

# OCH3R: Object-Centric Holistic 3D Reconstruction

Anonymous CVPR submission

Paper ID 14520

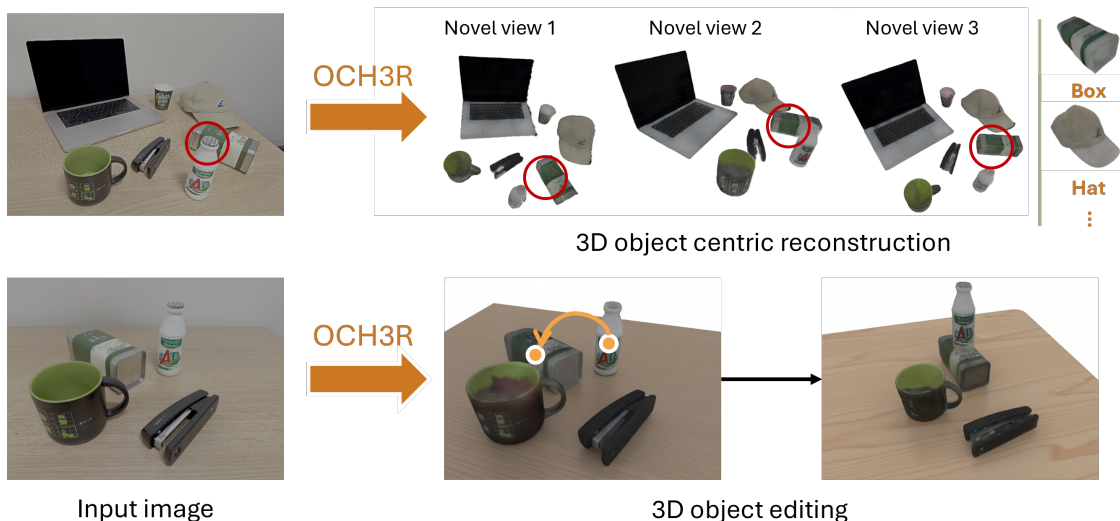


Figure 1. **OCH3R enables fully object-centric 3D scene reconstruction from a single RGB image.** Given one input view, OCH3R discovers all object instances, predicts their 6D poses, and reconstructs each object as a manipulable 3D Gaussian model in a single forward pass. Our feed-forward, per-pixel prediction framework supports selecting, moving, and rendering objects from arbitrary novel views without external segmentors or multi-stage pipelines. OCH3R produces amodally complete, editable 3D objects and generalizes to cluttered real scenes, enabling downstream tasks such as rearrangement and AR editing. The red circles highlight occluded regions that OCH3R successfully completes in 3D.

## Abstract

001 *Object-centric scene understanding is a fundamental chal-*  
 002 *lenge in computer vision. Existing approaches often rely*  
 003 *on multi-stage pipelines that first apply pre-trained segmen-*  
 004 *tors to extract individual objects, followed by per-object*  
 005 *3D reconstruction. Such methods are computationally ex-*  
 006 *pensive, fragile to segmentation errors, and scale poorly*  
 007 *with scene complexity. We introduce **OCH3R**, a unified*  
 008 *framework for **Object-Centric Holistic 3D Reconstruction***  
 009 *from a single RGB image. OCH3R performs one forward*  
 010 *pass to simultaneously predict all object instances with*  
 011 *their 6D poses and detailed 3D reconstructions. The key*  
 012 *idea is a transformer architecture that predicts per-pixel at-*  
 013 *tributes, including CLIP-based category embeddings, met-*  
 014 *ric depth, normalized object coordinates (NOCS), and a*  
 015 *fixed number of 3D Gaussians representing each object.*  
 016 *To supervise these Gaussian reconstructions, we transform*

*them into canonical space using the predicted 6D poses* 017  
*and align them with pre-rendered canonical ground truth,* 018  
*avoiding costly per-image Gaussian label generation. On* 019  
*standard indoor benchmarks, OCH3R achieves state-of-* 020  
*the-art performance across monocular depth estimation,* 021  
*open-vocabulary semantic segmentation, and RGB-only* 022  
*category-level 6D pose estimation, while producing high-* 023  
*fidelity, editable per-object reconstructions. Crucially, in-* 024  
*ference is fully feed-forward and scales independently of the* 025  
*number of objects, offering orders-of-magnitude speedups* 026  
*over conventional multi-stage pipelines in cluttered scenes.* 027

## 1. Introduction

Understanding a scene as a composition of discrete, posed 029  
 objects from a single image is a long-standing goal in 030  
 computer vision. Many downstream applications including 031  
 robotic manipulation, AR editing, and simulation rely on 032

033 object-centric outputs [73], where each object is represented  
034 with geometry, pose, and semantics that can be selected or  
035 manipulated. We study the following problem: given a single  
036 RGB image of an indoor tabletop scene, recover all objects  
037 together with their 6D poses and corresponding 3D  
038 Gaussians in one forward pass.

039 Prior work largely falls into two groups. Scene-level,  
040 feed-forward methods (*e.g.*, one-pass Gaussian predictors)  
041 [5, 54, 68, 74, 82] are fast and photorealistic, but  
042 typically produce an undifferentiated “soup” of geometry  
043 without instance-level structure, canonical frames, or object  
044 poses. Object-level approaches [24, 25, 73] instead  
045 rely on multi-stage pipelines that begin with external open-  
046 vocabulary segmentors and then perform per-object recon-  
047 struction, alignment, and correction. These systems often  
048 depend on RGB-D inputs or category-specific priors, and  
049 they are fragile to upstream errors, difficult to scale with the  
050 number of objects, and not trained end-to-end. As a result,  
051 their robustness and accuracy degrade in cluttered tabletop  
052 scenes.

