number of 200. The learning rate starts with 0.01 and decays to 0.001 in the end with a linear decay. The mAP performance of all models are given in Table 2. It shows that the YOLO framework achieves remarkable performance on the detection task. The mAP performance also improves when a larger network model is used.

Table 2: 2D Detection baseline results.

Model	mAP (IoU 0.5-0.95)	mAP (IoU 0.5)
YOLOv5n	75.3	93.6
YOLOv5s	77.8	94.2
YOLOv5m	82.7	95.9
YOLOv5l	86.7	96.4
YOLOv5x	89.1	96.8

### 5.2 3D CLASSIFICATION AND SEGMENTATION

In the past five years, a variety of network models have been proposed for point cloud data. Multi-view based Lawin et al. (2017); Boulch et al. (2017) and volumetric-based methods Maturana & Scherer (2015); Jiang et al. (2018) are mostly used in the early years. Since the pioneer work of PointNet Qi et al. (2017a), point-based methods which include point-wise MLP Qi et al. (2017b), point convolution-based methods Wu et al. (2019); Thomas et al. (2019) and graph-based methods Wang et al. (2019); Chen et al. (2021); Liang et al. (2020), gradually became the main choice. Recent work even adapt the idea of transformer Vaswani et al. (2017) for point cloud learning, *e.g.*, PCT Guo et al. (2021) and PT Zhao et al. (2020); Engel et al. (2021). Among all those methods, PointNet++ Qi et al. (2017b) and DGCNN Wang et al. (2019) are recognized as two key work in the domain. In this paper, we use DGCNN as the network backbone to produce baseline results.

We use an architecture of three successive edgeConv blocks (key block in DGCNN) to learn point-wise features, a linear layer and a max pooling layer is applied subsequently to compute a global feature. For the classification task, the global feature is processed with several followed dense layers to get a classification prediction; while for the segmentation task, the global feature is repeated and concatenated with point-wise features to process through several other dense layers to predict point-wise labels. An illustrative figure is given in Figure 4. Note that although two tail blocks are different, the encoder block is identical. It is possible to use a shared encoder for both classification and segmentation tasks. In this case, by adding the losses together and performing gradient back

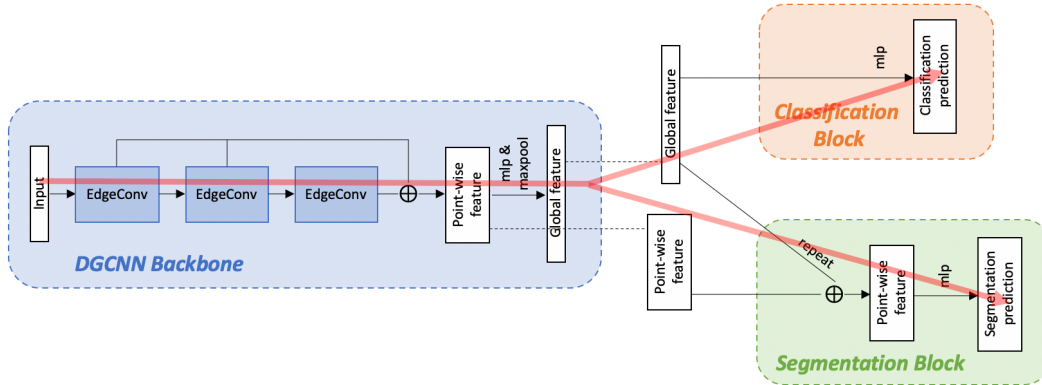


Figure 4: 3D classification and segmentation framework for baseline results.

Table 3: 3D Classification and segmentation baseline results.

Model	Classification $\uparrow$	Segmentation $\uparrow$	
	Accuracy	mIoU	Screw mIoU
DGCNN(cls)	100	-	-
DGCNN(seg)	-	93.17	83.05
DGCNN(cls+seg)	100	92.55	81.19

propagation at the same time, we can train on the classification task and the segmentation task simultaneously. This is a simple way of performing multi-task learning. The red line indicates how the information flows during one training step. The encoder block is co-trained with both tasks, while the respective task tail blocks are trained in parallel.

