

# A Benchmark Grocery Dataset of Realworld Point Clouds From Single View

Shivanand Venkanna Sheshappanavar, Tejas Anvekar, Shivanand Kundargi,  
 Yufan Wang and Chandra Kambhamettu  
 Video/Image Modeling and Synthesis (VIMS) Lab., Dept. of Computer and Information Sciences  
 University of Delaware, Newark, DE, USA 19716  
 {ssheshap, tanvekar, leff, chandrak}@udel.edu

## Abstract

*Fine-grained grocery object recognition is an important computer vision problem with broad applications in automatic checkout, in-store robotic navigation, and assistive technologies for the visually impaired. Existing datasets on groceries are mainly 2D images. Models trained on these datasets are limited to learning features from the regular 2D grids. While portable 3D sensors such as Kinect were commonly available for mobile phones, sensors such as LiDAR and TrueDepth, have recently been integrated into mobile phones. Despite the availability of mobile 3D sensors, there are currently no dedicated real-world large-scale benchmark 3D datasets for grocery. In addition, existing 3D datasets lack fine-grained grocery categories and have limited training samples. Furthermore, collecting data by going around the object versus the traditional photo capture makes data collection cumbersome. Thus, we introduce a large-scale grocery dataset called 3DGrocery100. It constitutes 100 classes, with a total of 87,898 3D point clouds created from 10,755 RGB-D single-view images. We benchmark our dataset on six recent state-of-the-art 3D point cloud classification models. Additionally, we also benchmark the dataset on few-shot and continual learning point cloud classification tasks. Project Page: <https://bigdatavision.org/3DGrocery100/>.*

## 1. Introduction

3D Computer vision is an active research area with broad applications in autonomous navigation [1], healthcare [2], and augmented/virtual reality [3]. Among these applications, autonomous robotic navigation [4] and assistance to the visually impaired [5, 6] require recognizing and localizing objects in real-world scenarios. Grocery recognition [7, 8] in stores is a challenging problem. Real-world grocery recognition helps recognize misplaced products, identify product depletion [9], automatic grocery checkout, and assistive technologies for the visually impaired for a

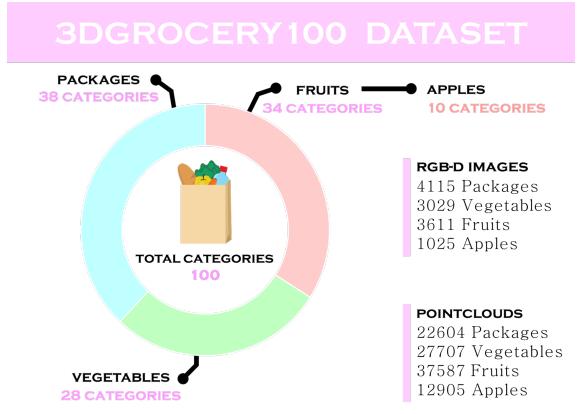


Figure 1. 3DGrocery100 Dataset Statistics. The dataset constitutes 10,755 RGB-D images and 87,898 point clouds spread across 100 classes. At a high level, the groceries are categorized into Fruits (10 apple and 24 non-apple classes), Vegetables (28), and Packages (38). Note: Apples are a subset of Fruits. Non-apples fruit RGB-D image count: 2,586; point cloud count: 24,682.

comfortable grocery experience. Developing robust grocery recognition systems would demand larger datasets to effectively train and deploy deep learning models.

The past two decades have witnessed real-world grocery datasets growth (see Table 1) with varying sample sizes and image resolutions for retail product recognition. 2D representations inherently suffer from the loss of shape information and mainly cater to the needs of deep learning models [10] that rely on 2D images. Whereas with the availability of shape information in our grocery dataset, intrinsic characteristics of grocery objects would be captured which aids the recognition applications. Furthermore, unlike the 2D image-based counterpart, 3D point cloud-based recognition systems [11] suffer due to the lack of large-sized benchmark datasets further limiting our ability to build, evaluate, and compare the strengths of different methods, especially the recent data-hungry deep learning techniques.

Table 1. Details of grocery datasets. OW - Open World, S - Studio, Acc - Accuracy, AP - Average Precision (*italicized*). “-” not available.

Dataset	Classes	Samples	Train	Test	OW	Acc/AP	Model/Paper
2D Datasets							
<b>SOIL-47</b> [12](2002)	47	987	-	-	X	71.0	MNS [13]
<b>GroZi-120</b> [14](2007)	120	11194	676	11194	X	18	SIFT [15]
<b>Supermarket</b> [16](2010)	15	2633	-	-	✓	98	SVM-Fusion [16]
<b>GroZi-3.2K</b> [17](2014)	80	9101	8421	680	✓	23.49	Exemplar-based MLIC [17]
<b>Grocery</b> [18](2015)	10	354	-	-	X	92.3	SVM [18]
<b>Freiburg</b> [19](2016)	25	4947	4000	1000	✓	78.9	CaffeNet[20]
<b>MVTec D2S</b> [21](2018)	60	21000	4380	13020	X	79.9	Mask R-CNN [22]
<b>Grocery Store</b> [23](2019)	81	5125	2640	2485	✓	84.0	DenseNet-169 [24]
<b>RPC</b> [25](2019)	200	<b>83739</b>	-	-	X	56.68	Syn [26, 27]+Render [28]
<b>Magdeburg</b> [29] (2019)	942	65315	<b>23360</b>	<b>41955</b>	X	91.83	VGG-16 [30]
<b>TGFS</b> [31](2019)	24	30000	22815	15212	✓	65.5	FCIOD [31]
<b>SKU-110K</b> [32](2019)	<b>110712</b>	11762	8233+588	2941	✓	49.2	Deep IoU Detection [32]
<b>RP2K</b> [33](2021)	2388	10385	-	-	✓	<b>95.18</b>	ResNet-34 [10]
3D Datasets							
<b>BigBird</b> [34](2014)	100	600	-	-	X	-	-
<b>HOPE</b> [35](2022)	28	238+914	-	-	Toy	-	-
<b>Object-Verse</b> [36](2023)	<b>21,000</b>	<b>818,000</b>	-	-	X	28.3	GOL+3DCP [36]
<b>3DGrocery63 (Ours)</b>	63	87898	66032	21866	✓	-	-
<b>3DGrocery (Ours)</b>	100	87898	66032	21866	✓	<b>50.50</b>	LocalFeatures [37, 38]

3D representations such as RGB-D and point clouds are richer in geometric information that closely represents the real world. The recent development of 3D deep learning models [39–44] exploits the representational power of 3D data. However, these models are often trained and evaluated on limited datasets. Typically a CAD-based synthetic dataset, i.e., ModelNet40 [45] with 12,311 point cloud samples spread across 40 classes, and a real-world dataset with variants, i.e., ScanObjectNN [46] with 14,510 samples spread across 15 classes. Most recently, very large 3D datasets such as Objaverse [36] and Omni3D [47] are released, but they either artist-curated or do not contain sufficient classes of grocery, respectively.

Technological innovations over the past decade in 3D cameras/sensors, such as Light Detection And Ranging (LiDAR[48]), TrueDepth[49], Kinect[50], etc., have increased the availability and accessibility of 3D sensors resulting in improved data acquisition. Additionally, the availability of high computational devices for 3D data processing has paved the way for many real-time 3D applications solving complex 3D problems. Furthermore, the availability of 3D sensors in handheld devices has opened new avenues of research that can address high-impact problems. However, despite the availability of 3D sensors in mobile devices, real-world 3D datasets are scarce. Notably, as listed in Table 1, the availability of 3D datasets for grocery recognition is limited. Collecting 3D data using mobile devices requires going around the object and scanning

it to create a 3D point cloud. It is often not practical in a grocery store to go around objects as they are placed on racks/shelves. These issues encourage us to explore and capture the single-view RGB and Depth images which can be processed into point cloud instances of groceries (dataset statistics in Figure 1). Though these point clouds are incomplete, they carry geometric details of the grocery objects. Single-view RGB-D image acquisition is more straightforward and close to real-world scenarios mimicking photo capturing but provides both objects’ geometric structures and point colors. 3D point clouds are a set of XYZ coordinates in 3D space with properties of permutation invariance.

From the advent of PointNet [39] to the most recent PointNeXt [44], a few hundred novel networks have been proposed for point cloud classification. Deep learning models inspired by these methods process raw point clouds and are among the popular choices for deployment among self-driving applications [51]. The performance of deep learning models depends heavily on the quality and quantity of the datasets used in training them. However, as stated in PatchAugment [52], most of the available 3D datasets are synthetic. Only a handful of real-world 3D point cloud datasets [46] are available for researchers to train deep learning models for real applications. Furthermore, real-world benchmark 3D datasets, especially those collected using mobile devices, are limited. We present the largest 3D point cloud grocery dataset from RGB-D images collected using mobile phones (iPhone 12 Pro and Pro Max).

Fresh produce grocery items, though packaged, are subject to pricing errors and usually require touching/picking by customers, which is discouraged during pandemics. Large items, such as watermelons, pose additional challenges for barcode scanning due to their size. These issues complicate grocery recognition tasks, highlighting the need for no-contact recognition methods. Additionally, the random orientation of produce and often hidden barcode labels further challenge effective recognition. Besides, image-based object recognition [53] relies heavily on texture, color, and appearance cues. 2D images miss 3D geometric details, making 3D datasets essential for utilizing the 3D features of groceries in classification. In this regard, we present the 3DGrocery100 point cloud dataset obtained from RGB-D images as shown in Figure 2. Our large-scale 3D grocery dataset will enable the grocery recognition community to apply, develop, and adapt various deep learning techniques for 3D grocery classification. The key contributions of our paper are four-fold:

- **3DGrocery100:** We present a novel benchmark 3D point cloud grocery dataset with 87,898 instances obtained from 10,755 RGB-D images across 100 categories.
- **Classification Evaluation:** We benchmark our 3D dataset and its five subsets, each with two variants (color and no color), on six recent state-of-the-art (SOTA) models for 3D point cloud classification tasks. We also test and compare the recent SOTA models for robustness.
- **Few-shot Evaluation:** We merge similar shape classes to create a 63-class variant of the dataset, i.e., 3DGrocery63, for Few-shot evaluation and propose a strong baseline that evaluates the true generalization power of meta-learners on point cloud few-shot classification.
- **Class Incremental Learning Evaluation:** We set a 3D grocery benchmark using a class-incremental learning baseline. We extend 2D LWF [54] with a dynamic multi-head classifier to benchmark on our 3D grocery dataset.

## 2. Related Works

### 2.1. 2D Grocery Datasets

Grocery datasets in the 2D domain (images) have been available since the mid-90s; veggievision [55] consists of 5000 images and 150 categories. Over the next decade, 2D grocery datasets [12, 14] mainly improved the RGB images' resolution. The improvement in computing hardware during the past 15 years has paved the way to larger datasets [16, 17, 19, 21, 23, 25, 29, 31–33], gradually increasing the number of categories and images per category. The Supermarket Produce [16] dataset features varied product view angles but uses a white canvas background, unlike real grocery store settings. Our dataset captures products in actual store environments, preserving the natural variations in view angles, elevations, and distances to objects.

The most recent dataset, RP2K [33], offers 2388 product categories, but it is limited to an average of only 37 RGB images per class. Like earlier datasets with challenges from specular reflections in packaged products, our images also feature classes with reflective packaging, complicating 3D analysis as discussed later. Recently, [5] carefully considered the navigation context and divided the dataset into four abstraction levels, i.e., product, shelf, trail, and others. Our dataset supports the first two abstraction levels of grocery recognition. The comprehensive review work [56] on retail grocery categorization while listing the absence of RGB-D images for groceries encourages the exploration of RGB-D to address critical problems in grocery recognition.

### 2.2. 3D Grocery Datasets

While there are several 2D datasets in the grocery domain, the 3D grocery datasets among the representations such as RGB-D, point cloud, mesh, and voxels are rare [34, 35, 57]. Among these datasets, though [57], are RGB-D images of objects, it is not entirely a grocery dataset but only contains a few classes that belong to the grocery. For example, the HOPE [35] dataset provides the RGB-D video recordings of 28 categories of toy grocery objects in different lighting conditions. However, these are toy grocery items and do not overlap with any of the existing real-world brands. Further, the toy grocery items are placed in a non-store setting limiting the models trained in demonstrating in a real-store scenario. BigBird [34] provides 3D point clouds along with RGB-D images but is a small dataset with only 600 images.

### 2.3. 3D Point Cloud Classification

Point cloud analysis has gained traction since the pioneering work of PointNet [39] that directly processes raw points through Multi-Layer Perceptrons (MLPs). However, PointNet loses valuable local geometric information while aggregating global features using max-pooling. Overcoming this limitation, PointNet++ [40] captures local neighborhoods to learn local semantic information. DGCNN [41] introduced EdgeConv to capture local geometry by generating edge features to distinguish a point from its neighbors.

PCT [42] learns features through the attention mechanism of Transformers [58]. PointMLP [43] employs a simple feed-forward residual MLP network aggregating hierarchically extracted local features. PointNeXt [44] improves PointNet++ [40] by using better training strategies and an inverted residual bottleneck design with separable MLPs. These methods use ModelNet40 [45] (synthetic) and ScanObjectNN [46] (real-world) datasets with 40 and 15 classes, respectively. Despite groceries being 3D objects, there are no benchmark 3D grocery datasets for training point-based deep learning models. Our dataset enables the training of models to classify groceries in real-world stores.