053 To address these limitations, we introduce **OCH3R**, a  
054 unified, object-centric, holistic 3D reconstructor that con-  
055 verts a single RGB image into a set of posed 3D objects in  
056 one pass. The key in our design is a 48-layer transformer  
057 that predicts dense, pixel-aligned attributes: CLIP [45]-  
058 based category embeddings, metric depth, normalized ob-  
059 ject coordinates (NOCS) [62], and a small set of 3D Gaus-  
060 sians [29] per pixel. During inference, object instances and  
061 their poses are recovered by clustering the semantic embed-  
062 dings and estimating each object’s SIM(3) pose using the  
063 predicted NOCS field.

064 To train the Gaussian representation, we allow Gaussians  
065 at each pixel to move freely off the pixel rays to compensate  
066 for (self-) occlusion. Rather than supervising Gaussians per  
067 training image [52, 53], we adopt Canonical-Space Super-  
068 vision: per-object Gaussians are transformed into canonical  
069 space using its SIM(3) pose, and their renderings are opti-  
070 mized against pre-rendered ground truth in the canonical  
071 frame. This eliminates the need for costly per-image Gaus-  
072 sian labels and promotes amodal shape completion.

073 We train on a curated, large-scale dataset that integrates  
074 PACE [75], Omni6DPose [81], GSO [12], and Hyper-  
075 sim [48], offering broad coverage across object categories,  
076 poses, and occlusion patterns.

077 Our experiments show that OCH3R substantially out-  
078 performs previous baselines across all evaluated tabletop  
079 object-centric benchmarks. OCH3R delivers consistently  
080 higher geometric accuracy, better semantic alignment, and  
081 significantly more complete amodal reconstructions. Im-  
082 portantly, because OCH3R reconstructs all objects in a single  
083 forward pass, it achieves orders-of-magnitude faster infer-  
084 ence than multi-stage pipelines while avoiding their brit-  
085 tleness to segmentation or pose-estimation errors. Together,

086 these results highlight the effectiveness of our unified for-  
087 mulation and its practical advantages for real-world object-  
088 centric applications.

089 To summarize, our contributions are as follows:

- 090 1. We construct a large scale dataset for holistic object cen-  
091 tric 3D scene representation. We assemble, relabel, and  
092 align PACE [75], Omni6DPose [81], GSO [12], and Hy-  
093 persim [48] into a unified dataset designed for object  
094 centric 3D tasks, providing per instance masks, segmen-  
095 tation labels, SIM(3) poses, and 3D models.
- 096 2. We propose a model that yields high-fidelity 3D recon-  
097 structions, recovering fine-grained geometry and amodal  
098 structure, while jointly predicting semantics, monocular  
099 depth, and object poses in a single pass.
- 100 3. Experiments show that our model reconstructs real-  
101 world tabletop scenes with arbitrary numbers of ob-  
102 jects, delivering more photorealistic and amodally com-  
103 plete results while running far faster than multi-stage  
104 pipelines.

## 105 2. Related Work

106 **Feed-forward 3D reconstruction.** Feed-forward 3D re-  
107 construction maps one or a few images directly to a ren-  
108 derable 3D scene. Previous works have explored 3D rep-  
109 resentations including voxel grids [57, 58], multi-plane im-  
110 ages [35, 59], meshes [17, 18], surfel [16], and radiance  
111 fields [21, 76]. More recently, 3D Gaussians [30] have  
112 emerged as a dominant representation for feed-forward re-  
113 gression thanks to their real-time differentiable rendering  
114 and compatibility with high-capacity 2D backbones.

115 Early feed-forward Gaussian predictors focus on single,  
116 centered objects, assigning one Gaussian to each input pixel  
117 and directly regressing its parameters without test-time op-  
118 timization [52, 55, 71, 80]. These models deliver fast and  
119 high-quality reconstructions, but they assume clean, un-  
120 cluttered inputs and cannot handle occlusions, multiple in-  
121 stances, or reassemble per-object predictions into a scene.

122 In contrast, scene-level Gaussian models predict a dense  
123 Gaussian field for an entire scene from one [53] or mul-  
124 tiple images [5, 8, 65, 68, 74, 82]. Techniques includ-  
125 ing probabilistic splatting [5], cost-volume aggregation [8],  
126 depth conditioning [68], pose-free formulations [74] and  
127 large transformer backbones [82] have been explored to im-  
128 prove performance. While effective for novel-view synthe-  
129 sis, these scene-level approaches treat the world as a single  
130 undifferentiated cloud without instance decomposition, pre-  
131 venting downstream reasoning or interaction.

132 **Object-centric scene reconstruction.** A separate line of  
133 work performs object-centric scene reconstruction, explic-  
134 itly recovering a set of 3D object instances and their lay-  
135 out. IM2CAD [26], Total3DUnderstanding [41], Zhang *et*  
136 *al.* [79], and CoReNet [44] reconstruct indoor scenes from

137 a single image by detecting furniture and room layout, then  
 138 retrieving or predicting per-object geometry and enforcing  
 139 consistency in a shared 3D frame. CAD- and RGB-D-based  
 140 pipelines such as Mask2CAD [32], ROCA [19], Center-  
 141 Snap [24], and ShAPO [25] further combine instance detec-  
 142 tion and depth with CAD retrieval or learned latent shape  
 143 codes for each object, making them sensitive to upstream  
 144 errors and computationally costly as the number of objects  
 145 increases.

146 More recent methods introduce strong generative pri-  
 147 ors but largely retain this compositional, multi-stage de-  
 148 sign: Gen3DSR [1], CAST [73], and DepR [83] first ap-  
 149 ply monocular depth estimation and instance segmenta-  
 150 tion, then run object-level image-to-3D or diffusion mod-  
 151 els and compose the resulting objects into a coherent scene;  
 152 MIDI [23] extends pre-trained image-to-3D generators to  
 153 a multi-instance diffusion model that still takes segmented  
 154 object crops as input. Consequently, computational cost and  
 155 brittleness scale with the number and quality of segmented  
 156 instances. With a sufficiently large dataset and a sufficiently  
 157 powerful model, we show that single-view, object-aware 3D  
 158 reconstruction can be approached as a direct, feed-forward  
 159 prediction problem, rather than a fragile sequence of seg-  
 160 mentation, retrieval, optimization, or generative refinement.  
 161 In practice, this shift yields reconstructions that are not only  
 162 orders-of-magnitude faster but also higher-fidelity.