The settings of these three experiments are identical. We use a sub-point cloud size of 2048 points for batch training. The batch size is 16. The optimizer AdamW is used with a epoch number of 100. The learning rate starts with 0.001 and decays to 0.00001 with a *cos* decay. In the EdgeConv blocks, we use  $K = 32$  when selecting neighbor points. Table 3 indicates that the DGCNN global feature allows for a complete identification of the motor types since it is a relatively easy task. However, the segmentation results from the multi-task training model are not as good as that from individual training. Using loss weights to balance the gradient information from two tails blocks may improve the results. We would like to leave this problem for other researchers that are interested in this topic.

### 5.3 3D OBJECT ATTRIBUTE REGRESSION

Same as previous subsection, DGCNN is used as backbone for encoding point-wise features and a global feature. Moreover, as illustrated in Figure 5, we consider the classification-segmentation parallel training introduced in the previous subsection as a meta-block. Note that since the one-hot attributes  $T$  and  $N_s$  indicate the validity of some other attributes, apart from the block of motor type classification, another block for classifying the number of cover screws has also been included in the meta-block. The regression tasks compose a second tail block. In our case, we use a simple method of concatenating the global feature with one-hot features that have been encoded with MLPs, and then directly perform MLP on the enhanced global feature to get a 28-dimensional vector, which indicates the regression results of 28 attributes.

With the proposed architecture, there are several possible ways to train those blocks. (i) **Totally separate training**: the encoder and the meta-block are trained as one network first, and then another same-structure encoder and the regression block are trained as another totally separate network. No information is shared between two trainings. (ii) **Use meta-block for pre-training**: this method is similar to the last one, but the encoder weights in the second step will be initialized with the weights from the first step. The meta-block is used for pre-training. (iii) **Encoder-shared yet tail blocks trained in parallel**: the encoder is shared between two tail blocks and three blocks compose one joint network. All tasks are trained in parallel. (iv) **Encoder-shared and tail blocks trained iteratively**: in each training step, the encoder is firstly connected with the meta-block. We compute



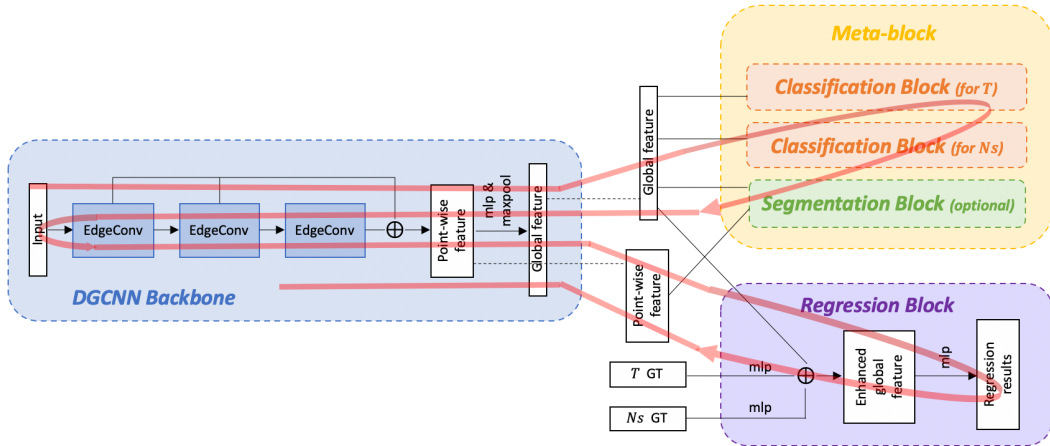


Figure 5: 3D regression framework for baseline results, the red line indicates that the tails blocks are trained in an iterative way.