## 2.4. Point Cloud Few Shot Learning

Recent progress in Few-Shot Learning (FSL) for 2D image processing has led to two main approaches: metric-based methods, which improve class discrimination through feature space, and optimization-based methods, which focus on model adaptability to new classes. One such metric-based approach, Prototypical Net [59], introduces class-prototype, which is the mean of the support features of each class. Meanwhile, optimization-based FSL techniques aim for quick adaptation to new tasks through minimal gradient updates, with ongoing research enhancing 2D FSL. Meta-baseline [60] introduces a cosine metric classifier with learnable weights, demonstrating improved performance compared to the Squared Euclidean Distance. SimpleTrans [61] proposes a straightforward transform function to adjust the weights of different channels, mitigating the channel bias problem in FSL. In this paper, we follow GPrNet [62] for point-cloud FSL using ProtoNet and observe the grocery dataset as a strong generalization benchmark.

## 2.5. Point Cloud Class-Incremental Learning

Many essential robotic perception applications heavily rely on real-world data. However, real-world data is dynamic and arrives in a continuous stream. Incremental learning strategies have been developed to learn from this incoming data effectively. These strategies can be broadly classified into Task-Incremental [63–66], Domain-Incremental [67, 68], and Class-Incremental learning (CIL) [69–71]. CIL closely emulates the dynamic arrival of data in real-world scenarios. Baseline approaches in CIL like [54, 72] address the problem of catastrophic forgetting in neural networks, which is essential when dealing with incremental real-world data. [73] explore learning without forgetting on-point clouds using semantic representation. Advancements in depth sensor technology have spurred growth in 3D and 2.5D data applications. CIL tackles the challenge of processing novel streaming data without task identifiers at test time, reflecting real-world conditions. This makes incremental learning on new data increasingly important.

## 3. Dataset creation

This section outlines our dataset creation process, as depicted in Figure 2, highlighting three main steps.

### 3.1. Data collection

Data was collected over four months from 18 different local grocery stores. Based on the availability of items in the stores, the collection spanned from a few days to weeks. Grocery stores are well-lit environments with similar arrangements for each class. For example, most packages are placed on shelves on both sides of an aisle, while a few vegetables and fruits are arranged on a table.

#### 3.1.1 iOS app for collection

Recently, Apple iPhones have offered multiple ways to capture 3D objects. The latest iPhones Pro and Pro Max (version numbers 12, 13, 14, and 15) come with a LiDAR sensor, a go-to sensor to capture the scene’s depth in a range of up to 5 meters. The LiDAR camera gives a low resolution ( $192 \times 256$ ) depth map. The RGB image resolution is  $3024 \times 4032$ . In the case of LiDAR, the lower resolutions of depth maps result in sparse point clouds losing significant 3D information. Another way to capture depth is stereo vision using the built-in back-facing dual camera. The depth map and camera calibration metrics are saved as a 32-bit floating-point array. The resolution of these depth maps is  $576 \times 768$ . The point clouds generated with stereo depth preserve the local geometry of the object, making it look more realistic (refer to Figure 8 and its description in the supplementary). The better quality of 3D point clouds from the depth captured using the stereo vision built-in back-facing dual camera [74] encouraged us to adopt the back-facing dual camera to capture depth data along with RGB images in our customized iOS app for data collection.

#### 3.1.2 Data Hierarchy and RGB-D dataset

Following [23], Figure 1 shows the high-level hierarchy as three main categories of groceries, i.e., Fruits, Vegetables, and Packages, along with the count of 2D RGB-D images and 3D point clouds from them. The Fruits category is divided into ten classes of apples and 24 classes of non-apple fruits. Vegetables and Packages have 28 and 38 classes, respectively. We combine these categories to form the “Full” dataset of 100 classes. Figures 3 and 4 show visual examples of classes from Apple10 and (Fruits, Vegetables, and Packages) respectively. We provide two variants of the point clouds i.e., with and without colors. These grocery items were positioned in racks (most packages and some vegetables), on tables (most fruits and a few vegetables), inside refrigerators (often packages), and in the natural environment of local grocery stores with similar lighting settings. The resolution of RGB images is  $3024 \times 4032$ , and the resolution of the depth maps is  $576 \times 768$ . During data collection, we used portrait mode, making the height of images larger than the width. In addition, our dataset also has challenging images, which have products with reflective surfaces, images taken at oblique views, and images with darker backgrounds.

#### 3.1.3 Challenges

Data collection using mobile phones would often result in shaky and blurred images. Such images are harder to annotate, resulting in heavily distorted or noisy 3D point clouds. We have carefully discarded such images from our

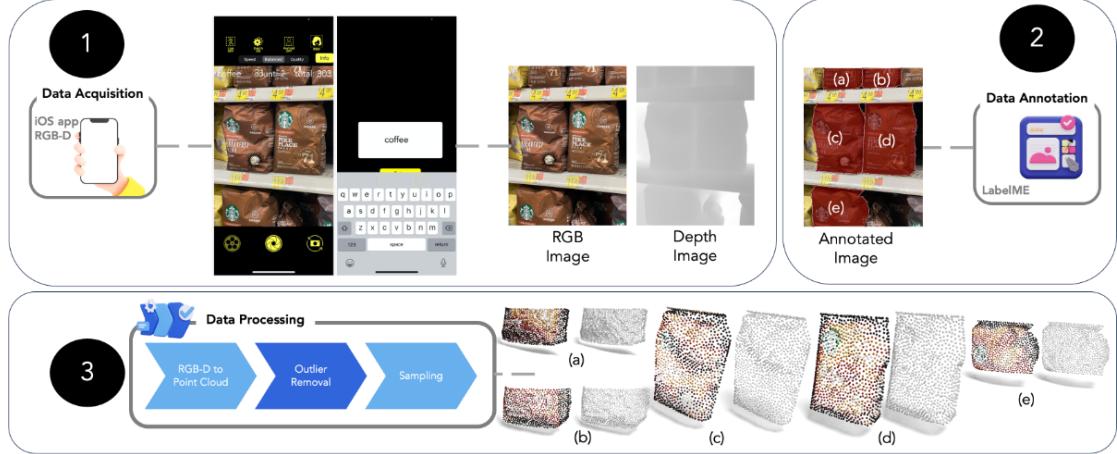


Figure 2. (1) iOS app to capture an RGB image and a Depth Image (darker regions nearer to the camera) for a class (2) annotations on the RGB image - coffee instances numbered [a-e], (3) [a-e] 1024 Farthest point sampled 3D points of 5 coffee objects with and without colors.

dataset. In some stores, a few grocery items, such as milk and eggs, are placed in open racks with sufficient cooling, while other stores place these items in a refrigerator or a cooler. We have collected and annotated them as separate classes to enable recognition systems to learn the differences between such classes. Although the doors of these refrigerators are slightly reflective, we observed no depth distortion for the milk-in-cooler and eggs-in-cooler classes. However, the colors are slightly different compared to their non-refrigerated counterpart classes. We left a gap of at least 24 hours between the visits to each grocery store to capture variations in the stores’ inventory, especially those of fruits and vegetables. A few items take longer to sell, and changes in their inventory often take more time. We planned data collection of such items with longer gaps.

### 3.2. Data Annotation

To annotate such a large dataset of 10755 grocery images, we started annotation in parallel while the data collection

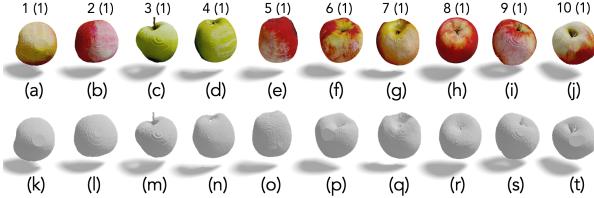


Figure 3. (a-j) point clouds of apple-evercrisp, apple-fuji, apple-golden-delicious, apple-granny-smith, apple-honeycrisp, apple-pazazz, apple-pink-lady, apple-red-delicious, apple-royal-gala, and apple-wild-twist with color. (k-t) point clouds in the same order without color. Colors play a significant role in distinguishing objects from different classes. Without color, all ten classes of apples appear very similar in the 3D point cloud representation.

was still in progress. We used LabelMe [75] to annotate collected 2D RGB images. We drew a polygon for each object in the picture to approximate its boundary as shown in Figure 2 (data annotation section). Since the annotated 2D RGB images and depth maps generate the 3D dataset, any extra pixel (due to human error) that does not belong to the object of interest would have a very different depth value leading to background noise or outliers. Therefore, we annotated the polygon carefully on or within the object boundary to avoid additional background noise. As objects in an image tend to get occluded or are partly out of the frame, we annotated objects whose visibility was significant, i.e., approximately at least 25% visible.

### 3.3. Data Processing

We used Pinhole Camera Intrinsic Parameters to convert RGB-D images to point clouds using Open3D [76] library. Processing RGB-D images for point cloud objects is challenging and produces outliers, especially at the object boundaries, due to minor annotation issues that overlap the background or neighboring objects. We compute the median of the point cloud and use it as a threshold to remove background points from the point cloud. The reflective and transparent surfaces of the grocery items often result in outlier depth values resulting in outlier 3D points. We use PointCleanNet [77] for removing outliers from the point clouds (refer to Figure 9 in supplementary). PointCleanNet cleans point cloud datasets in a two-step process: first, it finds and removes outliers, then denoises the remaining points by estimating a per-point displacement vector. We use their pre-trained outlier-removal model to process our RGB-D point cloud data. This deep learning approach requires less labor and does not make any assumptions about the noise model or object surface.



Figure 4. 3D point cloud representations of 24 (of the 34) classes of fruits (non-apples) with colors. The sets (bartlett-pear, donjau-pear, red-pear, pear-bosc), (cantaloupe, watermelon, honeydew-melon), (lemon, lime), (nectarine, peach, plum), and (grapefruit, navel-orange) are similar in shapes, and these classes get merged in case of 3DGrocery63. 28 classes of Vegetables with colors. 38 classes Packages with colors. \*-bp is bell-pepper and \*-s is spaghetti. The first number represents the class ID in 3DGrocery100, and the number in the parenthesis represents the class ID in 3DGrocery63 obtained after merging similar-shaped classes in 3DGrocery100.

We first separate XYZ coordinates from the XYZRGB data. Then, pass the XYZ data into the pre-trained Point-CleanNet [77], to estimate the per-point probabilities to check if a point is an outlier. In our case, a relatively lower threshold of 0.4 yields better results when compared to PointCleanNet [77], which uses a threshold of 0.5. Finally, we combine the RGB information from the original data and XYZ coordinates based on the estimated probability values to produce “cleaned” data in XYZRGB format. These XYZ values are normalized to the range [-1,1].

Even though we selectively annotated significantly visible objects, many of them ended up with only a few points in their 3D point cloud representations due to their size or placement within the image. We retained objects with at least 10,000 points and discarded the rest. After applying PointCleanNet [77] for outlier removal, point clouds from each class are similar to those shown in Figures 3 and 4. We used the farthest point sampling (FPS) method to sample 1024 points (with and without colors as in Figure 2) for our experiments. Visual examples of five samples per class with 1024 FPS points are shown in the supplementary. Table 2 shows the train and test split of all the subsets of our dataset. We divide the RGBD images into 75% training and 25% testing and use all the point clouds obtained from the 75-25 split as train and test point clouds, respectively.

## 4. Applications

Our dataset constitutes three categories, i.e., Fruits, Vegetables, and Packages with 34, 28, and 38 classes, respectively. The Fruits category is further divided into ten classes of apples (Apple10) and 24 non-apple classes. In this section, we analyze our experimental results on the complete dataset (Full) and four subsets (Apple10, Fruits, Vegetables, and Packages). Table 2 shows the number of images and the corresponding point clouds created from the images.

### 4.1. Real-World Point Cloud Classification

We consider six representative works on 3D object classification namely, PointNet [39], PointNet++[40], DGCNN[41], PCT [42], PointMLP[43], and PointNeXt[44] to benchmark our dataset. We retain the training parameters, such as learning rate, optimizer, number of epochs, batch size, and weight decay for each of these methods as per the original papers. While PointNet [39] and PointNet++ [39] had shown poor performance on the real-world dataset ScanObjectNN [46], the most recent deep learning models PointMLP [43] and PointNeXt [44] have achieved SOTA classification results. However, the hardest variant of ScanObjectNN [46] dataset has less than 15,000 samples and 15 classes. Table 2 shows the overall accuracy on both variants of the four subsets, and the Full dataset.

Table 2. 3D point cloud classification results on 3DGrocery100. Params in M (million) and FLOPs in G (giga).