### 163 3. Preliminaries

164 **3D Gaussian Splatting.** Gaussian Splatting [30] renders a  
 165 scene represented by a finite set of anisotropic 3D Gaussian  
 166 primitives by projecting each primitive to the image plane  
 167 as a 2D Gaussian and alpha-compositing them in visibility  
 168 order, yielding a fast, differentiable approximation to emis-  
 169 sion-absorption volume rendering. Compared with ray-  
 170 sampled neural fields [40], splatting enables real-time ren-  
 171 dering and efficient gradient backpropagation, and is widely  
 172 used as the rendering backbone in recent feed-forward re-  
 173 construction methods [5, 6, 8, 52, 54, 55, 65, 71, 74, 82];  
 174 we adopt the same renderer throughout.

175 **Normalized Object Coordinate Space.** Normalized Ob-  
 176 ject Coordinate Space (NOCS) [63] assigns each 3D point  
 177 on an object instance a category-level, pose-invariant coor-  
 178 dinate  $\mathbf{c} \in [0, 1]^3$  within a unit canonical cube whose axes  
 179 are consistently aligned across instances of that category.  
 180 We denote this unit cube as the *canonical space* and to its  
 181 associated rigid coordinate system as the *canonical frame*.

182 Dense per-pixel NOCS predictions  $\hat{\mathbf{c}}_{u,v}$  provide pixel-  
 183 to-canonical correspondences that, together with the pre-  
 184 dicted 3D point map, are sufficient to recover an instance’s  
 185 category-level pose  $\Pi = (s, R, \mathbf{t}) \in \text{SIM}(3)$  (Sec. 4.2). By  
 186 definition,  $\Pi$  transforms canonical coordinates to the cam-  
 187 era (or scene) frame:  $\mathbf{x}^{\text{cam}} = \Pi(\mathbf{x}^{\text{can}}) = sR\mathbf{x}^{\text{can}} + \mathbf{t}$ , with

inverse mapping given by  $\Pi^{-1}(\mathbf{x}^{\text{cam}}) = s^{-1}R^\top \cdot (\mathbf{x}^{\text{cam}} - \mathbf{t})$ . 188

## 189 4. Method

### 190 4.1. Problem formulation and notation

191 Given a single RGB image  $I \in \mathbb{R}^{H \times W \times 3}$  with un-  
 192 known intrinsics  $K$ , OCH3R converts the image into an  
 193 object-centric scene: a set of instances, each with a  
 194 category-level  $\text{SIM}(3)$  pose and a high-fidelity 3D repre-  
 195 sentation. Achieving this in one pass requires pixel-aligned  
 196 predictions that are sufficient to (i) discover instances and  
 197 semantics, (ii) recover a metric similarity transform for each  
 198 object, and (iii) assemble amodally complete Gaussian rep-  
 199 resentation of each object into an interactive scene.

200 Specifically, for each pixel  $(u, v)$ , our network  $\Phi$  out-  
 201 puts:

$$202 \Phi(I)_{u,v} = (\hat{\mathbf{e}}_{u,v}, \hat{d}_{u,v}, \hat{\mathbf{c}}_{u,v}, \hat{\mathcal{G}}_{u,v}), \quad (1)$$

203 where  $\hat{\mathbf{e}}_{u,v} \in \mathbb{R}^{512}$  is the semantic label of the object that  
 204 this pixel belongs to. We define semantic label of an ob-  
 205 ject as the CLIP [45] embedding of the object’s category  
 206 name [34].  $\hat{d}_{u,v} \in \mathbb{R}^+$  is the predicted metric depth. It  
 207 enables back-projecting the pixel into 3D space via  $\hat{\mathbf{p}}_{u,v} =$   
 $\hat{d}_{u,v} \cdot K^{-1} [u \ v \ 1]^\top \in \mathbb{R}^3$ .  $\hat{\mathbf{c}}_{u,v} \in [0, 1]^3$  is the pre-  
 208 dicted NOCS [63] coordinate, which enables  $\text{SIM}(3)$  pose  
 209 recovery.

210  $\hat{\mathcal{G}}_{u,v} = \{g_{u,v}^{(i)}\}_{i=1}^k$  is a small set of anisotropic 3D Gaus-  
 211 sian primitives (we use  $k = 2$ ) that will be aggregated into  
 212 per-object reconstruction: 213

$$214 g_{u,v}^{(i)} = (\boldsymbol{\mu}_{u,v}^{(i)}, \Sigma_{u,v}^{(i)}, \alpha_{u,v}^{(i)}, \mathbf{S}_{u,v}^{(i)}), \quad (2)$$

215 with mean  $\boldsymbol{\mu}_{u,v}^{(i)} \in \mathbb{R}^3$ , covariance  $\Sigma_{u,v}^{(i)} \in \mathbb{S}_{++}^3$ , opacity  
 216  $\alpha_{u,v}^{(i)} \in (0, 1)$ , and RGB spherical harmonics (SH) coeffi-  
 217 cients  $\mathbf{S}_{u,v}^{(i)} \in \mathbb{R}^{3(L+1)^2}$  (order  $L$ ). The Gaussian param-  
 218 eters are defined and predicted in camera frame, and will be  
 219 transformed into each object’s canonical frame for supervi-  
 220 sion and inference (Sec. 4.3).

221 Following VGGT [64], we also predict the camera field  
 222 of view  $(\hat{\theta}_w, \hat{\theta}_h)$  of the input image and construct  $K$  with  
 223  $f_w = \frac{W}{2 \tan(\hat{\theta}_w/2)}$ ,  $f_h = \frac{H}{2 \tan(\hat{\theta}_h/2)}$  and principal point at  
 224 image center.

