

3D-COCO: EXTENSION OF MS-COCO DATASET FOR SCENE UNDERSTANDING AND 3D RECONSTRUCTION

BIDEAUX Maxence, PHE Alice, CHAOUCH Mohamed, LUVISON Bertrand, PHAM Quoc-Cuong

Université Paris-Saclay, CEA, List, F-91120, Palaiseau, France

ABSTRACT

We introduce 3D-COCO, an extension of the original MS-COCO [1] dataset providing 3D models and 2D-3D alignment annotations. 3D-COCO was designed to achieve computer vision tasks such as 3D reconstruction or image detection configurable with textual, 2D image, and 3D CAD model queries. We complete the existing MS-COCO [1] dataset with 28K 3D models collected on ShapeNet [2] and Objaverse [3]. By using an IoU-based method, we match each MS-COCO [1] annotation with the best 3D models to provide a 2D-3D alignment. The open-source nature of 3D-COCO is a premiere that should pave the way for new research on 3D-related topics. The dataset and its source codes is available at <https://kalisteo.cea.fr/index.php/coco3d-object-detection-and-reconstruction/>

Index Terms— Dataset, Detection, Reconstruction, 3D models, 2D-3D alignment

1. INTRODUCTION

For almost a decade, object detection has become a central topic in computer vision. This growing interest finds its roots in the new challenges of autonomous driving, crowd counting, anomaly detection, and smart video surveillance. As a result, many innovative neural networks such as Faster R-CNN [4], YOLO [5], SSD [6], and DETR [7] have been developed over the years. Most of these architectures’ performances are evaluated and compared through some widespread datasets like Pascal VOC [8], Open Images [9] and MS-COCO [1].

Iterations were carried out to improve these architectures, thus making it possible to achieve optimal performance for the detection of objects appearing during training. A new area of research involves the detection of new semantic classes that do not appear during training. This innovation would allow object detectors to meet a wider field of application without the need for re-training. As an example, OV-DETR [10] uses

a foundation model as a backbone to transform the DETR [7] architecture into an open vocabulary detector configurable with text and images. Then, it could be interesting to develop detection networks configurable with 3D models, but traditional detection datasets do not include 3D modality.

At the same time, 3D reconstruction approaches based on neural networks have been developed. These architectures can be used, for example, in industrial or virtual reality applications. Recently, promising 3D reconstruction methods, such as 3D-C2FT [11], LegoFormer [12] or VPFusion [13] hav emerged. Their performances are usually evaluated on ShapeNet [2]. Although this dataset includes a wide range of commonly encountered objects, it could be supplemented with new semantic classes present in detection datasets such as MS-COCO [1]. In addition, ShapeNet [2] only offers synthetic renderings, which limits the application of 3D reconstruction networks to real-world situations.

Thus, we propose 3D-COCO, an extension of the widely used MS-COCO [1] dataset, adapted for object detection configurable with text, 2D images, or 3D CAD (Computer-Aided Design) model queries and for single or multi-view 3D reconstruction. The 3D-COCO dataset opens new perspectives to image detection by providing 3D models that are automatically aligned with 2D annotations. It also opens the way to the integration of real images for 3D reconstruction which until now remained confined to synthetic images. Moreover, 3D-COCO provides a wider variety of semantic classes for 3D reconstruction in addition to ShapeNet’s ones [2]. The Objaverse [3] 3D models database is used with ShapeNet [2] to provide a sufficient number of objects for each of the 80 MS-COCO [1] semantic classes. Alignment between collected 3D models and MS-COCO [1] 2D annotations is made using a simple, yet effective automatic class-driven retrieval method based on IoU.

To summarize, we make the following contributions:

- We propose 3D-COCO, a dataset adapted both to 2D-to-3D configurable detection and single or multi-view 3D reconstruction. This dataset uses as a basis the original MS-COCO [1] detection dataset and extends it with 3D models collected from ShapeNet [2] and Objaverse [3].
- We present an automatic class-driven method based on

Funded by the European Union under Grant Agreement Number 101073951. Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union or European Research Executive Agency. Neither the European Union nor the granting authority can be held responsible for them.