Models	Apples	Fruits	Veggies	Packages	Full	#Params	FLOPs
#Classes	10	24	28	38	100		
#Images	1025	2586	3029	4115	10755		
Train	772	1944	2264	3103	8083		
Test	253	642	765	1012	2672		
#Patches	12905	24682	27707	22604	87898		
Train	9706	18406	20720	17214	66032		
Test	3199	6276	6987	5390	21866		
<b>INPUT: RGB Image</b>							
Swin T [78]	85.10	98.50	98.40	98.3	98.5	28.00	04.50
<b>INPUT: 1024 POINTS + COLORS</b>							
PointNet [39]	78.35	96.42	97.00	97.85	92.62	00.70	00.31
PointNet++ [40]	<b>84.00</b>	<b>97.58</b>	<b>98.25</b>	<b>98.46</b>	<b>95.21</b>	01.48	01.71
DGCNN [41]	83.59	98.01	97.88	98.53	96.20	01.81	05.39
PCT [42]	82.14	97.89	97.63	97.97	94.19	02.88	04.34
PointMLP [43]	90.18	97.75	97.68	98.38	<u>96.82</u>	13.24	31.35
PointNext [44]	<u>84.81</u>	<u>97.36</u>	<u>97.37</u>	<u>98.13</u>	94.44	04.52	06.49
<b>INPUT: ONLY 1024 POINTS</b>							
PointNet [39]	18.93	33.16	35.15	71.61	37.61	00.70	00.31
PointNet++ [40]	<b>26.20</b>	<b>47.98</b>	<b>54.11</b>	<b>82.70</b>	<b>48.65</b>	01.48	01.71
DGCNN [41]	19.04	37.01	42.16	79.44	44.19	01.81	05.39
PCT [42]	18.95	32.39	43.42	75.34	39.67	02.88	04.34
PointMLP [43]	19.16	36.38	47.29	79.94	<u>45.72</u>	13.24	31.35
PointNext [44]	<u>21.63</u>	<u>40.58</u>	<u>48.40</u>	<u>81.35</u>	43.73	04.52	06.49

## 4.2. Benchmarking Real-World Scenarios with Minimal Data

Building upon insights from Meta-Dataset [79], existing 3D few-shot classification benchmarks [80–83] primarily emphasize intra-dataset generalization. However, our objective is to benchmark models for effective generalization across entirely new distributions, even unseen datasets. To address this challenge, we introduce cross-domain benchmarking strategies for Meta-learners, spotlighting their generalization capabilities on 3D datasets.

The significance of 3D Grocery arises from its unique attributes. Comprising real-world 2.5D point clouds, the proposed 3DGrocery dataset is meticulously refined to eliminate inter-class intersections across categories, yielding 63 distinct classes as illustrated in Figure 4. This dataset remains distinct from well-known point cloud datasets like ShapeNet [45], ModelNet40 [45], and ScanObjectNN [46], exhibiting no overlaps. The absence of overlapping classes and differentiation from other popular datasets positions it perfectly for robust generalization benchmarks for Meta-Learners. Through 3D Grocery, we comprehensively assess how well Meta-Learning models can adapt and generalize to novel and distinct data distributions, offering insights into their genuine resilience in real-world challenges.

**Experimental Setup:** We conducted two experiments

*Baseline Few shot Evaluation* and *Few shot Meta-Learning* to demonstrate that our proposed dataset 3DGrocery63 can act as a very strong benchmarking dataset to perform weak vs. strong generalization tasks.

**(1) Baseline Few shot Evaluation:** We take pre-trained point cloud classifiers [39–43] on ModelNet40[45], and perform  $k$ -way,  $m$ -shot few-shot evaluation using features learnt by the classifier

Table 3. Quantitative analysis for Few-shot 3D point cloud classification on 3DGrocery63 dataset. All models are pre-trained on ModelNet40[45].

Model	Weight Init	5-ways		10-ways	
		10-shots	20-shots	10-shots	20-shots
<b>PointNet</b> [39]	Random	54.30	$\pm 09.51$	59.40	$\pm 09.59$
	MN40	56.26	$\pm 09.81$	61.58	$\pm 10.22$
<b>PointNet++</b> [40]	Random	55.18	$\pm 09.37$	60.80	$\pm 10.42$
	MN40	61.60	$\pm 09.23$	68.76	$\pm 09.05$
<b>DGCNN</b> [41]	Random	55.20	$\pm 10.26$	63.02	$\pm 09.43$
	MN40	<b>65.44</b>	$\pm 10.64$	<b>72.18</b>	$\pm 09.59$
<b>PointMLP</b> [43]	Random	40.34	$\pm 06.17$	47.14	$\pm 07.97$
	MN40	63.96	$\pm 11.75$	69.18	$\pm 09.74$
<b>PCT</b> [42]	Random	57.36	$\pm 09.12$	53.86	$\pm 08.73$
	MN40	62.45	$\pm 09.64$	66.52	$\pm 09.55$
<b>PointNext</b> [44]	Random	55.82	$\pm 08.48$	62.38	$\pm 09.46$
	MN40	60.64	$\pm 10.60$	65.12	$\pm 09.46$

as expressed by authors in [84]. The few-shot evaluation results on ScanObjectNN [46] are reported in Table 12 in supplementary and similarly results on proposed 3DGrocery63 are reported in Table 3; comparing both the tables, it is evident that classifiers generalize well on ScanObjectNN dataset while failing on 3DGrocery63, the main reason for this phenomenon is a data-inductive bias that is common in ScanObjectNN. **(2) Few shot Meta-Learning:** We further evaluate the aforementioned phenomenon by benchmarking on Few-shot Meta-Learning using ProtoNet [59] for point-cloud FSL as described in GPR-Net [62]. We propose six data splits for this setup: Train, Val, and Test (weak1, weak2, weak3, and strong). For the train, we propose to use ShapeNet55 [45] and exclude 15 categories that intersect with ScanObjectNN; the remaining are used for Val split. For weak1 → weak3 test splits, we use ScanobjectNN (ONLY OBJ, OBJ+BG, PB75) dataset and 3DGrocery63 for strong test split. We follow an episodic paradigm to train the ProtoNet version of 3D classifiers. The settings are given in the supplementary. Our findings are reported in Table 4 and Table 13 (in supplementary), which depicts the curse of data-inductive bias; every model generalizes well on weak generalization tasks but fails on proposed strong generalization. Finally, we conclude that the proposed 3DGrocery63 is a strong baseline to validate the true generalization of meta-learners on point cloud few-shot classification considering data-inductive biases.

## 4.3. Benchmarking Real-World Scenarios of Continual Data

Our extensive grocery dataset has the largest collection of point clouds captured in the real world. As we anticipate the possibility of new classes being introduced to our real-world data, we are actively investigating CIL methods, leveraging which training from scratch on entire data can be avoided. These methods will serve as a benchmark for our data and enable us to effectively handle the inclusion of novel classes without forgetting the information of the past as our dataset continues to grow. We follow PointCLIMB [85] for developing, the problem setting of Point-cloud CIL on the pro-

Table 4. Quantitative analysis for Meta-Learning on few-shot 3D point cloud classification. Avg Weak is the average of accuracies in % from three splits in a weak generalization task with 15 novel classes ([46]) and Strong is a strong generalization task with 63 novel classes.

Models	Avg Weak								Strong							
	5-ways				10-ways				5-ways				10-ways			
	5-shots		10-shots		5-shots		10-shots		5-shots		10-shots		5-shots		10-shots	
PointNet [39]	53.50	$\pm 0.63$	56.36	$\pm 0.60$	36.40	$\pm 0.37$	40.11	$\pm 0.37$	27.38	$\pm 0.41$	28.91	$\pm 0.43$	14.99	$\pm 0.23$	17.10	$\pm 0.24$
PointNet++ [40]	50.16	$\pm 0.60$	55.92	$\pm 0.62$	37.03	$\pm 0.36$	39.73	$\pm 0.36$	24.11	$\pm 0.37$	29.39	$\pm 0.43$	15.64	$\pm 0.23$	18.15	$\pm 0.25$
DGCNN [41]	51.30	$\pm 0.62$	56.47	$\pm 0.61$	36.80	$\pm 0.37$	40.08	$\pm 0.37$	<b>28.39</b>	<b><math>\pm 0.43</math></b>	<b>30.56</b>	<b><math>\pm 0.43</math></b>	<b>17.35</b>	<b><math>\pm 0.24</math></b>	18.24	$\pm 0.25$
PCT [42]	57.05	$\pm 0.64$	56.58	$\pm 0.64$	39.66	$\pm 0.37$	43.24	$\pm 0.37$	27.57	$\pm 0.41$	29.22	$\pm 0.43$	16.58	$\pm 0.25$	<b>18.59</b>	<b><math>\pm 0.26</math></b>
PointMLP [43]	46.13	$\pm 0.62$	50.14	$\pm 0.59$	32.76	$\pm 0.35$	40.45	$\pm 0.37$	24.70	$\pm 0.38$	25.52	$\pm 0.39$	13.45	$\pm 0.20$	14.76	$\pm 0.21$
PointNeXt [44]	<b>58.63</b>	<b><math>\pm 0.64</math></b>	<b>58.94</b>	<b><math>\pm 0.65</math></b>	<b>41.70</b>	<b><math>\pm 0.39</math></b>	<b>44.86</b>	<b><math>\pm 0.38</math></b>	28.23	$\pm 0.45$	28.95	$\pm 0.42$	16.29	$\pm 0.23$	17.81	$\pm 0.24$

posed dataset.

**problem setting:** Class-Incremental learning problem  $\mathcal{T}$  consists of sequence of  $k$  tasks:

$$\mathcal{T} = [(C^1, D^1), (C^2, D^2), \dots, (C^k, D^k)] \quad (1)$$

where each task  $\parallel$  consists of a set of classes  $C^k = \{c_1^k, c_2^k, \dots, c_{m^k}^k\}$  and  $D^k$  is the training data. The point cloud class-incremental problem in which  $D^k = \{(p_1, y_1), (p_2, y_2), \dots, (p_{l^k}, y_{l^k})\}$ , where  $p$  is point cloud with  $n$  points such that  $p \in \mathbb{R}^{n \times 3}$ . During training for task  $k$ , the learner only has access to  $C^k, D^k$ , whereas, during inference, the evaluation is done for the union of all previous tasks  $\bigcup_{i=1}^k C^i, D^i$ . For instance if we encounter task  $k = 2$ , the learner has access to  $(C^2, D^2)$  whereas evaluation is done for  $\{(C^1, D^1), (C^2, D^2)\}$ .

**Class-Incremental learning on 3D-Grocery:** we propose to benchmark a baseline approach known as “Learning Without Forgetting” (LWF) [54] used in the 2D realm to address the issue of catastrophic forgetting. We extend the methodology used in the 2D realm to 3D by adopting features from SOTA point cloud processing architectures such as PointNet [39], PointNet++ [40], DGCNN [41], PointMLP [43], and PCT [42], and using a dynamic multi-head classifier which adjusts itself automatically according to the novel classes that arrive in stream.

The combination of baseline and these advanced architectures results in our extended method called “LWF” and is compared with baselines Fine-tuning (FT) / Lower-bound, and Joint-training / Upper-bound. LWF aims to mitigate catastrophic forgetting in point clouds, ensuring our dataset’s adaptability to new and evolving classes while preserving the learned knowledge from previous data. For benchmarking on the CIL setting, we split our dataset into five tasks, with the base task involving 39 classes and the remaining tasks with six classes each arriving incrementally. We train a joint head classifier, which adjusts according to the arrived novel classes. We benchmark point cloud CIL with each task trained on 40 epochs. PointNeXt [44] and DGCNN [41] perform better on incremental tasks compared to other backbones, as shown in Table 5.

Table 5. Performance of different backbone’s on 3DGrocery63 dataset in a CIL scenario. **Joint**: Upper bound, **FT**: Fine-Tuning.

Backbone	# Classes $\rightarrow$	39	6	6	6	6	
		Loss $\downarrow$	Acc	Acc	Acc	Acc	
		FT	41.40	07.61	04.62	06.27	03.49
PointNet [39]	<b>LwF</b>	41.40	07.55	06.34	06.42	04.77	
		<b>Joint</b>	41.40	42.46	43.50	43.99	44.07
		FT	55.33	07.99	04.89	06.70	04.34
PointNet++ [40]	<b>LwF</b>	55.33	13.45	07.17	07.17	05.94	
		<b>Joint</b>	55.33	55.87	56.49	57.05	57.65
		FT	49.02	08.52	06.40	06.61	04.58
DGCNN [41]	<b>LwF</b>	49.02	17.75	<b>07.55</b>	<b>07.33</b>	<b>06.25</b>	
		<b>Joint</b>	49.02	50.91	50.10	51.04	51.23
		FT	51.16	07.74	04.85	06.44	04.29
PointMLP-E [43]	<b>LwF</b>	51.16	10.79	06.71	06.62	05.56	
		<b>Joint</b>	51.16	52.14	50.69	51.03	51.06
		FT	22.30	04.37	03.10	04.19	01.80
PCT [42]	<b>LwF</b>	22.30	04.96	04.26	04.43	02.58	
		<b>Joint</b>	22.30	22.09	21.26	22.53	25.79
		FT	54.65	07.97	04.84	06.72	03.76
PointNeXT [44]	<b>LwF</b>	54.84	<b>21.41</b>	<b>06.51</b>	<b>07.56</b>	<b>05.87</b>	
		<b>Joint</b>	54.43	55.80	56.15	57.78	58.28

## 5. Conclusion

This paper introduces the largest real-world 3D dataset on groceries called 3DGrocery100. One of the key contributions of this dataset is its wide and fine-grained variety of grocery categories. It contains 100 classes divided into 10, 24, 28, and 38: Apples, Fruits (non-apples), Vegetables, and Packages. High-resolution 10,755 RGB-D images were collected using mobile phones with 3D sensors and were processed to create the largest real-world point cloud dataset of 87,898 objects. In addition, this dataset is diverse due to the presence of point cloud objects of varying sizes under natural occlusions. We benchmarked six representative state-of-the-art methods on all five subsets and two color-based variants of our dataset. Our dataset stands out with its distinctive fine-grained features, making it an excellent benchmark for few-shot classification tasks, especially in strong generalization. Given the challenge of encountering novel classes, we focus on evaluating class-incremental learning approaches for classification. This evaluation will help us assess the dataset’s ability to handle the incorporation of new classes and maintain robust classification performance while building on the knowledge acquired from existing data.