225 In Sec. 4.2, we show how  $\hat{\mathbf{e}}$ ,  $\hat{d}$ ,  $\hat{\mathbf{c}}$ , and  $\hat{\mathcal{G}}$  are used to dis-  
 226 cover instances, estimate object poses, and assemble recon-  
 227 structed objects. Sec. 4.3 introduces Canonical Space Su-  
 228 pervision (CSS) that trains Gaussians to be object-aligned  
 229 and amodally complete without per-image Gaussian labels.  
 230 Sec. 4.4 summarizes architectural and training details. Our  
 231 full pipeline is given in Fig. 2.

### 232 4.2. Assembling objects from dense predictions

233 **Instance discovery.** At inference time, for each pixel,  
 234 we first compute the cosine similarity between the pre-  
 235 dicted embedding  $\hat{\mathbf{e}}_{u,v}$  and a set of predefined category

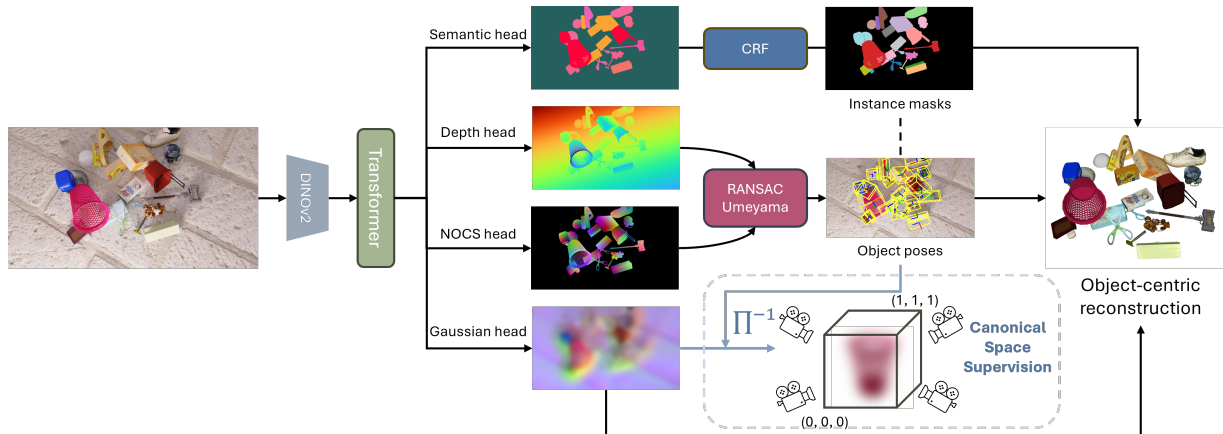


Figure 2. **Overview of our single-view object-centric 3D reconstruction pipeline.** Given a single RGB input, we extract dense DINOv2 features and feed them to a transformer that predicts per-pixel depth, CLIP-space semantic embeddings, NOCS coordinates, and Gaussian primitives. A CRF refines semantic affinities to produce coherent instance masks. For each instance, we estimate a category-level SIM(3) pose via RANSAC-Umeyama using the predicted NOCS-to-3D correspondences, enabling a transformation from camera space into the canonical object frame. The per-pixel Gaussians are then grouped and transformed into canonical space, where Canonical-Space Supervision (CSS) trains them to form amodally complete, compact 3D Gaussians. Aggregating all reconstructed objects yields an interactive, object-aligned scene representation from a single image.

236 name CLIP embeddings  $\{l_c\}$ . We then apply a fully connected  
 237 conditional random field (CRF) [31], using unary  
 238 potentials defined as  $-\log \frac{\exp(\cos(\hat{e}_{u,v}, l_c)/\tau)}{\sum_{c'} \exp(\cos(\hat{e}_{u,v}, l_{c'})/\tau)}$  for each  
 239 category  $c$ , where  $\tau$  denotes the temperature parameter in  
 240 the softmax function. Pairwise potentials are defined as  
 241  $\cos(\hat{e}_{u,v}, \hat{e}_{u',v'})$ . For more details about CRF, we refer the  
 242 reader to [31]. This process yields groups  $\{\hat{\mathcal{P}}_j\}$ , where each  
 243  $\hat{\mathcal{P}}_j$  represents the set of pixels corresponding to object  $j$ .

244 **Pose estimation.** With the pixels for each object instance  
 245 identified, we use their predicted NOCS [62] coordinates  
 246  $\mathbf{c}_{u,v}$  to determine the object’s precise SIM(3) pose in the  
 247 scene. The NOCS coordinates establish a correspondence  
 248 between a point’s observed position in the scene and its  
 249 standardized position within a unit canonical cube. We  
 250 can therefore solve for the similarity transformation  $\hat{\Pi}_j =$   
 251  $(\hat{s}_j, \hat{R}_j, \hat{\mathbf{t}}_j)$ , representing the scale, rotation, and transla-  
 252 tion, which maps the canonical space of object  $j$  to the cam-  
 253 era space. This is achieved by minimizing the alignment  
 254 error between the back-projected 3D points and the trans-  
 255 formed NOCS coordinates over all pixels belonging to that  
 256 instance:

$$257 \quad \hat{\Pi}_j = \arg \min_{\Pi} \sum_{(u,v) \in \hat{\mathcal{P}}_j} \|\hat{\mathbf{p}}_{u,v} - \Pi(\hat{\mathbf{c}}_{u,v})\|^2, \quad (3)$$

258 where  $\Pi(\hat{\mathbf{c}}_{u,v}) = sR \cdot \hat{\mathbf{c}}_{u,v} + \mathbf{t}$ . This optimization can be  
 259 solved using Umeyama algorithm [61] with RANSAC [15].  
 260 The resulting inverse transformation,  $\hat{\Pi}_j^{-1}$ , gives us a di-  
 261 rect mapping from the cluttered scene into the clean canoni-  
 262 cal space for each object. Notably, this NOCS prediction  
 263 can also be used to differentiate object instances with the  
 264 same category name but that are adjacent in the mask, where

265 CRF alone may not be enough. We run multiple RANSACs  
 266 within each CRF-generated mask, and output objects when  
 267 there are still enough inliers.