This publication was made possible by the use of the FactoryIA supercomputer, financially supported by the Ile-De-France Regional Council.

Dataset	# Classes	# COCO	# Images	# Ann	# Models	Real	3D	Det	Rec
COCO 2017 [1]	80	80	164K	897K	N/A	✓	✗	✓	✗
3DObject [14]	10	7	6.7K	N/C	N/A	✓	✗	✓	✗
EPFL Car [15]	1	1	2.3K	N/C	N/A	✓	✗	✓	✗
NYU Depth [16]	894	N/C	1.4K	35.1K	N/A	✓	✗	✓	✗
SUN RGB-D [17]	≈800	N/C	10.3K	211.2K	N/A	✓	✗	✓	✗
KITTI [18]	2	2	15K	200K	N/A	✓	✗	✓	✗
IKEA [19]	11	5	0.8K	72K	0.2K	✓	✓	✓	✓
PASCAL3D+ [20]	12	12	30.9K	13.8K	0.1K	✓	✓	✓	✓
ObjectNet3D [21]	100	≈40	90.1K	204.6K	44.1K	✓	✓	✓	✓
ABO [22]	63	≈15	398.2K	6.3K	8K	✓	✓	✓	✓
ShapeNetCore [2]	55	27	N/A	N/A	51.3K	✗	✓	✗	✓
3D-Future [23]	34	≈5	20.2K	37.4K	10.0K	✗	✓	✓	✓
Google Scans [24]	≈17	≈4	N/A	N/A	1K	✗	✓	✓	✓
CO3D [25]	50	50	1,500K	18.6K	18.6K	✓	✗	✓	✓
Pix3D [26]	9	4	10.1K	10.1K	0.4K	✓	✓	✓	✓
PhotoShape [27]	1	1	N/A	N/A	5.8K	✗	✓	✗	✓
Objaverse [3]	N/A	80	N/A	N/A	≈800K	✗	✓	✗	✓
Objaverse XL [28]	N/A	80	N/A	N/A	≈10,000K	✗	✓	✗	✓
3D-COCO (Ours)	80	80	164K	897K	28K	✓	✓	✓	✓

Table 1: Properties of different detection and 3D reconstruction datasets. The different columns correspond to the dataset name, the number of semantic classes, the number of semantic classes shared with MS COCO [1], the number of images, the number of detection annotations, the number of 3D models, the synthetic or realistic nature of images, the availability of 3D models, the possibility to use the dataset for detection configurable with textual or image queries and the possibility to use the dataset for 3D reconstruction

IoU retrieval to match each MS-COCO [1] 2D annotation with the best 3D model in the dataset in terms of shape and geometry similarity.

2. RELATED WORK

In computer vision, combining image modalities with 3D presents significant interest due to its potential to enhance the accuracy of scene understanding and generation tasks. By integrating these complementary modalities, computer vision systems gain improved spatial awareness and object recognition capabilities, effectively addressing challenges such as occlusions, variable lighting, and perspective distortions, which are commonly encountered in analyses based on 2D images.

In the context of object detection configurable with queries, many studies have already been led to propose datasets. We can for example cite MS-COCO [1], 3DObject [14], EPFL Car [15], or NYU Depth [16]. Indeed, these datasets provide images with annotation files containing bounding boxes and labels which can be used for simple detection tasks or detection tasks with text queries and 2D image queries respectively extracted from labels and bounding boxes. Some other detection datasets also provide 3D CAD models such as ObjectNet3D [21], ABO [22], etc.

Meanwhile, other datasets for 3D reconstruction tasks have been proposed, such as ShapeNet [2], PASCAL3D+

[20], or more recently the extensive databases Objaverse [3], and ObjaverseXL [28].