Table 4: 3D regression baseline results. SRE stands for size relative error, GLE stands for gear location error, MRE stands for motor rotation error, SLE stands for screw location error. Method details are given in 5.3

Method	Classification $\uparrow$		Segmentation $\uparrow$		Regression $\downarrow$					
	$T$ accuracy	$N_s$ accuracy	mIoU	Screw mIoU	SRE	GLE $_{xyz}$	GLE $_{xz}$	MRE	SLE $_{xyz}$	SLE $_{xz}$
Separate	100	82.54	<b>92.55</b>	<b>81.19</b>	6.30%	0.3632	0.3626	0.8528	0.5925	0.5903
Pre-train	100	82.54	<b>92.55</b>	<b>81.19</b>	6.27%	0.3335	0.3326	<b>0.8403</b>	0.5521	0.5509
Parallel	100	84.67	90.23	80.22	6.21%	0.3423	0.3411	0.8522	0.5636	0.5622
Iterative	100	<b>86.15</b>	92.23	80.52	<b>5.90%</b>	<b>0.3251</b>	<b>0.3245</b>	0.8445	<b>0.5416</b>	<b>0.5395</b>

the loss, perform gradient back propagation and update model weights with these two blocks. Then the same encoder connects with the regression block in a switch manner, the input is reprocessed with the weight-updated encoder to get new encoded representations which is used for computing the regression loss. We then again perform gradient back propagation and weight update in these two blocks. This action performs iteratively. The red line in Figure 5 indicates how the information flows during one training step.

The experiments of four methods used same settings. We set a sub-point cloud size of 2048 points for batch training. The batch size is 16. The optimizer AdamW is used with a. epoch number of 100. The learning rate starts with 0.001 and decays to 0.00001 with a *cos* decay. In the EdgeConv blocks, we use  $K = 32$  when selecting neighbor points. As shown in Table 4, performing attribute regression with semantic segmentation decreases the regression error. Regression can also help in determining the number of cover screws. Among them, the iterative training-based method achieves best performance in most metrics. There are surely some other better ways in designing the architecture, *e.g.*, post-processing on the segmentation results for a better  $N_s$  classification results, or use 2D detection results as the supplementary input. We would like to leave this as an open question for other interested researchers.

## 6 CONCLUSION

In this paper, a benchmark is proposed using synthetic 2D and 3D dataset, in which motors are the key objects. Motor attributes are saved during the dataset generation and are suitable for the less explored attribute regression task. Apart from the common computer vision tasks including classification, detection, and segmentation, several metrics have been designed especially for the regression task. Baseline results on several tasks are also provided. We hope this work could contribute to the attribute regression or multi-task learning domain in the computer vision community and inspire the development of other novel algorithms in the future.