## References

- [1] Daniel Sales, Diogo Correa, Fernando S Osório, and Denis F Wolf. 3d vision-based autonomous navigation system using ann and kinect sensor. In *Engineering Applications of Neural Networks: 13th International Conference, EANN 2012, London, UK, September 20-23, 2012. Proceedings 13*, pages 305–314. Springer, 2012. 1
- [2] Satya P Singh, Lipo Wang, Sukrit Gupta, Haveesh Goli, Parasuraman Padmanabhan, and Balázs Gulyás. 3D deep learning on medical images: a review. *Sensors*, 20(18):5097, 2020. 1
- [3] Max Krichenbauer, Goshiro Yamamoto, Takafumi Taketom, Christian Sandor, and Hirokazu Kato. Augmented reality versus virtual reality for 3d object manipulation. *IEEE transactions on visualization and computer graphics*, 24(2):1038–1048, 2017. 1
- [4] Charu C Aggarwal et al. *Data mining: the textbook*, volume 1. Springer, 2015. 1
- [5] Kostas Georgiadis, Fotis Kalaganis, Panagiotis Migkotzidis, Elisavet Chatzilari, Spiros Nikolopoulos, and Ioannis Komatsiari. A computer vision system supporting blind people—the supermarket case. In *Computer Vision Systems: 12th International Conference, ICVS 2019, Thessaloniki, Greece, September 23–25, 2019, Proceedings 12*, pages 305–315. Springer, 2019. 1, 3
- [6] Daniela Massiceti, Luisa Zintgraf, John Bronskill, Lida Theodorou, Matthew Tobias Harris, Edward Cutrell, Cecily Morrison, Katja Hofmann, and Simone Stumpf. Orbit: A real-world few-shot dataset for teachable object recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10818–10828, 2021. 1
- [7] Weidong Geng, Feilin Han, Jiangke Lin, Liuyi Zhu, Jieming Bai, Suzhen Wang, Lin He, Qiang Xiao, and Zhangjiong Lai. Fine-grained grocery product recognition by one-shot learning. In *Proceedings of the 26th ACM international conference on Multimedia*, pages 1706–1714, 2018. 1
- [8] Marco Leo, Pierluigi Carcagnì, and Cosimo Distante. A systematic investigation on end-to-end deep recognition of grocery products in the wild. In *2020 25th International Conference on Pattern Recognition (ICPR)*, pages 7234–7241. IEEE, 2021. 1
- [9] Carlos Ruiz, Joao Falcao, Shijia Pan, Hae Young Noh, and Pei Zhang. Aim3s: Autonomous inventory monitoring through multi-modal sensing for cashier-less convenience stores. In *Proceedings of the 6th ACM International Conference on Systems for Energy-Efficient Buildings, Cities, and Transportation*, pages 135–144, 2019. 1
- [10] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 1, 2
- [11] Qirui Yu, Huijun Yang, Yangbo Gao, Xinrui Ma, Guochao Chen, and Xin Wang. LFPNet: Lightweight network on real point sets for fruit classification and segmentation. *Computers and Electronics in Agriculture*, 194:106691, 2022. 1
- [12] D Koubaroulis, J Matas, J Kittler, and CTU CMP. Evaluating colour-based object recognition algorithms using the soil-47 database. In *Asian Conference on Computer Vision*, volume 2, 2002. 2, 3
- [13] Jiri George Matas, Dimitri Koubaroulis, and Josef Kittler. Colour image retrieval and object recognition using the multimodal neighbourhood signature. In *Computer Vision-ECCV 2000: 6th European Conference on Computer Vision Dublin, Ireland, June 26–July 1, 2000 Proceedings, Part I 6*, pages 48–64. Springer, 2000. 2
- [14] Michele Merler, Carolina Galleguillos, and Serge Belongie. Recognizing groceries in situ using in vitro training data. In *2007 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8. IEEE, 2007. 2, 3
- [15] David G Lowe. Distinctive image features from scale-invariant keypoints. *International journal of computer vision*, 60:91–110, 2004. 2
- [16] Anderson Rocha, Daniel C Hauagge, Jacques Wainer, and Siome Goldenstein. Automatic fruit and vegetable classification from images. *Computers and Electronics in Agriculture*, 70(1):96–104, 2010. 2, 3
- [17] Marian George and Christian Floerkemeier. Recognizing products: A per-exemplar multi-label image classification approach. In *Computer Vision-ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part II 13*, pages 440–455. Springer, 2014. 2, 3
- [18] Güл Varol and Ridvan Salih Kuzu. Toward retail product recognition on grocery shelves. In *Sixth International Conference on Graphic and Image Processing (ICGIP 2014)*, volume 9443, pages 46–52. SPIE, 2015. 2
- [19] Philipp Jund, Nichola Abdo, Andreas Eitel, and Wolfram Burgard. The freiburg groceries dataset. *arXiv preprint arXiv:1611.05799*, 2016. 2, 3
- [20] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio Guadarrama, and Trevor Darrell. Caffe: Convolutional architecture for fast feature embedding. In *Proceedings of the 22nd ACM international conference on Multimedia*, pages 675–678, 2014. 2
- [21] Patrick Follmann, Tobias Bottger, Philipp Hartinger, Rebecca Konig, and Markus Ulrich. MVTec D2S: densely segmented supermarket dataset. In *Proceedings of the European conference on computer vision (ECCV)*, pages 569–585, 2018. 2, 3
- [22] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. 2
- [23] Marcus Klasson, Cheng Zhang, and Hedvig Kjellström. A hierarchical grocery store image dataset with visual and semantic labels. In *2019 IEEE winter conference on applications of computer vision (WACV)*, pages 491–500. IEEE, 2019. 2, 3, 4
- [24] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017. 2
- [25] Xiu-Shen Wei, Quan Cui, Lei Yang, Peng Wang, and Lingqiao Liu. RPC: A large-scale retail product checkout dataset. *arXiv preprint arXiv:1901.07249*, 2019. 2, 3

- [26] Ping Hu, Weiqiang Wang, Chi Zhang, and Ke Lu. Detecting salient objects via color and texture compactness hypotheses. *IEEE Transactions on Image Processing*, 25(10):4653–4664, 2016. 2
- [27] Philipp Krähenbühl and Vladlen Koltun. Efficient inference in fully connected crfs with gaussian edge potentials. *Advances in neural information processing systems*, 24, 2011. 2
- [28] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232, 2017. 2
- [29] Marco Filax, Tim Gonschorek, and Frank Ortmeier. Data for Image Recognition Tasks: An Efficient Tool for Fine-Grained Annotations. In *ICPRAM*, pages 900–907, 2019. 2, 3
- [30] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 2
- [31] Yu Hao, Yanwei Fu, and Yu-Gang Jiang. Take goods from shelves: A dataset for class-incremental object detection. In *Proceedings of the 2019 on International Conference on Multimedia Retrieval*, pages 271–278, 2019. 2, 3
- [32] Eran Goldman, Roei Herzig, Aviv Eisenshtat, Jacob Goldberger, and Tal Hassner. Precise detection in densely packed scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5227–5236, 2019. 2
- [33] Jingtian Peng, Chang Xiao, and Yifan Li. RP2K: A large-scale retail product dataset for fine-grained image classification. *arXiv preprint arXiv:2006.12634*, 2020. 2, 3
- [34] Arjun Singh, James Sha, Karthik S Narayan, Tudor Achim, and Pieter Abbeel. Bigbird: A large-scale 3d database of object instances. In *2014 IEEE international conference on robotics and automation (ICRA)*, pages 509–516. IEEE, 2014. 2, 3
- [35] Stephen Tyree, Jonathan Tremblay, Thang To, Jia Cheng, Terry Mosier, Jeffrey Smith, and Stan Birchfield. 6-DoF pose estimation of household objects for robotic manipulation: An accessible dataset and benchmark. In *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 13081–13088. IEEE, 2022. 2, 3
- [36] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3d objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13142–13153, 2023. 2
- [37] Shivanand Venkanna Sheshappanavar and Chandra Kambhamettu. Local Neighborhood Features for 3D Classification. In *Scandinavian Conference on Image Analysis*, pages 386–395. Springer, 2023. 2
- [38] Shivanand Venkanna Sheshappanavar. *Learning from Neighborhoods for 3D Point Cloud Classification*. PhD thesis, University of Delaware, 2023. 2
- [39] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 652–660, 2017. 2, 3, 6, 7, 8, 9, 10, 11, 12
- [40] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *Advances in neural information processing systems*, 30, 2017. 3, 6, 7, 8, 9, 10, 11, 12
- [41] 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. 3, 6, 7, 8, 9, 10, 11, 12
- [42] Meng-Hao Guo, Jun-Xiong Cai, Zheng-Ning Liu, Tai-Jiang Mu, Ralph R Martin, and Shi-Min Hu. Pct: Point cloud transformer. *Computational Visual Media*, 7:187–199, 2021. 3, 6, 7, 8, 9, 10, 11, 12
- [43] Xu Ma, Can Qin, Haoxuan You, Haoxi Ran, and Yun Fu. Rethinking network design and local geometry in point cloud: A simple residual MLP framework. *arXiv preprint arXiv:2202.07123*, 2022. 3, 6, 7, 8, 9, 10, 11, 12
- [44] Guocheng Qian, Yuchen Li, Houwen Peng, Jinjie Mai, Hasan Hammoud, Mohamed Elhoseiny, and Bernard Ghanem. Pointnext: Revisiting pointnet++ with improved training and scaling strategies. *Advances in Neural Information Processing Systems*, 35:23192–23204, 2022. 2, 3, 6, 7, 8, 1, 9, 10, 11, 12
- [45] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Lin-guang Zhang, Xiaou Tang, and Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shapes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1912–1920, 2015. 2, 3, 7, 12
- [46] Mikaela Angelina Uy, Quang-Hieu Pham, Binh-Son Hua, Thanh Nguyen, and Sai-Kit Yeung. Revisiting point cloud classification: A new benchmark dataset and classification model on real-world data. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1588–1597, 2019. 2, 3, 6, 7, 8, 10, 11, 12
- [47] Garrick Brazil, Abhinav Kumar, Julian Straub, Nikhila Ravi, Justin Johnson, and Georgia Gkioxari. Omni3D: A Large Benchmark and Model for 3D Object Detection in the Wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13154–13164, 2023. 2
- [48] You Li and Javier Ibanez-Guzman. Lidar for autonomous driving: The principles, challenges, and trends for automotive lidar and perception systems. *IEEE Signal Processing Magazine*, 37(4):50–61, 2020. 2
- [49] Deepu Rajan and Subhasis Chaudhuri. Simultaneous estimation of super-resolved scene and depth map from low resolution defocused observations. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25(9):1102–1117, 2003. 2
- [50] Zhengyou Zhang. Microsoft Kinect Sensor and Its Effect. *IEEE MultiMedia*, 19(2):4–10, 2012. 2
- [51] Duarte Fernandes, António Silva, Rafael Névoa, Cláudia Simões, Dibet Gonzalez, Miguel Guevara, Paulo Novais, João Monteiro, and Pedro Melo-Pinto. Point-cloud based 3D object detection and classification methods for self-driving