Among all these datasets, 3D models can be provided in many different formats : multi-view images in KITTI [18] ; RGB-D images in SUN-RGBD [17] ; point clouds in Google Scans [24] and CO3D [25] ; meshes in IKEA [19], PASCAL3D+ [20], ObjectNet3D [21], ABO [22], 3DFuture [23], Pix3D [26] and PhotoShape [27] ; voxel grids in ShapeNet [2].

Moreover, it can be noticed that datasets represent either widespread concepts such as MS COCO [1] or ObjectNet3D [21] or very specialized categories of objects such as EPFL Car [15] or KITTI [18].

Some relevant information about all these datasets is summarized in Table 1. The motivation behind 3D-COCO is to provide a common object dataset that addresses most of the scene understanding and 3D reconstruction tasks. To reach such a goal, the MS-COCO [1] detection dataset is used as a baseline. Indeed, this dataset provides 164K realistic images with many detection annotations (about 897K) representing instances of 80 semantic classes. Moreover, this dataset is used as a reference for detection, segmentation, and pose estimation tasks.

The 3D-COCO dataset is equivalent in format to the ObjectNet3D [21] as it provides all the necessary data to train classical detection networks, with 3D models and a 2D-3D alignment to integrate 3D in detection. ShapeNet Core [2],

COCO Label	ID	Syn ID	# Models	SN	OV	COCO Label	ID	Syn ID	# Models	SN	OV
person	1	05224944	12		✓	wine glass	46	03443167	59		✓
bicycle	2	02834778	37		✓	cup	47	03152175	25		✓
car	3	02958343	3514	✓		fork	48	03388794	20		✓
motorcycle	4	03790512	337	✓		knife	49	03624134	424	✓	
airplane	5	02691156	4045	✓		spoon	50	04291140	8		✓
bus	6	02924116	939	✓		bowl	51	02880940	185	✓	
train	7	04468005	389	✓		banana	52	07769568	52		✓
truck	8	04497386	20		✓	apple	53	07755101	51		✓
boat	9	02861626	22		✓	sandwich	54	07711710	16		✓
traffic light	10	06887235	36		✓	orange	55	07763583	13		✓
fire hydrant	11	03351744	82		✓	broccoli	56	07730735	13		✓
stop sign	13	04224949	31		✓	carrot	57	07746183	12		✓
parking meter	14	03897029	11		✓	hot dog	58	07713282	19		✓
bench	15	02828884	1776	✓		pizza	59	07889783	44		✓
bird	16	01505702	24		✓	donut	60	07654678	37		✓
cat	17	02124272	27		✓	cake	61	07644479	27		✓
dog	18	02086723	28		✓	chair	62	03001627	6778	✓	
horse	19	02377103	22		✓	couch	63	04256520	3027	✓	
sheep	20	02414351	6		✓	potted plant	64	00017402	31		✓
cow	21	01890428	20		✓	bed	65	02818832	219	✓	
elephant	22	02506148	34		✓	dining table	67	03205892	32		✓
bear	23	02134305	14		✓	toilet	70	04453655	81		✓
zebra	24	02393701	7		✓	tv	72	03211117	1093	✓	
giraffe	25	02441664	33		✓	laptop	73	03642806	451	✓	
backpack	27	02772753	24		✓	mouse	74	03799022	25		✓
umbrella	28	04514450	21		✓	remote	75	04074963	66	✓	
handbag	31	02777157	25		✓	keyboard	76	03085013	64	✓	
tie	32	03821155	5		✓	cell phone	77	02992529	830	✓	
suitcase	33	02773838	83	✓		microwave	78	03761084	152	✓	
frisbee	34	03402783	7		✓	oven	79	03868196	21		✓
skis	35	04235116	8		✓	toaster	80	04449446	9		✓
snowboard	36	04258901	7		✓	sink	81	04230655	19		✓
sports ball	37	02781674	115		✓	refrigerator	82	04077839	40		✓
kite	38	03626682	7		✓	book	84	02873453	17		✓
baseball bat	39	02802334	34		✓	clock	85	03046257	651	✓	
baseball glove	40	02803372	10		✓	vase	86	03593526	596	✓	
skateboard	41	04225987	152	✓		scissors	87	04155119	24		✓
surfboard	42	04370646	15		✓	teddy bear	88	04406517	24		✓
tennis racket	43	04416941	18		✓	hair drier	89	03488399	12		✓
bottle	44	02876657	483	✓		toothbrush	90	04460427	13		✓