## REFERENCES

- Iro Armeni, Ozan Sener, Amir Roshan Zamir, Helen Jiang, Ioannis K. Brilakis, Martin Fischer, and Silvio Savarese. 3d semantic parsing of large-scale indoor spaces. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1534–1543, 2016.
- Christian Bonatti, Sylvain Crovisier, Lorenzo Diaz, and Amie Wilkinson. What is... a blender? *arXiv preprint arXiv:1608.02848*, 2016.
- Alexandre Boulch, Bertrand Le Saux, and Nicolas Audebert. Unstructured point cloud semantic labeling using deep segmentation networks. *3DOR@ Eurographics*, 3, 2017.
- Angel X. Chang, Thomas A. Funkhouser, Leonidas J. Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, L. Yi, and Fisher Yu. Shapenet: An information-rich 3d model repository. *ArXiv*, abs/1512.03012, 2015.
- Can Chen, Luca Zanotti Fragonara, and Antonios Tsourdos. Gapnet: Graph attention based point neural network for exploiting local feature of point cloud. *Neurocomputing*, 438:122–132, 2021.
- Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas A. Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2432–2443, 2017.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, K. Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 248–255, 2009.
- Maximilian Denninger, Martin Sundermeyer, Dominik Winkelbauer, Youssef Zidan, Dmitry Olefir, Mohamad Elbadrawy, Ahsan Lodhi, and Harinandan Katam. Blenderproc. *arXiv preprint arXiv:1911.01911*, 2019.
- David Eigen and Rob Fergus. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. *2015 IEEE International Conference on Computer Vision (ICCV)*, pp. 2650–2658, 2015.
- Nico Engel, Vasileios Belagiannis, and Klaus C. J. Dietmayer. Point transformer. *IEEE Access*, 9: 134826–134840, 2021.
- Mark Everingham, Luc Van Gool, Christopher K. I. Williams, John M. Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International Journal of Computer Vision*, 88:303–338, 2009.
- Yuan Gao, Qi She, Jiayi Ma, Mingbo Zhao, W. Liu, and Alan Loddon Yuille. Nddr-cnn: Layer-wise feature fusing in multi-task cnns by neural discriminative dimensionality reduction. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3200–3209, 2019.
- Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3354–3361, 2012.
- David Griffiths and Jan Boehm. Synthcity: A large scale synthetic point cloud. *ArXiv*, abs/1907.04758, 2019.
- Michael Gschwandtner, Roland Kwitt, Andreas Uhl, and Wolfgang Pree. Blensor: Blender sensor simulation toolbox. In *International Symposium on Visual Computing*, pp. 199–208. Springer, 2011.
- Meng-Hao Guo, Junxiong Cai, Zheng-Ning Liu, Tai-Jiang Mu, Ralph Robert Martin, and Shimin Hu. Pct: Point cloud transformer. *Comput. Vis. Media*, 7:187–199, 2021.
- Gary B. Huang, Marwan A. Mattar, Tamara L. Berg, and Eric Learned-Miller. Labeled faces in the wild: A database for studying face recognition in unconstrained environments. 2008.

- Mingyang Jiang, Yiran Wu, Tianqi Zhao, Zelin Zhao, and Cewu Lu. Pointsift: A sift-like network module for 3d point cloud semantic segmentation. *arXiv preprint arXiv:1807.00652*, 2018.
- Mahdi M. Kalayeh, Boqing Gong, and Mubarak Shah. Improving facial attribute prediction using semantic segmentation. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 4227–4235, 2017.
- Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7482–7491, 2018.
- Samin Khan, Buu Phan, Rick Salay, and Krzysztof Czarnecki. Procsy: Procedural synthetic dataset generation towards influence factor studies of semantic segmentation networks. In *CVPRW*, pp. 88–96, 2019.
- Sebastian Koch, Albert Matveev, Zhongshi Jiang, Francis Williams, A. Artemov, Evgeny Burnaev, M. Alexa, D. Zorin, and Daniele Panozzo. Abc: A big cad model dataset for geometric deep learning. *CVPR*, pp. 9593–9603, 2019.
- Iasonas Kokkinos. Ubernet: Training a universal convolutional neural network for low-, mid-, and high-level vision using diverse datasets and limited memory. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 5454–5463, 2017.
- Felix Järemo Lawin, Martin Danelljan, Patrik Tosteberg, Goutam Bhat, Fahad Shahbaz Khan, and Michael Felsberg. Deep projective 3d semantic segmentation. In *International Conference on Computer Analysis of Images and Patterns*, pp. 95–107. Springer, 2017.
- Le Hoang-An, Thomas Mensink, Partha Das, Sezer Karaoglu, and Theo Gevers. Eden: Multimodal synthetic dataset of enclosed garden scenes. *2021 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pp. 1578–1588, 2021.
- Dangwei Li, Xiaotang Chen, and Kaiqi Huang. Multi-attribute learning for pedestrian attribute recognition in surveillance scenarios. *2015 3rd IAPR Asian Conference on Pattern Recognition (ACPR)*, pp. 111–115, 2015.
- Dangwei Li, Xiaotang Chen, Z. Zhang, and Kaiqi Huang. Pose guided deep model for pedestrian attribute recognition in surveillance scenarios. *2018 IEEE International Conference on Multimedia and Expo (ICME)*, pp. 1–6, 2018.
- Zhidong Liang, Ming Yang, Hao Li, and Chunxiang Wang. 3d instance embedding learning with a structure-aware loss function for point cloud segmentation. *IEEE Robotics and Automation Letters*, 5:4915–4922, 2020.
- Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014.
- Yu Liu, Fangyin Wei, Jing Shao, Lu Sheng, Junjie Yan, and Xiaogang Wang. Exploring disentangled feature representation beyond face identification. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2080–2089, 2018.
- Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. *2015 IEEE International Conference on Computer Vision (ICCV)*, pp. 3730–3738, 2015.
- Daniel Maturana and Sebastian Scherer. Voxnet: A 3d convolutional neural network for real-time object recognition. In *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 922–928. IEEE, 2015.
- Ishan Misra, Abhinav Shrivastava, Abhinav Kumar Gupta, and Martial Hebert. Cross-stitch networks for multi-task learning. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3994–4003, 2016.
- Kaichun Mo, Shilin Zhu, Angel X. Chang, L. Yi, Subarna Tripathi, Leonidas J. Guibas, and Hao Su. Partnet: A large-scale benchmark for fine-grained and hierarchical part-level 3d object understanding. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 909–918, 2019.