268 **Object Gaussians.** With the instance mask and esti-  
 269 mated pose in hand, we obtain each object’s canonical-  
 270 space Gaussian representation by transforming every pre-  
 271 dicted Gaussian mean as  $\mu_{u,v}^{\text{can},(i)} = \hat{\Pi}_j^{-1}(\mu_{u,v}^{(i)})$  for all  
 272  $(u,v) \in \hat{\mathcal{P}}_j, i \in \{1, \dots, k\}$ . The resulting set of trans-  
 273 formed Gaussians forms the complete canonical represen-  
 274 tation of object  $j$ .

275 **Efficiency.** Since OCH3R predicts all per-pixel quantities  
 276 in one forward pass, every object is reconstructed at once.  
 277 Our CUDA CRF runs in roughly 200 ms per image, and  
 278 our CUDA RANSAC Umeyama adds under 10 ms per ob-  
 279 ject, making its cost negligible. Consequently, runtime is  
 280 nearly invariant to scene complexity and remains far be-  
 281 low prior pipelines [56, 73, 83], which synthesize each ob-  
 282 ject through iterative diffusion denoising and often require  
 283 relation-graph optimization that grows quadratically with  
 284 the number of objects.

### 4.3. Canonical-Space Supervision 285

286 One key challenge for our Gaussian prediction network is  
 287 that it must infer a full, amodal set of object Gaussians from  
 288 only the pixels that are actually visible. A natural idea is  
 289 to use pre-optimized object Gaussians and place them in  
 290 the camera frame so they can serve as ground-truth super-  
 291 vision. However, there lacks one-to-one correspondence  
 292 between visible pixels and ground-truth object Gaussians.  
 293 To address this, we introduce *Canonical Space Supervision*  
 294 (Fig. 3), a strategy that transfers training signals to the ob-

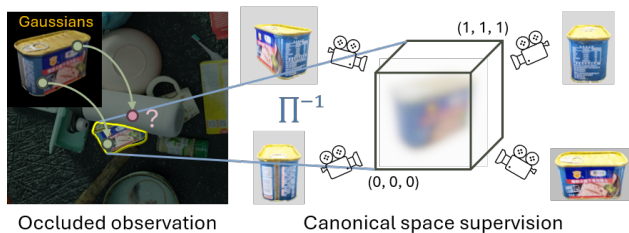


Figure 3. **Canonical Space Supervision (CSS)**. Predicted per-pixel Gaussians are transformed into the object’s canonical frame via the ground-truth pose  $\Pi^{-1}$ . In canonical space, they are supervised against pre-rendered multi-view ground-truth images, providing clean amodal signals that resolve occlusions and enforce compact, object-aligned Gaussian reconstructions.

ject’s canonical frame, where clean targets are available.

Concretely, we place each training object mesh in the canonical frame, and pre-render a set of  $N$  views  $\mathcal{V} = \{I_n^{\text{gt}}\}_{n=1}^N$ , where  $I_n^{\text{gt}} \in \mathbb{R}^{H_{\text{can}} \times W_{\text{can}} \times 3}$ . We set  $N = 42$ ,  $H_{\text{can}} = W_{\text{can}} = 512$ . This is done once per object and reused across all images containing that object.

During training, we use ground truth object masks  $\{\mathcal{P}_j\}$  to extract pixels of each object  $j$ . We transform the predicted Gaussians per object from the camera space into the canonical space with the ground truth object pose.

With the transformed Gaussians  $\{g_{u,v}^{\text{can},(i)}\}$  of the object, we can render images of that object in its canonical space using the same camera angles as  $\mathcal{V}$  with a differentiable Gaussian rasterizer [29]. Let  $\{\hat{I}_n\}_{n=1}^N$  denote the rendered images. Then the CSS loss is calculated by

$$\mathcal{L}_{\text{CSS}} = \sum_{n=1}^N \left( \|I_n^{\text{gt}} - \hat{I}_n\|_1 + \lambda_{\text{SSIM}}(1 - \text{SSIM}(I_n^{\text{gt}}, \hat{I}_n)) \right). \quad (4)$$

**Occlusion handling via off-ray offsets.** Each Gaussian mean is anchored at  $\hat{\mathbf{p}}_{u,v}$  with a predicted camera-space offset  $\Delta_{u,v}^{(i)}$ , allowing a visible pixel to spawn Gaussians behind the first surface. Because CSS supervises in canonical space, occluded Gaussians receive gradients even when not visible in the input. Optionally, we regularize with an annealed small-offset prior:  $\mathcal{L}_{\text{reg}} = \sum_{(u,v) \in \mathcal{P}_j} \sum_{i=1}^k \text{ReLU}(\|\Delta_{u,v}^{(i)} - \tau_{\text{offset}}\|_1)$ .

#### 4.4. Architecture and training details

**Architecture.** Inspired by VGGT [64], OCH3R uses a DINOv2 backbone [42] followed by a 48-layer ViT encoder and DPT-style decoder heads [46]. The input image is patchified by DINOv2 and processed by global self-attention through the encoder.

For dense tasks, each head takes features from four intermediate encoder layers (lateral skips), projects them to a common width, and upsamples to the image resolution with convolutional fusion.