Table 2: Information about MS-COCO [1] labels (COCO Label), their ID in MS-COCO [1], their WordNet ID (Syn ID), the number of collected models (# Models) and their source (ShapeNet [2] (SN) or Objaverse [3] (OV))

has been used to collect CAD models for the 22 classes shared with MS COCO [1]. The collection of the remaining 58 classes has been made possible thanks to the very large diversity of semantic classes in Objaverse [3].

Finally, to address a wide variety of applications, 3D-COCO provides 3D CAD models in multiple formats. Thus, these normalized inputs will be common to every user.

3. DATASET CREATION METHOD

3.1. Collection of 3D models

First, 3D models are collected to create the extension of MS-COCO [1], which does not contain any CAD model. As represented on the left side of Figure 1, the 80 MS-COCO [1] labels are first matched with the 55 ShapeNet [2] labels. Indeed,

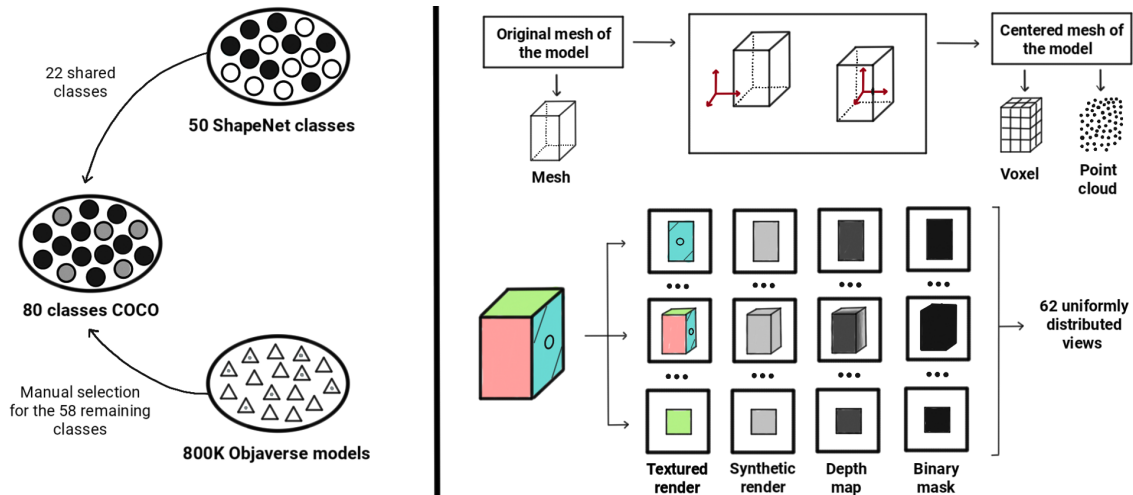


Fig. 1. 3D CAD models data collection from ShapeNet [2] and Objaverse [3] (left) and pre-processing steps including centering, conversion (upper right), and 2D rendering (lower right). Textured renderings display the model with colors and textures, synthetic renderings display the model in a uniform gray color, depth map renderings display the nearest model points in darker colors and binary mask renderings display the silhouette of the model.

ShapeNet [2] provides a wide variety of common objects’ good-quality models, which is important for our new dataset. The 22 matching categories, represented by the dark circles of Figure 1 constitute the first contribution to 3D-COCO with 26,254 models provided.