- applications: A survey and taxonomy. *Information Fusion*, 68:161–191, 2021. 2
- [52] Shivanand Venkanna Sheshappanavar, Vinit Veerendraveer Singh, and Chandra Kambhamettu. Patchaugment: Local neighborhood augmentation in point cloud classification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2118–2127, 2021. 2
- [53] Wendy P Fernandcz, Yang Xian, and Yingli Tian. Image-based barcode detection and recognition to assist visually impaired persons. In *2017 IEEE 7th Annual International Conference on CYBER Technology in Automation, Control, and Intelligent Systems (CYBER)*, pages 1241–1245. IEEE, 2017. 3
- [54] Michael McCloskey and Neal J Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation*, volume 24, pages 109–165. Elsevier, 1989. 3, 4, 8
- [55] Ruud M Bolle, Jonathan H Connell, Norman Haas, Rakesh Mohan, and Gabriel Taubin. Veggievision: A produce recognition system. In *Proceedings Third IEEE Workshop on Applications of Computer Vision. WACV'96*, pages 244–251. IEEE, 1996. 3
- [56] Bikash Santra and Dipti Prasad Mukherjee. A comprehensive survey on computer vision based approaches for automatic identification of products in retail store. *Image and Vision Computing*, 86:45–63, 2019. 3
- [57] Kevin Lai, Liefeng Bo, Xiaofeng Ren, and Dieter Fox. A large-scale hierarchical multi-view rgb-d object dataset. In *2011 IEEE international conference on robotics and automation*, pages 1817–1824. IEEE, 2011. 3
- [58] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 3
- [59] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. *Advances in neural information processing systems*, 30, 2017. 4, 7, 11
- [60] Yinbo Chen, Zhuang Liu, Huijuan Xu, Trevor Darrell, and Xiaolong Wang. Meta-baseline: Exploring simple meta-learning for few-shot learning. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9062–9071, 2021. 4
- [61] Xu Luo, Jing Xu, and Zenglin Xu. Channel importance matters in few-shot image classification. In *International conference on machine learning*, pages 14542–14559. PMLR, 2022. 4
- [62] Tejas Anvekar and Dena Bazazian. GPr-Net: Geometric Prototypical Network for Point Cloud Few-Shot Learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pages 4178–4187, June 2023. 4, 7
- [63] Mozhgan PourKeshavarzi, Guoying Zhao, and Mohammad Sabokrou. Looking back on learned experiences for class/task incremental learning. In *International Conference on Learning Representations*, 2021. 4
- [64] Matthew McLeod, Chunlok Lo, Matthew Schlegel, Andrew Jacobsen, Raksha Kumaraswamy, Martha White, and Adam White. Continual auxiliary task learning. *Advances in Neural Information Processing Systems*, 34:12549–12562, 2021.
- [65] Davide Abati, Jakub Tomczak, Tijmen Blankevoort, Simone Calderara, Rita Cucchiara, and Babak Ehteshami Bejnordi. Conditional channel gated networks for task-aware continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3931–3940, 2020.
- [66] Arthur Douillard, Matthieu Cord, Charles Ollion, Thomas Robert, and Eduardo Valle. Podnet: Pooled outputs distillation for small-tasks incremental learning. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XX 16*, pages 86–102. Springer, 2020. 4
- [67] Fei Ye and Adrian G Bors. Learning latent representations across multiple data domains using lifelong VAEGAN. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XX 16*, pages 777–795. Springer, 2020. 4
- [68] Christian Simon, Masoud Faraki, Yi-Hsuan Tsai, Xiang Yu, Samuel Schulter, Yumin Suh, Mehrtash Harandi, and Manmohan Chandraker. On generalizing beyond domains in cross-domain continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9265–9274, 2022. 4
- [69] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl: Incremental classifier and representation learning. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 2001–2010, 2017. 4
- [70] Chen He, Ruiping Wang, Shiguang Shan, and Xilin Chen. Exemplar-supported generative reproduction for class incremental learning. In *BMVC*, page 98, 2018.
- [71] Eden Belouadah and Adrian Popescu. Il2m: Class incremental learning with dual memory. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 583–592, 2019. 4
- [72] Liyang Liu, Zhanghui Kuang, Yimin Chen, Jing-Hao Xue, Wenming Yang, and Wayne Zhang. IncDet: In Defense of Elastic Weight Consolidation for Incremental Object Detection. *IEEE Transactions on Neural Networks and Learning Systems*, 32(6):2306–2319, 2021. 4
- [73] Townim Chowdhury, Mahira Jalisha, Ali Cheraghian, and Shafin Rahman. Learning without forgetting for 3d point cloud objects. In *Advances in Computational Intelligence: 16th International Work-Conference on Artificial Neural Networks, IWANN 2021, Virtual Event, June 16–18, 2021, Proceedings, Part I 16*, pages 484–497. Springer, 2021. 4
- [74] Capturing Photos with Depth. [https://developer.apple.com/documentation/avfoundation/additional\\_data\\_capture/capturing\\_photos\\_with\\_depth](https://developer.apple.com/documentation/avfoundation/additional_data_capture/capturing_photos_with_depth). 4
- [75] LabelMe. <https://github.com/wkentaro/labelme>. 5
- [76] Qian-Yi Zhou, Jaesik Park, and Vladlen Koltun. Open3D: A modern library for 3D data processing. *arXiv preprint arXiv:1801.09847*, 2018. 5

- [77] Marie-Julie Rakotosaona, Vittorio La Barbera, Paul Guerero, Niloy J Mitra, and Maks Ovsjanikov. Pointcleannet: Learning to denoise and remove outliers from dense point clouds. In *Computer graphics forum*, volume 39, pages 185–203. Wiley Online Library, 2020. 5, 6, 2
- [78] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021. 7, 10, 11
- [79] Eleni Triantafillou, Tyler Zhu, Vincent Dumoulin, Pascal Lamblin, Utku Evci, Kelvin Xu, Ross Goroshin, Carles Gelada, Kevin Swersky, Pierre-Antoine Manzagol, et al. Meta-dataset: A dataset of datasets for learning to learn from few examples. *arXiv preprint arXiv:1903.03096*, 2019. 7
- [80] Charu Sharma and Manohar Kaul. Self-supervised few-shot learning on point clouds. *Advances in Neural Information Processing Systems*, 33:7212–7221, 2020. 7
- [81] Hengxin Feng, Weifeng Liu, Yanjiang Wang, and Baodi Liu. Enrich features for few-shot point cloud classification. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2285–2289. IEEE, 2022.
- [82] Chuangguan Ye, Hongyuan Zhu, Yongbin Liao, Yanggang Zhang, Tao Chen, and Jiayuan Fan. What makes for effective few-shot point cloud classification? In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 1829–1838, 2022.
- [83] Chuangguan Ye, Hongyuan Zhu, Bo Zhang, and Tao Chen. A Closer Look at Few-Shot 3D Point Cloud Classification. *International Journal of Computer Vision*, 131(3):772–795, 2023. 7, 10
- [84] Mohamed Afham, Isuru Dissanayake, Dinithi Dissanayake, Amaya Dharmasiri, Kanchana Thilakarathna, and Ranga Rodriguez. Crosspoint: Self-supervised cross-modal contrastive learning for 3d point cloud understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9902–9912, 2022. 7
- [85] Shivanand Kundargi, Tejas Anvekar, Ramesh Ashok Tabib, and Uma Mudenagudi. PointCLIMB: An Exemplar-Free Point Cloud Class Incremental Benchmark, 2023. 7
- [86] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International conference on machine learning*, pages 1126–1135. PMLR, 2017. 10

# A Benchmark Grocery Dataset of Realworld Point Clouds From Single View

## Supplementary Material

We show the use of the packages subset dataset to pre-train the PointNeXt model to achieve the state-of-the-art result on the hardest variant of the ScanObjectNN dataset. We provide counts of RGB-D images for each class. We show five visual samples from each of the 100 classes. We show the confusion matrices from evaluations of our dataset on six state-of-the-art object classification models.

### 6. Packages subset for Pre-training

Unlike Fruits or Vegetables, Packages' shapes are unique because they do not show much deformation at each instance level. For example, the shape of beans or Gatorade bottles has consistency across multiple instances, whereas different broccoli samples have different shapes. This shape uniqueness of packages and high performance of package classification, even in the absence of colors, make the subset dataset suitable for pretraining methods upon which small datasets can be fine-tuned to achieve a gain in classification accuracy as shown in Table 6.

Table 6. Pre-training PointNeXt [44] model on the packages subset (w/o colors) and fine-tuning with the ScanObjectNN's hardest variant improves classification accuracy by 0.64%.

Method	Acc. (%)
PointNeXt [44]	87.70
PointNeXt w/ pre-training	88.34 ( $\uparrow$ <b>0.64</b> )

### 7. Number of RGB-D Images

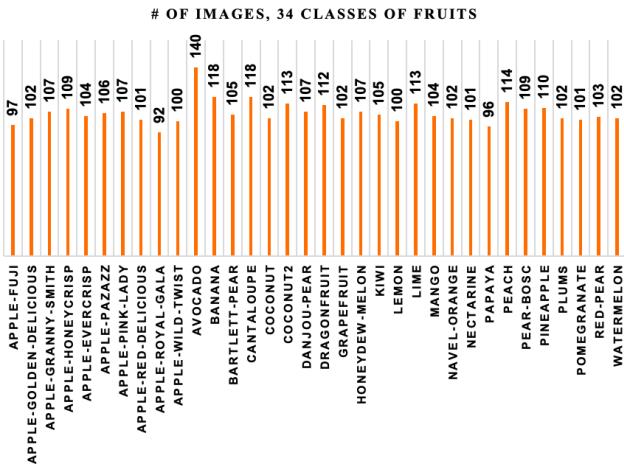


Figure 5. 34 Fruit Classes: each with the count of images.

Our dataset consists of 10,755 RGB-D images spread across 100 classes of groceries, as shown in Figure 5 for fruits, Figure 6 for vegetables, and Figure 7 for packages. Each of these images is annotated.

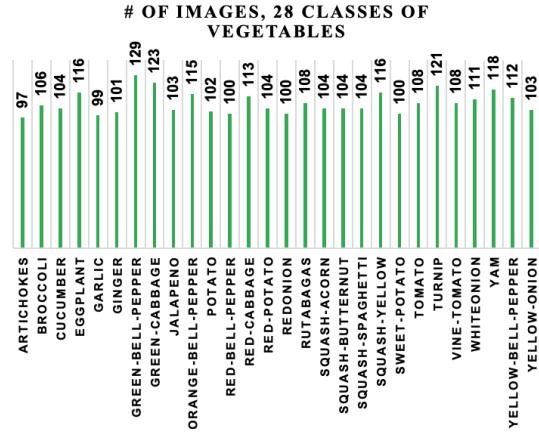


Figure 6. 28 Vegetable Classes: each with the count of images.



Figure 7. 38 Package Classes: each with the count of images.

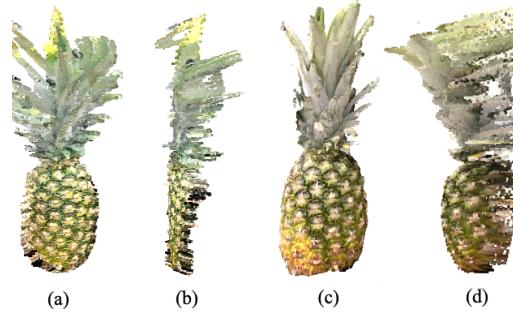


Figure 8. Comparison between LiDAR (a) front view (b) side view and Stereo Vision (c) front view (d) side view of pineapple instance after conversion to a point cloud.

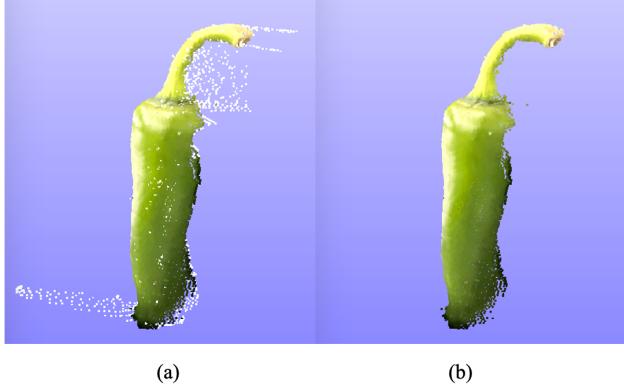


Figure 9. Jalapeno object example: a) point cloud with outliers (colored with white for visibility) b) point cloud after removing outliers using PointCleanNet [77] with a threshold of 0.4.

## 8. Visual samples of 100 classes

Figure 10 shows five samples from each of the ten apple classes. Each sample consists of 1024 points sampled using the farthest point sampling method. Figures 11, 12, and 13 together show five samples from each of the 24 non-apple fruit classes. Figures 14, 15, and 16 together show five samples from each of the 28 vegetable classes. Figures 17, 18, 19, 20 and 21 together show five samples from each of the 38 package classes. Zoom in to a sample for better visibility of 3D points and colors associated with the points.

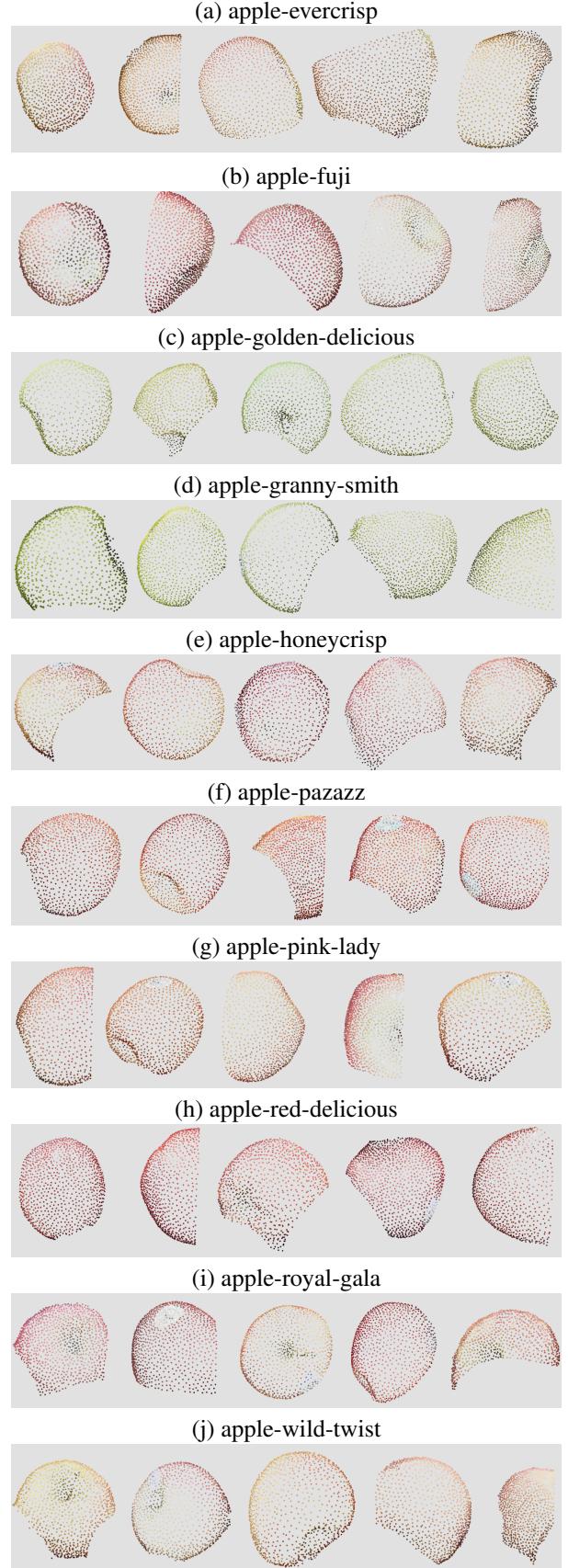


Figure 10. Visual samples (1024 points) from 10 Apple classes. Labels on top of the objects. Zoom in for better visibility.