For camera FOV prediction, we append a learnable camera token that is updated by a small stack of transformer blocks with adaptive layer-norm modulation and iterative refinement; a linear layer regresses the field-of-view angles  $(\hat{\theta}_w, \hat{\theta}_h)$ .<sup>1</sup>

For Gaussian prediction we decouple geometry and appearance. A geometry head outputs per-pixel off-ray offsets  $\Delta_{u,v}^{(i)}$  (for  $i=1, \dots, k$ ), which are added to the back-projected point  $\hat{\mathbf{p}}_{u,v}$  to obtain camera-frame means  $\mu_{u,v}^{(i)}$ . An appearance/shape head predicts canonical-frame scales  $\sigma_{u,v}^{(i)}$ , unit quaternions  $\mathbf{q}_{u,v}^{\text{can},(i)}$ , opacities  $\alpha_{u,v}^{(i)}$ , and SH coefficients  $\mathbf{S}_{u,v}^{(i)}$ . Following NoPoSplat [74], we provide an *RGB shortcut* to the appearance/shape head to improve fine texture details in 3D reconstruction. All parameters of each Gaussian and the  $k$  Gaussians of every pixel are simply concatenated.

**Training.** We optimize all tasks jointly with AdamW and a cosine learning-rate schedule. DINOv2 is initialized from public weights and fine-tuned end-to-end; the full list of hyperparameters is provided in the appendix.

**Depth.** We supervise canonical inverse depth as in Depth Pro [3]. Let  $f_w$  be the horizontal focal length (pixels) and  $W$  the image width, define  $C = \frac{f_w}{W \cdot d}$ . Our model outputs  $\hat{C}$  and is trained to minimize  $\mathcal{L}_{\text{depth}} = \|\hat{C} - C\|_2 + \lambda \nabla \left( \|\nabla_x(\hat{C} - C)\|_2 + \|\nabla_y(\hat{C} - C)\|_2 \right)$ . At test time, we recover metric depth by  $\hat{d}_{u,v} := \frac{f_{\text{px}}}{W \cdot \hat{C}_{u,v}}$ .

**Semantics.** During training time, we dynamically compute cosine similarities between the embedding of pixel and all the words that appear in the training image. We then encourage the embedding of the pixel to align with the ground-truth class by using the computed cosine similarities as the logits for softmax, with cross entropy loss.

**NOCS.** We reformulate NOCS coordinate regression as a bin classification task augmented with a learnable offset, which implicitly resolves ambiguities in symmetric objects. Each axis in NOCS (i.e.,  $xyz$ ) is discretized into  $M=64$  centered bins. We supervise the bin classification using a cross-entropy loss and the offset prediction using a mean squared error loss. The total loss is averaged over all foreground object pixels.

**Gaussians (CSS).** Gaussian supervision follows Canonical Space Supervision discussed in Sec. 4.3.

**Camera FOV.** We supervise  $(\hat{\theta}_w, \hat{\theta}_h)$  with a robust Huber loss on angles:  $\mathcal{L}_{\text{cam}} = \|\hat{\theta}_w - \theta_w\|_{\epsilon} + \|\hat{\theta}_h - \theta_h\|_{\epsilon}$ .

Task losses are combined with homoscedastic uncertainty weighting [28]; ablations are in the appendix.

<sup>1</sup>Our implementation retains VGGT’s iterative refinement; translation/rotation channels are present in the token state but only FOV is used at test time.

## 375 5. Dataset

376 Existing indoor benchmarks for monocular depth estima-  
377 tion and open vocabulary semantic segmentation [10, 50,  
378 51] primarily emphasize room layout and large furniture,  
379 while providing limited coverage of the object interaction  
380 scale that is central to everyday visual tasks. In contrast,  
381 progress in embodied perception [4, 14, 49, 66] and in mo-  
382 bile AR or MR systems [13, 20] requires accurate model-  
383 ing of small, manipulable, and semantically diverse objects  
384 such as cups, tools, and containers that humans routinely  
385 interact with.

386 To support this direction, we construct a new evaluation  
387 benchmark by integrating several real world datasets tai-  
388 lored to this domain. Specifically, we include the valida-  
389 tion split of HOPE [60] and the test splits of YCB Video  
390 [67], PACE [75], Omni6DPose [81] (OMNI), and NOCS  
391 [63]. For training, we curate and align four large scale  
392 sources: the training splits of PACE and Omni6DPose,  
393 Google Scanned Objects [12] renderings from Founda-  
394 tion-Pose [66], and HyperSim [48].

## 395 6. Experiments

396 We first evaluate holistic 3D object-centric reconstruction  
397 from a single RGB image in Sec. 6.1. Sec. 6.2 then demon-  
398 strates that beyond 3D reconstruction, our method also de-  
399 livers state-of-the-art zero-shot performance on depth esti-  
400 mation, segmentation, and object pose prediction. Sec. 6.3  
401 examines key design choices through targeted ablations.

### 402 6.1. 3D Reconstruction

403 We evaluate OCH3R against Gen3DSR [1], ACDC [11],  
404 and a unified glued baseline that uses instance masks from  
405 SAM2 [47] and GroundingDINO [39], object poses from  
406 MonoDiff9D [38], and depth from DepthPro [3]. We denote  
407 this pipeline as AoE (Army of Experts). Since the baselines  
408 are extremely slow, we randomly sample ten images from  
409 each dataset in our benchmark.

410 Following prior work [1, 73], we report chamfer dis-  
411 tance and F-1@0.1 between the predicted and ground truth  
412 meshes, and CLIP similarity between the rendered and  
413 ground truth images. All backgrounds are manually nor-  
414 malized to white for both predictions and ground truth.