The 58 remaining labels are then completed using the Objaverse [3] 3D database. Indeed, as previously mentioned, Objaverse [3] provides about 800K CAD models from a very important number of semantic classes. Thus, a manual selection is processed on the Objaverse website to complete 3D-COCO with relevant models. First, the label name is used to make the research on Objaverse [3] website. Then, the universal identifier (UID) of manually selected models is stored and used later with the Objaverse [3] python API to collect selected models in GLB format. Finally, collected 3D models are stored using a folder architecture described later in this paper.

As a result of the manual model collection, Objaverse [3] provided 1,506 models to 3D-COCO. Figure 1 illustrated the following methodology. All the information about the MS-COCO [1] semantic classes, their identifiers, and their models are summarized in Table 2.

Once collected, 3D meshes collected on Objaverse [3] are converted from GLB into OBJ format using the trimesh python module to match the ShapeNet meshes format. Then, for each model of 3D-COCO, a centering operation is made by calculating its vertices’ mean where each vertex coordinate is weighted by the sum of the faces containing this vertex. Following this operation, models are pre-processed to make them available in the following formats: voxels of size 32, point clouds with 10,000 elements and 62-view 224×224 rendering images of different natures (textured, grey synthetic,

depth map, and binary). Point clouds and voxels are generated using respectively open3d¹ and binvox² python modules. Rendering views are generated for each of the 4 render types by using Blender’s Python API³. The 62 rendering views are uniformly sampled in an Isdyakis triacontahedron structure (composed of 62 vertices). The images and voxel sizes were chosen to fit the sizes handled by most of the 3D reconstruction networks. These operations, which make the dataset usable in a wider range of applications, are illustrated on the right part of Figure 1.

3.2. 2D-3D matching

Then, a 2D-3D alignment is implemented based on an automatic class-driven retrieval method using IoU. Indeed, CAD models contain an important amount of information about shape, and IoU is very efficient in quantifying shape similarity between elements: thus, this metric is relevant to address the 2D-3D alignment task 3D-COCO needs. In that way, each MS-COCO [1] annotation is matched with the most representative 3D CAD models of 3D-COCO in terms of geometry and shape.

The IoU-based matching method described in Figure 2 requires some pre-processing both on MS-COCO [1] annotations and on 3D-COCO models. Indeed, MS-COCO [1] annotations and API are used to generate a binary mask for each annotation as illustrated on the left part of Figure 2. This mask is saved on a 224×224 image and normalized in scale to make the represented element touch the edges of the pic-