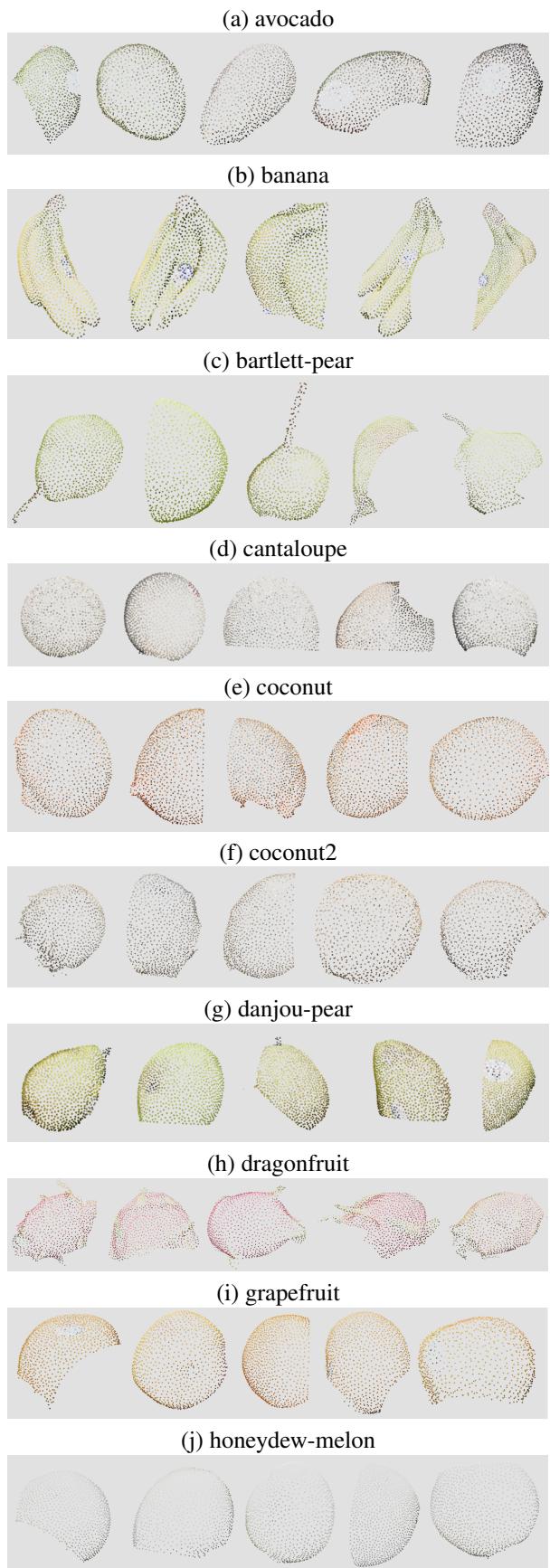


Figure 11. Visual samples (1024 points) from 10 of 24 non-apple fruit classes. Labels on top of the objects.

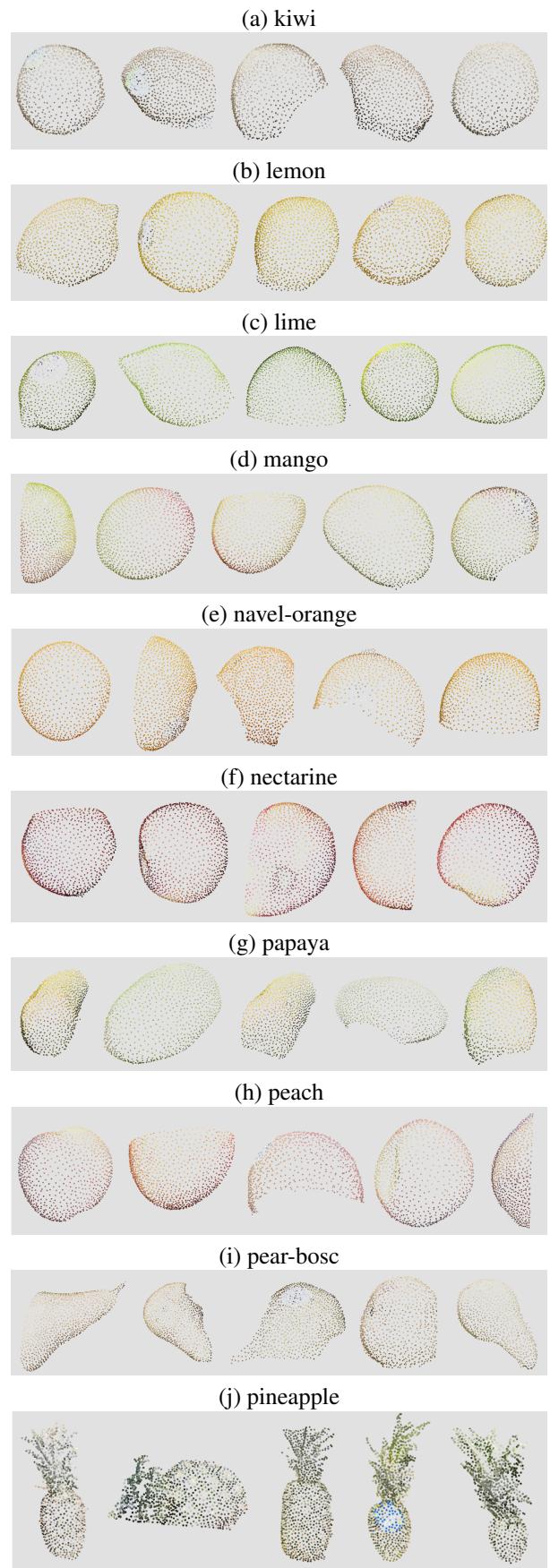


Figure 12. Visual samples (1024 points) from 10 of 24 non-apple fruit classes. Labels on top of the objects. Zoom in for better visibility.

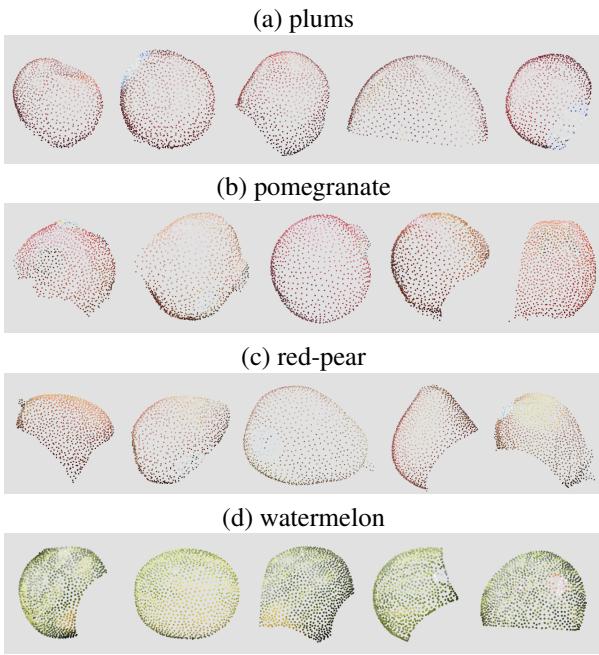


Figure 13. Visual samples (1024 points) from 4 or 24 non-apple fruit classes. Labels on top of the objects. Zoom in for better visibility.

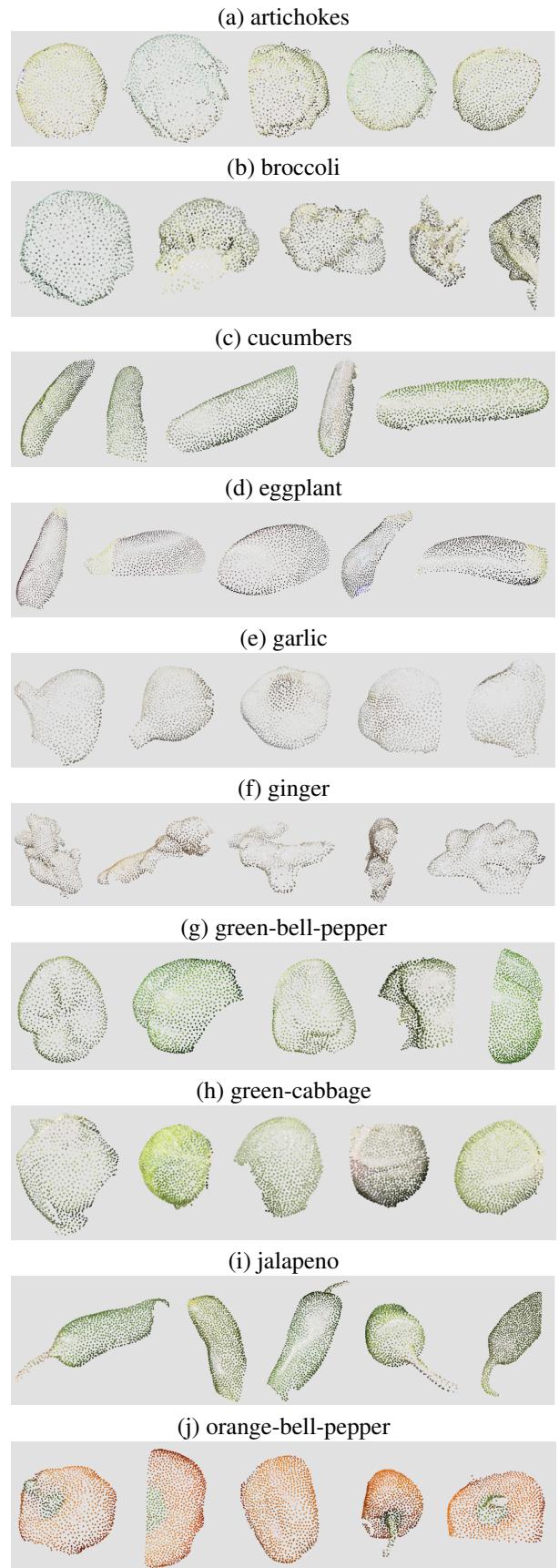


Figure 14. Visual samples (1024 points) from 10 vegetable classes. Labels on top of the objects. Zoom in for better visibility.

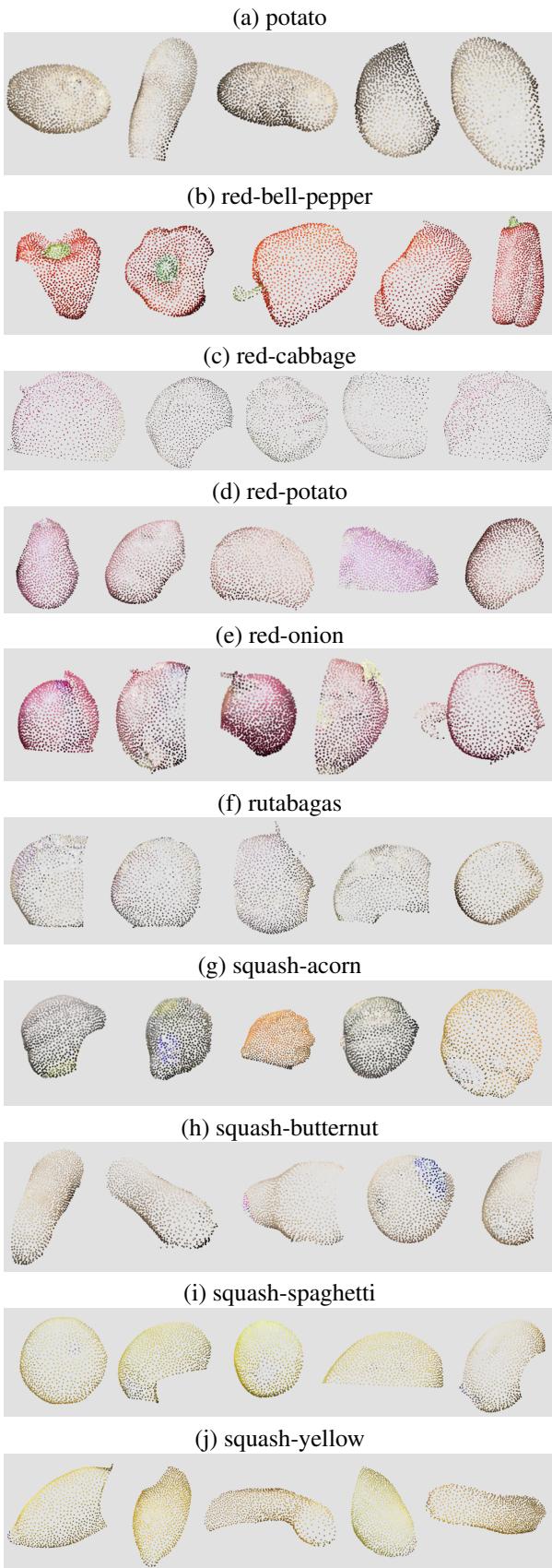


Figure 15. Visual samples (1024 points) from 10 vegetable classes. Labels on top of the objects. Zoom in for better visibility.