415 As shown in Tab. 1, our method establishes a clear mar-  
416 gin over all baselines. On PACE, OCH3R reduces the  
417 Chamfer Distance from 0.31 (Gen3DSR) and 0.35 (AoE)  
418 to 0.18, while more than doubling the best F-1 score (45.00  
419 versus AoE’s 21.39). Similar trends hold across all remain-  
420 ing datasets: on YCB-V, OCH3R achieves 0.17 CD (a 26  
421 percent improvement over AoE and a 48 percent improve-  
422 ment over Gen3DSR) and reaches 22.71 F-1, surpassing the  
423 strongest baseline by over 10 points. On HOPE, OCH3R  
424 attains 83.69 CLIP similarity, exceeding Gen3DSR by +6.1

and AoE by +19.9. For NOCS real, OCH3R’s gains are the  
most pronounced, improving CD from 0.15 to 0.07 and F-1  
from 38.01 to 76.77. These accuracy improvements come  
alongside a dramatic speedup: our 0.7 s inference time per  
image is roughly 2,000x faster than Gen3DSR (25.6 min)  
and ACDC (22.1 min), while also running more than 30x  
faster than AoE (21.6 s). It demonstrates the advantage of  
our unified per pixel prediction formulation. Some qualita-  
tive results are shown in Fig. 4.

### 6.2. Individual task performance

**Zero-shot metric depth.** Accurate metric depth from a sin-  
gle RGB image is essential to our pipeline, as it defines  
the anchor positions for our 3D Gaussians. We evaluate  
OCH3R on five datasets against seven state-of-the-art base-  
lines using three standard metrics:  $\delta_1$  [33], AbsRel, and  
RMSE. Additional metrics ( $\delta_2$ ,  $\delta_3$ ,  $\log_{10}$ ,  $\text{RMSE}_{\log}$ , SI-  
log) are provided in the Supplementary.

As shown in Tab. 2, OCH3R achieves leading perfor-  
mance on all three metrics for PACE, HOPE, and NOCS-  
real, and on  $\delta_1$  for YCB-V. Although Metric3D V2 [22] and  
Depth Pro [3] perform slightly better on OMNI, the margin  
is minimal; moreover, Metric3D V2 requires ground-truth  
camera intrinsics, giving it an inherent advantage, yet it is  
still surpassed by OCH3R on four of the five datasets. Depth  
Anything V2 [72], trained primarily for relative depth and  
fine-tuned for metric estimation, shows high domain sen-  
sitivity, excelling on YCB-V (narrowly ahead of OCH3R)  
but degrading substantially elsewhere. Overall, OCH3R at-  
tains the best results in 10 of 15 metric-dataset combina-  
tions and delivers the strongest average performance across  
benchmarks.

**Zero-shot semantic segmentation.** Open-vocabulary se-  
mantic segmentation assigns per-pixel labels drawn from a  
potentially open set of natural-language concepts. To build  
an evaluation vocabulary disjoint from training, we aggre-  
gate names from the test datasets, common indoor cate-  
gories from ADE20K [84, 85], and additional household  
items.

We report standard OVSS metrics: mIoU and FB-IoU.  
Given the difficulty of segmenting fine-grained, cluttered  
indoor scenes, we additionally report hit@5, which allows  
each method to produce up to five candidate labels per pixel.  
Since FB-IoU already captures the ability to separate fore-  
ground from background, the remaining metrics are com-  
puted on foreground regions only.

Tab. 3 compares OCH3R with seven OVSS baselines.  
Across five datasets and three metrics (15 settings), OCH3R  
ranks first in 12 and achieves the best average rank (1.27).  
It leads all three metrics on PACE, YCB-V, and NOCS-Real;  
mIoU and FB-IoU on HOPE; and FB-IoU on OMNI. Aver-  
aged across datasets, OCH3R obtains 11.04 mIoU, 83.18  
FB-IoU, and 83.56 hit@5, outperforming the strongest

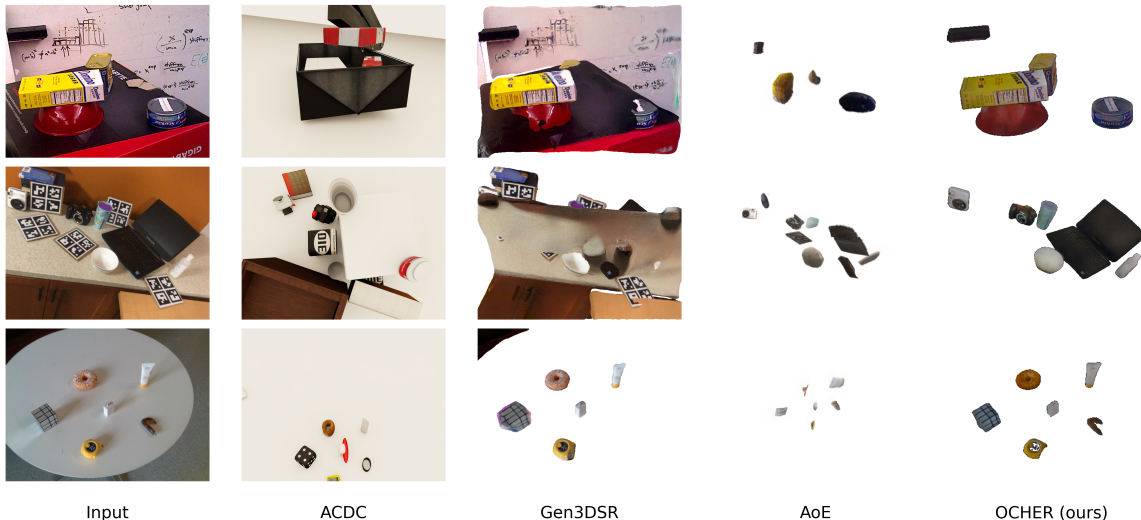


Figure 4. **Qualitative comparison of single-image 3D object-centric reconstruction.** Given a single RGB input, we compare our method (OCH3R) with ACDC, Gen3DSR, and AoE (Army of Experts: SAM2 + GroundingDINO + MonoDiff9D + DepthPro). Prior methods often yield incomplete geometry, distorted textures, or missing objects. OCH3R reconstructs sharper, more complete, and semantically consistent objects across diverse scenes.