¹<https://pypi.org/project/open3d/>

²<https://www.patrickmin.com/binvox/>

³<https://pypi.org/project/bpy/>

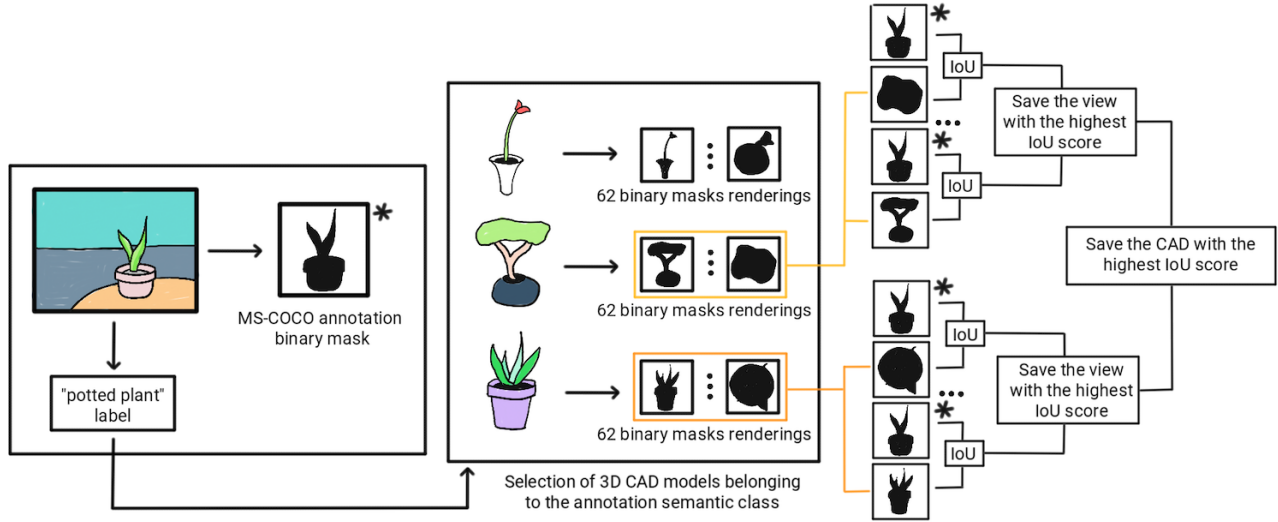


Fig. 2. Matching between MS-COCO [1] 2D annotations and 3D models using our automatic class-driven retrieval method. First, the MS-COCO [1] annotation binary mask image is extracted from the image (left). Then, the annotation label is used to select the 3D models of the same class and their 62 binary masks (middle). Finally, IoU is calculated between the MS-COCO [1] annotation mask and each CAD binary mask: models with the highest IoU score are saved as the best matching models in the 3D-COCO annotation file

ture. For each CAD model, we used the process described in 3.1 to get 62 binary masks represented on the bottom right side of Figure 1. This process allows to get silhouettes on images of similar sizes which will be compatible with an IoU calculation.

As illustrated on the right side of Figure 2, for each MS-COCO [1] annotation, IoU is computed between the binary mask of the annotation and the binary mask of all the render views of the 3D models sharing the same label. The best matching model is the one that provides the highest IoU. Thus, each MS-COCO [1] annotation is matched with its 3 most representative models of 3D-COCO.

3.3. Specific issues addressed in the annotation process

When observing MS-COCO [1] images and annotations, some situations may harm retrieval operations or during object detection:

- The annotation is too small (Figure 3a). This scenario is detected if the ratio between the number of pixels in the bounding box and the number of pixels in the image stays below a threshold (here 1%). Then, the annotation is flagged as *is_small*.
- The annotation is composed of several instances appearing in a single annotation (Figure 3b). This scenario is detected by using the existing MS-COCO [1] flag *is_crowd*.
- The annotation is truncated (Figure 3c). This scenario is detected if the ratio between the distance separating

the bounding box from the image edges and the image dimensions stays below a threshold (here 2%). Then, the annotation is flagged as *is_truncated*.

- The annotation is occluded by another annotation of the image (Figure 3d). This scenario is detected if the annotation mask intersects another mask in the image, which results in an IoU score different from 0. Then, the annotation is flagged as *is_occluded*.
- The instance is divided into multiple separated areas, (Figure 3e). This scenario is detected by connecting components labeling over the binary mask of the instance. Considering that each pixel shares the same label as the pixels it is connected to if more than one label appears after applying the method, the annotation is flagged as *is_divided*.
- There is a lack of accuracy in the labeling of MS-COCO [1] images and a lack of diversity in the collected 3D models (Figure 3f). In this example, the instance is labeled as "banana", but all the 3D models with this label represent an entire banana. This instance should then be labeled "piece of banana" or the CAD models database should be completed with meshes representing pieces of banana. This scenario is difficult to determine automatically.

Two examples of automatic class-driven IoU-based retrieval are presented in Figure 4.



Fig. 3. Example of difficult scenarios from COCO [1] images and annotations with their associated flag names

4. LICENSE AND ETHICS

From a license point of view, MS-COCO [1] and ShapeNet [2] are both CC-BY 4.0 licensed while Objaverse [3] is licensed ODC-BY. Then, 3D-COCO is licensed in a compatible and nonrestrictive way regarding the datasets that are used.