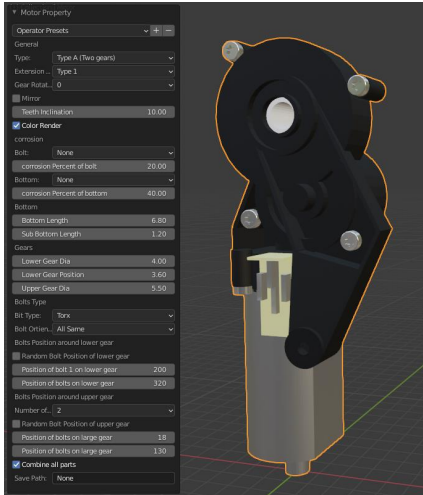
- R. Pierdicca, M. Mameli, E. Malinverni, M. Paolanti, and E. Frontoni. Automatic generation of point cloud synthetic dataset for historical building representation. In *AVR*, pp. 203–219, 2019.
- C. Qi, Hao Su, Kaichun Mo, and L. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. *CVPR*, pp. 77–85, 2017a.
- Charles R Qi, Li Yi, Hao Su, and Leonidas J Guibas. PointNet++: Deep hierarchical feature learning on point sets in a metric space. *arXiv preprint arXiv:1706.02413*, 2017b.
- G. Ros, Laura Sellart, Joanna Materzynska, David Vázquez, and Antonio M. López. The synthia dataset: A large collection of synthetic images for semantic segmentation of urban scenes. *CVPR*, pp. 3234–3243, 2016.
- Sebastian Ruder. An overview of multi-task learning in deep neural networks. *ArXiv*, abs/1706.05098, 2017.
- Sebastian Ruder, Joachim Bingel, Isabelle Augenstein, and Anders Søgaard. Latent multi-task architecture learning. In *AAAI*, 2019.
- Olga Russakovsky and Li Fei-Fei. Attribute learning in large-scale datasets. In *ECCV Workshops*, 2010.
- Nikolaos Sarafianos, Xiang Xu, and I. Kakadiaris. Deep imbalanced attribute classification using visual attention aggregation. *ArXiv*, abs/1807.03903, 2018.
- Pierre Sermanet, David Eigen, Xiang Zhang, Michaël Mathieu, Rob Fergus, and Yann LeCun. Overfeat: Integrated recognition, localization and detection using convolutional networks. *CoRR*, abs/1312.6229, 2014.
- Shuran Song, Samuel P. Lichtenberg, and Jianxiong Xiao. Sun rgb-d: A rgb-d scene understanding benchmark suite. *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 567–576, 2015.
- Marvin Teichmann, Michael Weber, Johann Marius Zöllner, Roberto Cipolla, and Raquel Urtasun. Multinet: Real-time joint semantic reasoning for autonomous driving. *2018 IEEE Intelligent Vehicles Symposium (IV)*, pp. 1013–1020, 2018.
- Hugues Thomas, Charles R Qi, Jean-Emmanuel Deschaud, Beatriz Marcotegui, François Goulette, and Leonidas J Guibas. Kpconv: Flexible and deformable convolution for point clouds. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 6411–6420, 2019.
- Jonathan Tremblay, Thang To, and Stan Birchfield. Falling things: A synthetic dataset for 3d object detection and pose estimation. *CVPRW*, pp. 2119–21193, 2018.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pp. 5998–6008, 2017.
- Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E Sarma, Michael M Bronstein, and Justin M Solomon. Dynamic graph cnn for learning on point clouds. *Acm Transactions On Graphics (tog)*, 38(5): 1–12, 2019.
- Wenxuan Wu, Zhongang Qi, and Fuxin Li. Pointconv: Deep convolutional networks on 3d point clouds. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 9613–9622, 2019.
- Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shapes. *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1912–1920, 2015.
- Dan Xu, Wanli Ouyang, Xiaogang Wang, and N. Sebe. Pad-net: Multi-tasks guided prediction-and-distillation network for simultaneous depth estimation and scene parsing. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 675–684, 2018.