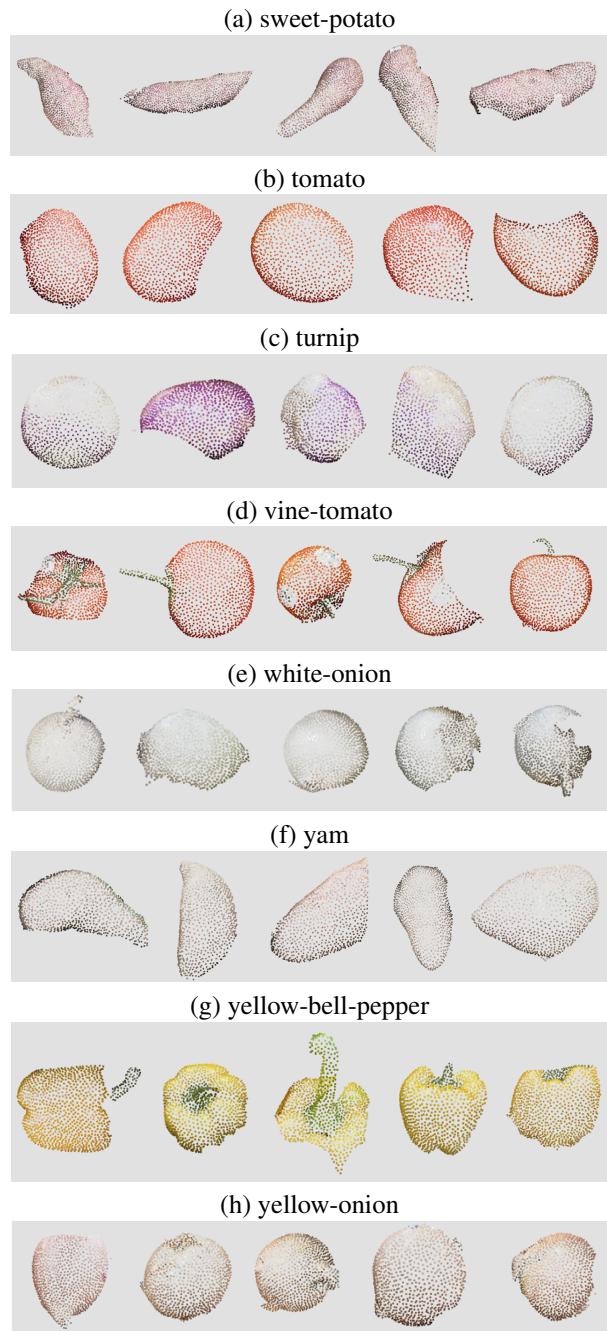


Figure 16. Visual samples (1024 points) from 8 vegetable classes. Labels on top of the objects. Zoom in for better visibility.

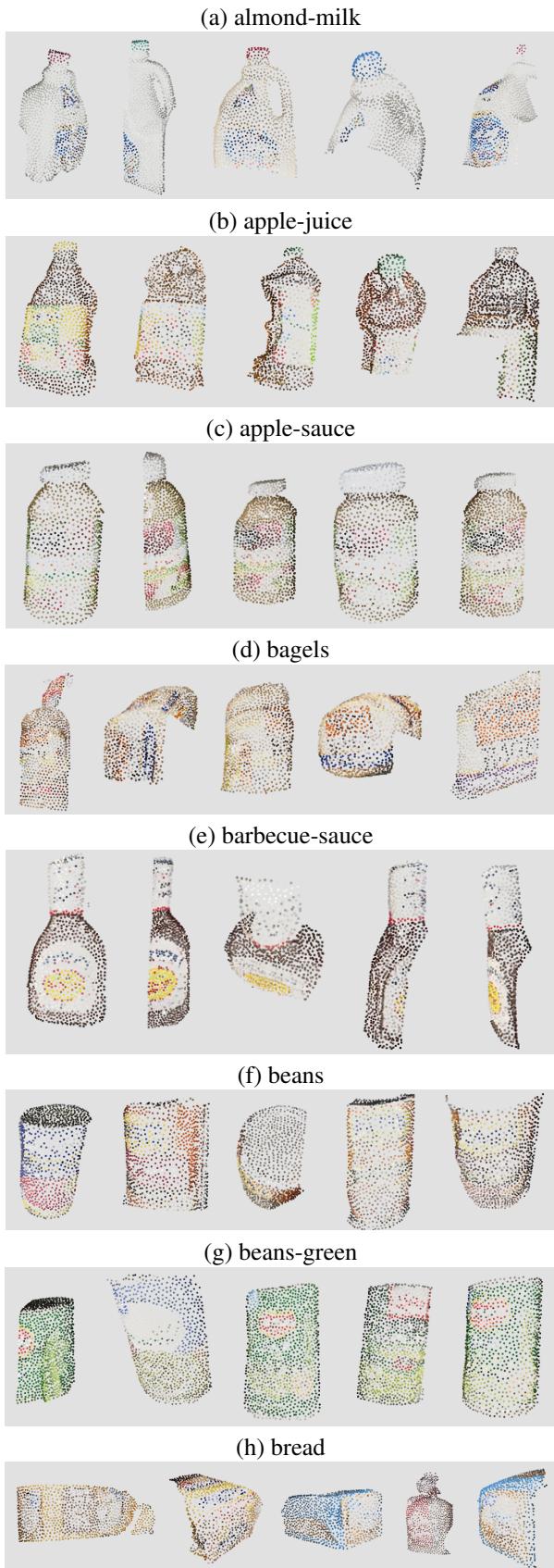


Figure 17. Visual samples (1024 points) from 8 package classes. Labels on top of the objects. Zoom in for better visibility.

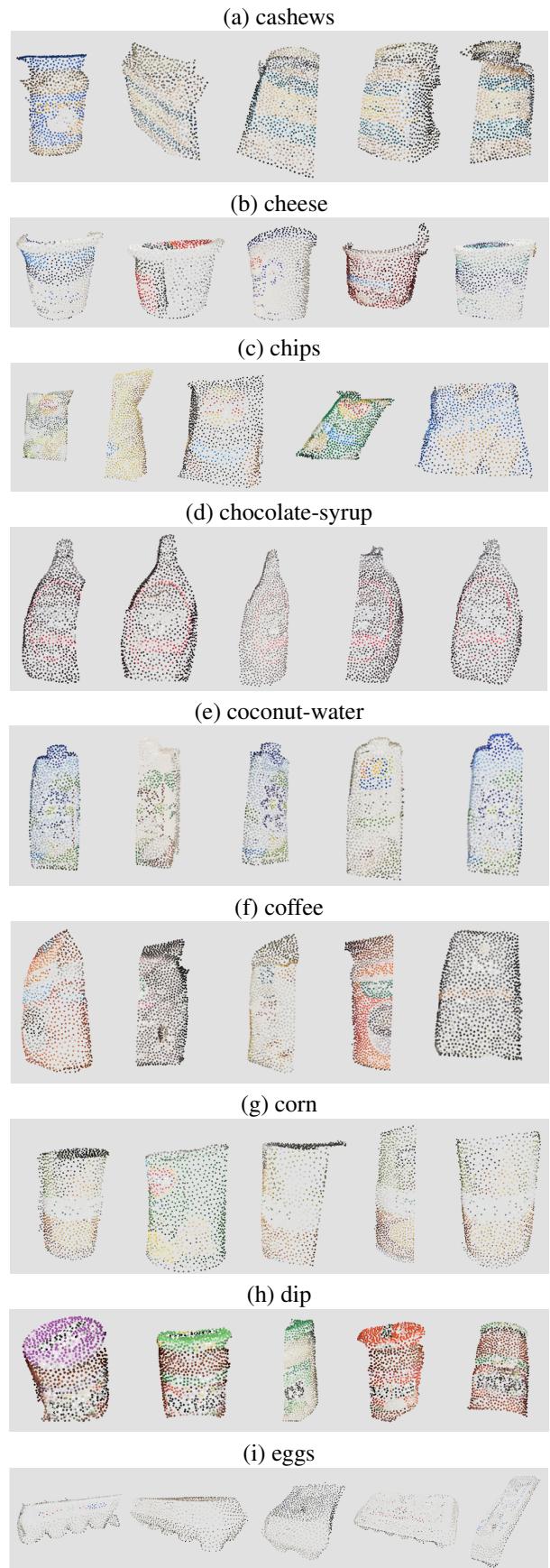


Figure 18. Visual samples (1024 points) from 9 package classes. Labels on top of the objects. Zoom in for better visibility.

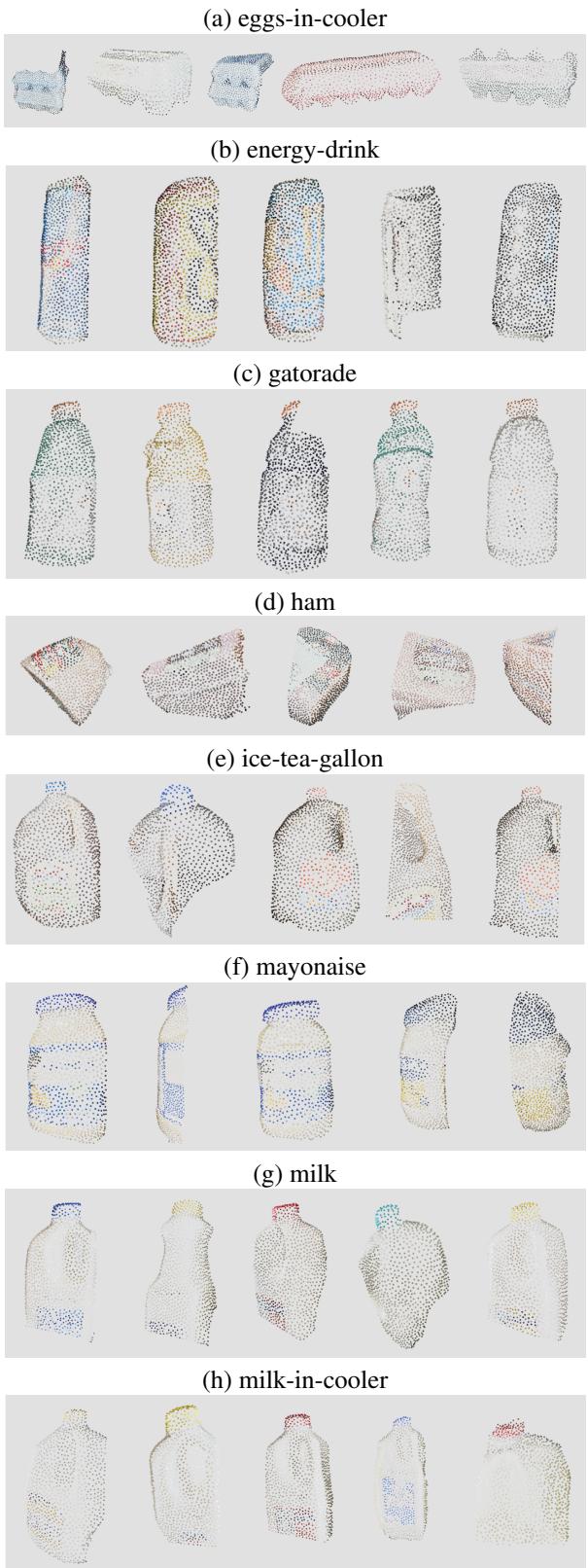


Figure 19. Visual samples (1024 points) from 8 package classes. Labels on top of the objects. Zoom in for better visibility.

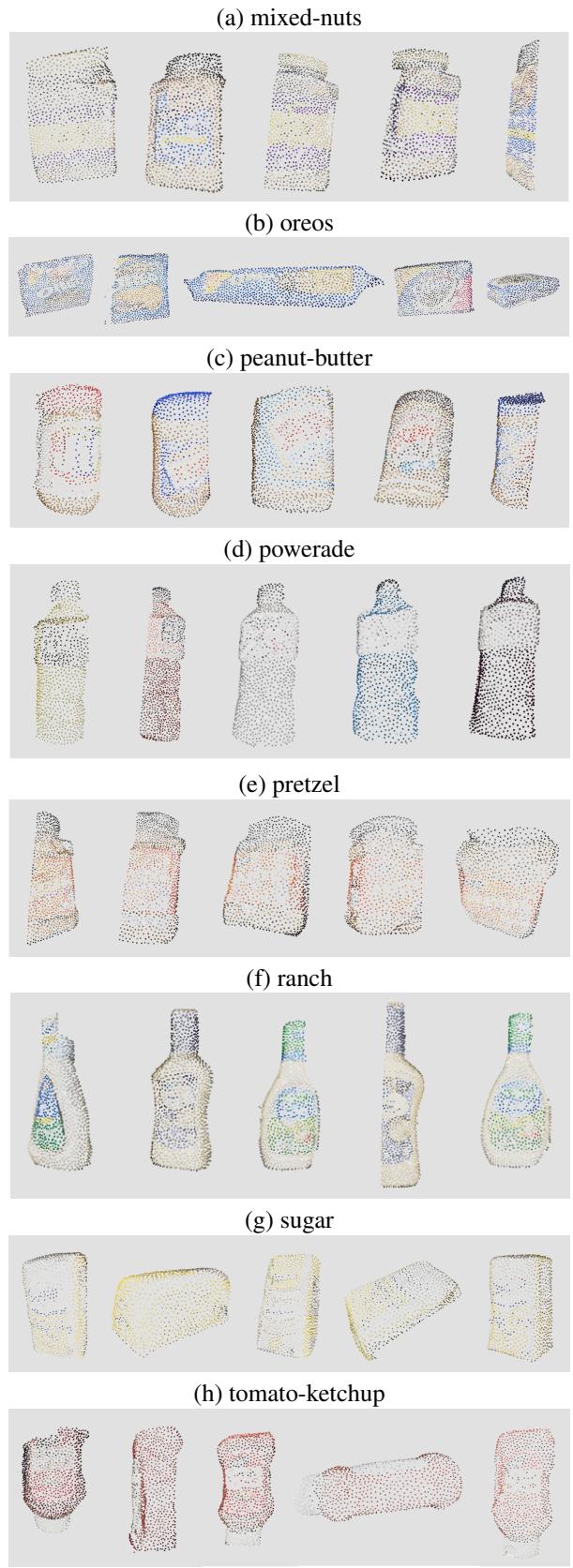


Figure 20. Visual samples (1024 points) from 8 package classes. Labels on top of the objects. Zoom in for better visibility.