Regarding ethical considerations, 3D-COCO’s contribution to the field extends solely to the addition of 3D CAD models and the implementation of 2D-3D alignment techniques. This augmentation of MS-COCO does not alter or affect the original dataset’s adherence to privacy and ethical standards.

5. CONCLUSION

To conclude, 3D-COCO was thought of as an extension of the original MS-COCO [1] dataset including 27,760 3D CAD models of 80 different semantic classes collected from ShapeNet [2] and Objaverse [3]. An automatic class-driven retrieval method based on IoU has been implemented to provide a 2D-3D alignment between the 860,001 training or the 36,781 validation annotations and the collected 3D models. This extension bridges the gap between MS-COCO [1] and the 3D world: new tasks such as detection networks configurable with 3D modes, synthetic multi-view 3D reconstruction networks, or real single-view 3D reconstruction networks could be addressed thanks to 3D-COCO.

The philosophy of 3D-COCO lies in its transparency, open access, and the possibility for users to iterate over the originally proposed dataset through code sharing.

Nevertheless, for future iteration of the dataset, 3D-COCO could be improved with a better 2D-3D alignment method for articulated 3D models such as humans or animals. Exploring other retrieval methods based on neural network feature extraction or integrating new 3D models to have a



1st match IoU = 0.49 **2nd match** IoU = 0.47 **3rd match** IoU = 0.46



1st match IoU = 0.66 **2nd match** IoU = 0.63 **3rd match** IoU = 0.62

Fig. 4. MS-COCO [1] images with *Truck* and *Horse* annotations followed by their 3 best matching models predicted by our IoU-based retrieval method

more balanced number of CAD models for each class are also relevant perspectives.

6. REFERENCES

- [1] T. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. Zitnick, “Microsoft coco: Common objects in context,” in *ECCV*, 2014, pp. 740–755.
- [2] A. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song, H. Su, et al., “Shapenet: An information-rich 3d model repository,” *arXiv:1512.03012*, 2015.
- [3] M. Deitke, D. Schwenk, J. Salvador, L. Weihs, O. Michel, E. VanderBilt, L. Schmidt, K. Ehsani, A. Kembhavi, and A. Farhadi, “Objaverse: A universe of annotated 3d objects,” in *CVPR. IEEE/CVF*, 2023, pp. 13142–13153.