Yu Zhang and Qiang Yang. A survey on multi-task learning. *ArXiv*, abs/1707.08114, 2017.

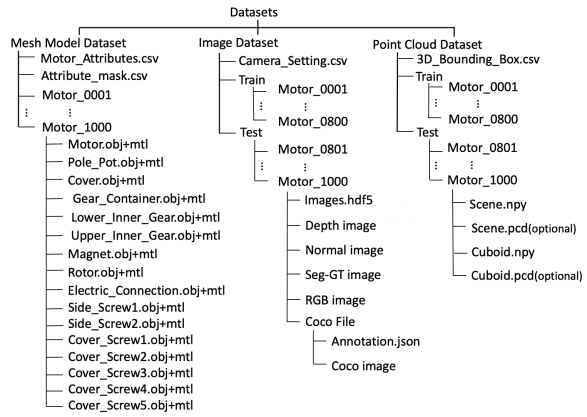
Hengshuang Zhao, Li Jiang, Jiaya Jia, Philip H. S. Torr, and Vladlen Koltun. Point transformer. *ArXiv*, abs/2012.09164, 2020.

Qingnan Zhou and A. Jacobson. Thingi10k: A dataset of 10, 000 3d-printing models. *ArXiv*, abs/1605.04797, 2016.

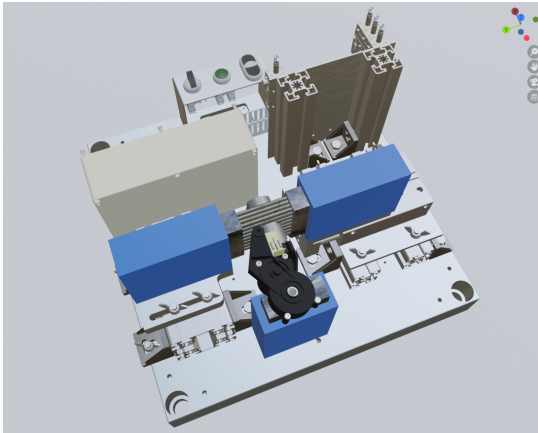
## A APPENDIX



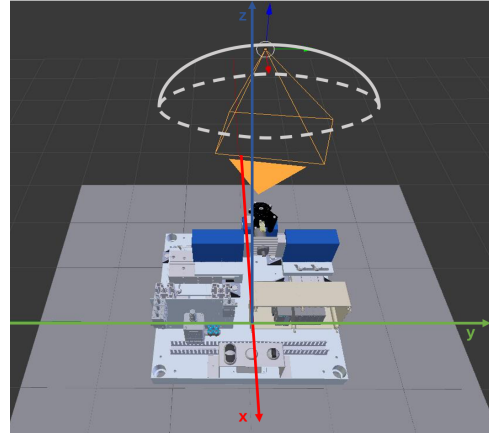
(a) Developed addon.



(b) Tree structure of our dataset.



(c) Blender scene.



(d) Camera setting.

Figure 6: Appendix figures.