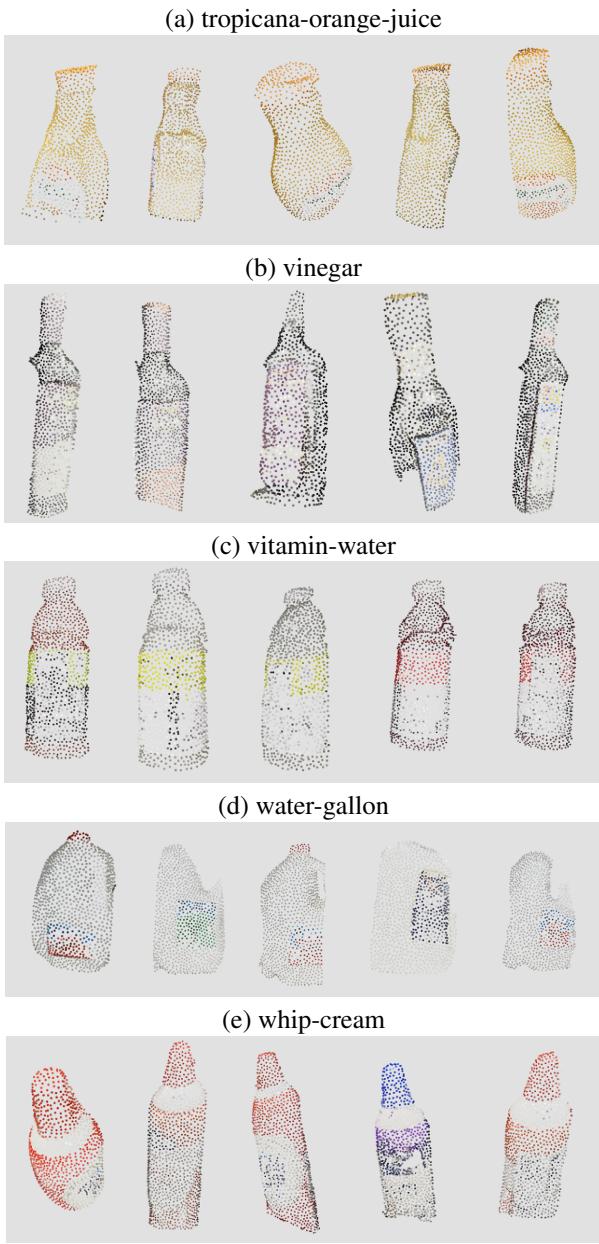


Figure 21. Visual samples (1024 points) from 5 package classes. Labels on top of the objects. Zoom in for better visibility.

## 9. 3D Point cloud Classification - detailed results

Figure 3 shows 3D samples from each of the ten apple classes with colors (a-j) and the respective 3D apples without colors(k-t). While apples with colors are often visibly distinguishable, it is challenging to recognize them without colors. Also, the appearance of apples due to color is quite similar between classes, while the shapes of different apples in the same class vary. For example, in the apple10 subset,

eight of the ten apple classes are light red to dark red, posing significant challenges to the models in learning discriminative features. Unlike Fruits and Vegetables, Packages often have a fixed shape for a given class and do not suffer surface deformations for different instances of the same class.

Several Packages, such as tomato ketchup, ranch, and mayonnaise, are stacked straight and upside down in a rack. In the case of fruits and vegetables, there is no specific order of arrangement. Some stores where the fruits and vegetables are placed in racks suffered from relatively proper illumination compared to objects at the front of the rack leading to variations in the amount of lighting received.

Table 7 shows the class-wise accuracy of the Apple10 subset, both with and without colors. Table 8 lists the class-wise accuracy of 24 fruit classes. Since apples without color are hard to distinguish from other apples, we do not include apples with other fruits in this experiment. Tables 9 and 10 show 3D point cloud classification results of vegetables and packages, respectively. Each table contains the number of point clouds used as train and test samples for each class. The maximum accuracy for each class without colors is highlighted in bold. Each model trained on point clouds with colors achieved 100% accuracy on multiple classes. Here we list some of our findings from those confusion matrices. From the confusion matrix of each state-of-the-art method on our full dataset without colors, we observe that none of the classes from packages (38 classes) is misclassified as a non-package class (62 classes - fruits and vegetables) or vice-versa. Table 8 consists of 3D Fruits (non-apple classes) results in both with and without colors on the six methods [39–44]. Colors as features contribute richly to learning discriminative features. In the absence of colors, fruits with a similar shape, such as watermelon, honeydew-melon, and cantaloupe, are misclassified amongst each other.

Table 9 shows 3D Vegetables results with and without colors on each of the six methods. The Vegetables without colors showed improved classification results compared to Fruits without colors. However, a few ambiguous misclassifications are observed. e.g., tomato as red-potato. For example, without colors, a tiny percentage of tomatoes are misclassified as vine-tomato attributed because of the red color of the objects. All six models in our experiments achieve higher classification performance without colors as additional features. Among the six state-of-the-art methods, PointNet++ [39] achieves the highest classification accuracy on the Packages (without colors).

## 10. Ablations

### 10.1. Quality vs. Quantity

We test the robustness of state-of-the-art models by giving fewer input points, i.e., 512, 256, and 128 points, dur-







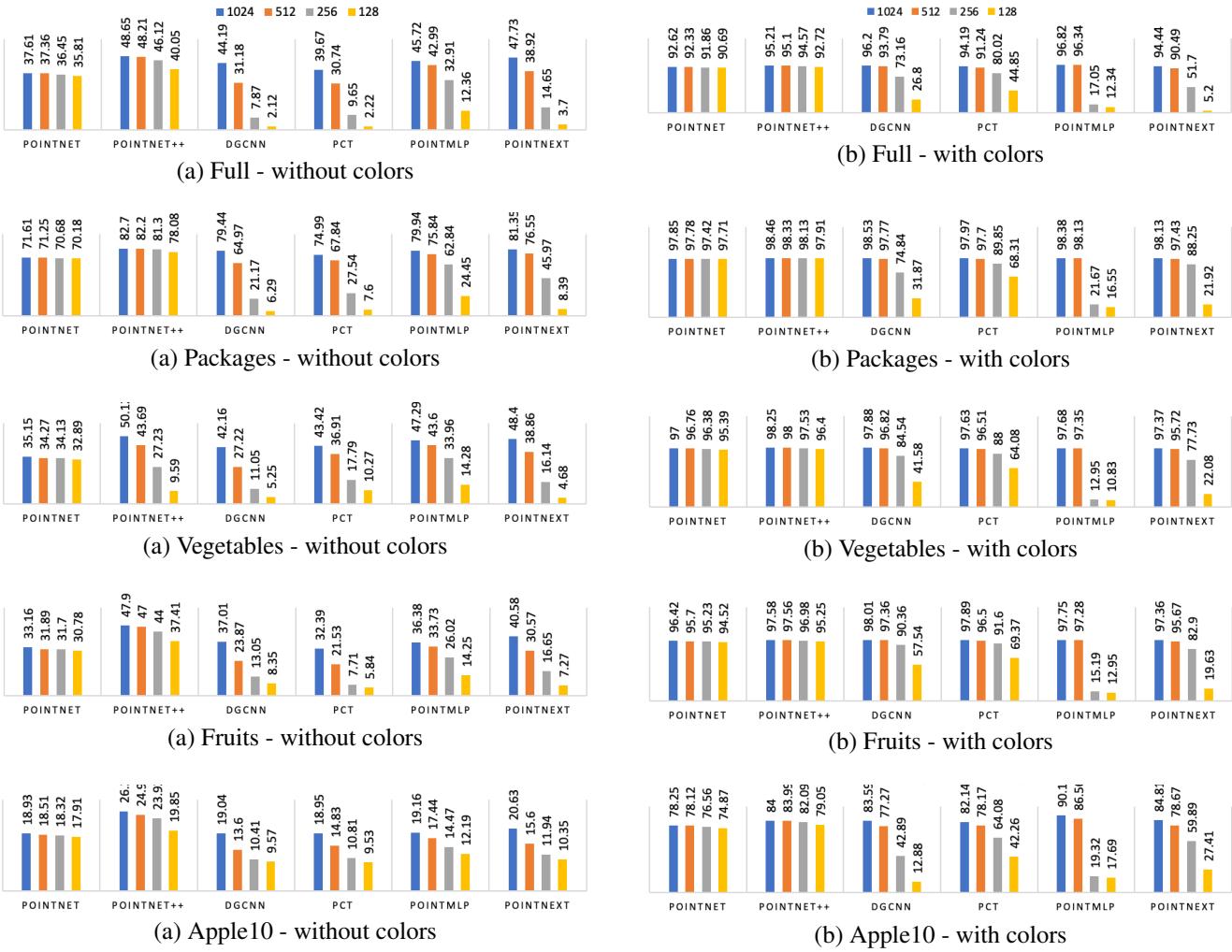


Figure 22. Performance of state-of-the-art methods with a different number of input points, i.e., 1024 (blue), 512 (orange), 256 (grey), and 128 (yellow).

Table 12. Quantitative analysis for Few-shot 3D point cloud classification evaluation on ScanObjectNN [46] dataset. **Note:** all models are pre-trained on ModelNet40 [45]

Model	Weight Init	5-way		10-way	
		10-shots	20-shots	10-shots	20-shots
<b>DGCNN[41]</b>	Random	67.60	$\pm 8.35$	74.46	$\pm 7.38$
	MN40	79.86	$\pm 6.67$	85.10	$\pm 5.76$
<b>PointMLP[43]</b>	Random	33.20	$\pm 6.11$	37.86	$\pm 7.63$
	MN40	<b>81.60</b>	$\pm 6.03$	<b>85.78</b>	$\pm 6.38$
<b>PointNet++[40]</b>	Random	67.68	$\pm 8.92$	72.98	$\pm 7.46$
	MN40	77.02	$\pm 6.98$	80.86	$\pm 5.81$
<b>PointNet[39]</b>	Random	69.78	$\pm 8.56$	74.76	$\pm 7.65$
	MN40	76.56	$\pm 8.26$	80.34	$\pm 7.48$
<b>PCT[42]</b>	Random	67.73	$\pm 7.72$	73.83	$\pm 7.33$
	MN40	72.76	$\pm 8.73$	79.54	$\pm 6.95$
<b>PointNeXT[44]</b>	Random	69.03	$\pm 8.66$	74.52	$\pm 7.40$
	MN40	79.02	$\pm 6.22$	83.94	$\pm 5.69$

Table 13. Quantitative analysis for Few-shot 3D point cloud classification on ScanObjectNN[46] for Weak generalizations evaluation. **Note:** W1 - Weak 1 (ONLY OBJ split), W2 - Weak 2 (OBJ + BG split), and W3 - Weak 3 (PB75 split).

Model	Split	5-ways		10-ways	
		5-shots	10-shots	5-shots	10-shots
<b>PointNet [39]</b>	W1	56.32	$\pm 0.64$	59.91	$\pm 0.61$
	W2	54.88	$\pm 0.65$	57.44	$\pm 0.61$
	W3	49.29	$\pm 0.59$	51.73	$\pm 0.57$
<b>PointNet++ [40]</b>	W1	51.35	$\pm 0.61$	58.19	$\pm 0.64$
	W2	52.04	$\pm 0.60$	58.13	$\pm 0.63$
	W3	47.09	$\pm 0.60$	51.44	$\pm 0.60$
<b>DGCNN [41]</b>	W1	46.43	$\pm 0.56$	52.97	$\pm 0.59$
	W2	56.47	$\pm 0.65$	<b>61.26</b>	$\pm 0.64$
	W3	51.00	$\pm 0.65$	55.14	$\pm 0.61$
<b>PointMLP [43]</b>	W1	47.67	$\pm 0.63$	52.04	$\pm 0.57$
	W2	48.10	$\pm 0.63$	52.94	$\pm 0.60$
	W3	42.61	$\pm 0.61$	45.44	$\pm 0.60$
<b>PointTransformer [42]</b>	W1	57.94	$\pm 0.64$	57.07	$\pm 0.66$
	W2	59.70	$\pm 0.65$	59.48	$\pm 0.66$
	W3	53.52	$\pm 0.64$	53.20	$\pm 0.60$
<b>PointNeXT [44]</b>	W1	<b>60.17</b>	$\pm 0.63$	<b>61.21</b>	$\pm 0.63$
	W2	<b>60.88</b>	$\pm 0.64$	59.92	$\pm 0.67$
	W3	54.84	$\pm 0.64$	55.69	$\pm 0.65$