- [4] S. Ren, K. He, R. Girshick, and J. Sun, “Faster r-cnn: Towards real-time object detection with region proposal networks,” *Advances in neural information processing systems*, vol. 28, 2015.
- [5] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, “You only look once: Unified, real-time object detection,” in *CVPR*. IEEE, 2016, pp. 779–788.
- [6] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C. Fu, and A. Berg, “Ssd: Single shot multibox detector,” in *ECCV*, 2016, pp. 21–37.
- [7] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, “End-to-end object detection with transformers,” in *ECCV*, 2020, pp. 213–229.
- [8] M. Everingham, S. Eslami, L. Van Gool, C. Williams, J. Winn, and A. Zisserman, “The pascal visual object classes challenge: A retrospective,” *IJCV*, vol. 111, pp. 98–136, 2015.
- [9] A. Kuznetsova, H. Rom, N. Alldrin, J. Uijlings, I. Krasin, J. Pont-Tuset, S. Kamali, S. Popov, M. Mallocci, A. Kolesnikov, et al., “The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale,” *IJCV*, vol. 128, no. 7, pp. 1956–1981, 2020.
- [10] Y. Zang, W. Li, K. Zhou, C. Huang, and C. Loy, “Open-vocabulary detr with conditional matching,” in *ECCV*, 2022, pp. 106–122.
- [11] L. Tiong, D. Sigmund, and A. Teoh, “3d-c2ft: Coarse-to-fine transformer for multi-view 3d reconstruction,” in *ACCV*, 2022, pp. 1438–1454.
- [12] F. Yagubbayli, Y. Wang, A. Tonioni, and F. Tombari, “Legoformer: Transformers for block-by-block multi-view 3d reconstruction,” *arXiv:2106.12102*, 2021.
- [13] J. Mahmud and J. Frahm, “Vpfusion: Joint 3d volume and pixel-aligned feature fusion for single and multi-view 3d reconstruction,” *arXiv:2203.07553*, 2022.
- [14] S. Savarese and L. Fei-Fei, “3d generic object categorization, localization and pose estimation,” in *2007 IEEE 11th International Conference on Computer Vision*. IEEE/CVF, 2007, pp. 1–8.
- [15] M. Ozuysal, V. Lepetit, and P. Fua, “Pose estimation for category specific multiview object localization,” IEEE, 06 2009.
- [16] N. Silberman, D. Hoiem, P. Kohli, and R. Fergus, “Indoor segmentation and support inference from rgb-d images,” 10 2012, pp. 746–760.
- [17] S. Song, S. Lichtenberg, and J. Xiao, “Sun rgb-d: A rgb-d scene understanding benchmark suite,” in *CVPR*. IEEE, June 2015.
- [18] A. Geiger, P. Lenz, and R. Urtasun, “Are we ready for autonomous driving? the kitti vision benchmark suite,” in *CVPR*. IEEE, 2012, pp. 3354–3361.
- [19] J. Lim, H. Pirsiavash, and A. Torralba, “Parsing ikea objects: Fine pose estimation,” in *ICCV*. IEEE, December 2013.
- [20] Y. Xiang, R. Mottaghi, and S. Savarese, “Beyond pascal: A benchmark for 3d object detection in the wild,” in *WACV*. IEEE, 2014, pp. 75–82.
- [21] Y. Xiang, W. Kim, W. Chen, J. Ji, C. Choy, H. Su, R. Mottaghi, L. Guibas, and S. Savarese, “Objectnet3d: A large scale database for 3d object recognition,” in *ECCV*, 2016.
- [22] J. Collins, S. Goel, K. Deng, A. Luthra, L. Xu, E. Gundogdu, X. Zhang, T. Vicente, T. Dideriksen, H. Arora, et al., “Abo: Dataset and benchmarks for real-world 3d object understanding,” in *CVPR*. IEEE/CVF, 2022, pp. 21126–21136.
- [23] H. Fu, R. Jia, L. Gao, M. Gong, B. Zhao, S. Maybank, and D. Tao, “3d-future: 3d furniture shape with texture,” *IJCV*, vol. 129, 12 2021.
- [24] L. Downs, A. Francis, N. Koenig, B. Kinman, R. Hickman, K. Reymann, T. McHugh, and V. Vanhoucke, “Google scanned objects: A high-quality dataset of 3d scanned household items,” in *ICRA*. IEEE, 2022, pp. 2553–2560.
- [25] J. Reizenstein, R. Shapovalov, P. Henzler, L. Sbordone, P. Labatut, and D. Novotny, “Common objects in 3d: Large-scale learning and evaluation of real-life 3d category reconstruction,” in *ICCV*. IEEE/CVF, 2021, pp. 10901–10911.
- [26] X. Sun, J. Wu, X. Zhang, Z. Zhang, C. Zhang, T. Xue, J. Tenenbaum, and W. Freeman, “Pix3d: Dataset and methods for single-image 3d shape modeling,” in *CVPR*. IEEE, 2018, pp. 2974–2983.
- [27] K. Park, K. Rematas, A. Farhadi, and S. Seitz, “Photoshape: Photorealistic materials for large-scale shape collections,” *arXiv:1809.09761*, 2018.
- [28] M. Deitke, R. Liu, M. Wallingford, H. Ngo, O. Michel, A. Kusupati, A. Fan, C. Laforte, V. Voleti, S. Gadre, et al., “Objaverse-xl: A universe of 10m+ 3d objects,” *arXiv:2307.05663*, 2023.