Model	PACE			OMNI			YCB-V			HOPE			NOCS real			Time per image↓
	CD↓	F-1↑	CLIP↑	CD↓	F-1↑	CLIP↑	CD↓	F-1↑	CLIP↑	CD↓	F-1↑	CLIP↑	CD↓	F-1↑	CLIP↑	
ACDC [11]	0.69	0.00	73.74	-	-	72.23	0.32	7.76	60.15	1.42	1.92	63.75	3.44	1.72	69.74	22.1 min
Gen3DSR [1]	<u>0.31</u>	12.93	<u>81.61</u>	-	-	<u>82.00</u>	0.33	8.85	<u>75.26</u>	0.35	12.05	<u>77.63</u>	0.26	17.73	<u>75.33</u>	25.6 min
AoE	0.35	<u>21.39</u>	77.97	-	-	69.31	<u>0.23</u>	<u>12.45</u>	66.05	<u>0.18</u>	<b>47.71</b>	63.80	<u>0.15</u>	<u>38.01</u>	70.62	21.6 s
<b>OCH3R (ours)</b>	<b>0.18</b>	<b>45.00</b>	<b>83.15</b>	-	-	<b>82.78</b>	<b>0.17</b>	<b>22.71</b>	<b>78.80</b>	<b>0.14</b>	<u>40.08</u>	<b>83.69</b>	<b>0.07</b>	<b>76.77</b>	<b>85.90</b>	<b>0.7 s</b>

Table 1. Comparison of 3D reconstruction and semantic consistency across datasets using CD (Chamfer Distance, lower is better), F-1 score, and CLIP similarity.

477 baseline by +2.44 mIoU (MAFT+ [27]), +36.80 FB-IoU  
478 (SAN [70]), and +2.47 hit@5 (MAFT+).

479 **Zero-shot pose estimation.** We further evaluate OCH3R  
480 's ability to recover category-level 6D object poses from  
481 a single RGB image in a zero-shot setting. Due to  
482 space constraints, the full quantitative comparison with AG-  
483 Pose [37], SecondPose [7], and MonoDiff9D [38] across  
484 five indoor benchmarks is provided in the Supplementary.  
485 AG-Pose and SecondPose require RGB-D inputs, so we  
486 supply them with our predicted depths. Following MonoD-  
487 iff9D, we report accuracy within 10 cm, within 10°, and  
488 under the joint 10°/10 cm criterion.

489 Across all datasets, OCH3R shows consistent gains on  
490 the stricter angular and joint metrics. It attains the high-  
491 est 10° and joint accuracy on PACE and HOPE and im-  
492 proves the joint metric on NOCS-real. For example, on  
493 PACE, OCH3R improves the 10° rate from 15.1 to 25.9 and  
494 the joint 10°/10 cm rate from 8.6 to 14.0 compared to AG-  
495 Pose. These results indicate that the unified 3D representa-  
496 tion learned by OCH3R naturally supports precise, canoni-  
497 cally aligned object poses without any dataset-specific fine-  
498 tuning.

In summary, our model's strong performance across  
monocular depth estimation, open-vocabulary semantic  
segmentation, and pose estimation jointly enables state-of-  
the-art 3D reconstruction quality, producing geometrically  
precise, semantically coherent, and canonically aligned  
scene representations.

### 6.3. Ablation

**Multi-task learning.** To demonstrate the advantage of  
using a unified model for multiple traditionally separated  
tasks, we retrain the model while removing one head at a  
time. Experiments shows that dropping semantic embed-  
dings causes large pose and Gaussian degradations and also  
hurts depth. More broadly, removing any head weakens  
the remaining tasks, and the full four-head variant performs  
best. Full per-dataset quantitative results are provided in the  
Supplementary.

**Offset in camera space vs. directly in canonical space.**  
We also tried predicting Gaussian parameters directly in  
the canonical frame, bypassing the offset-along-ray formu-  
lation. As shown in Fig. 5, removing the geometric scaf-  
fold of camera rays causes the network to collapse into

Model	PACE			OMNI			YCB-V			HOPE			NOCS real		
	$\delta_1\uparrow$	AbsRel $\downarrow$	RMSE $\downarrow$	$\delta_1\uparrow$	AbsRel $\downarrow$	RMSE $\downarrow$	$\delta_1\uparrow$	AbsRel $\downarrow$	RMSE $\downarrow$	$\delta_1\uparrow$	AbsRel $\downarrow$	RMSE $\downarrow$	$\delta_1\uparrow$	AbsRel $\downarrow$	RMSE $\downarrow$
NeWCRFs [78]	17.00	0.7356	0.5308	21.41	0.6432	0.9934	52.71	0.2740	0.2699	5.37	0.8073	0.5029	36.00	0.3352	0.3586
ZoeDepth [2]	8.20	0.8508	0.6146	20.73	0.5794	0.9406	52.84	0.2704	0.2750	0.53	0.9814	0.6055	34.08	0.4399	0.5275
Metric3D V2 [22]	1.32	0.4695	0.3995	<b>73.73</b>	<b>0.1788</b>	<u>0.7963</u>	15.44	0.2791	0.2538	0.01	0.4039	0.2968	8.50	0.3273	0.3571
Depth Anything V2 [72]	13.07	0.7373	0.5906	5.83	0.9009	1.1728	50.74	<u>0.1899</u>	<b>0.2029</b>	0.07	1.0364	0.6518	2.44	0.5507	0.6088
Depth Pro [3]	33.63	0.5603	0.5610	<u>57.37</u>	0.2742	<b>0.7940</b>	49.04	0.2078	0.2148	23.09	<u>0.4035</u>	0.3013	81.58	0.1370	0.1831
VGGT [64]	<u>55.58</u>	<u>0.2415</u>	<u>0.2096</u>	57.14	<u>0.2272</u>	1.0365	<b>63.88</b>	<b>0.1822</b>	0.2131	<u>37.17</u>	0.4313	<u>0.2844</u>	<u>92.26</u>	<u>0.0965</u>	<u>0.1394</u>
UniDepth V2 [43]	35.05	0.5076	0.4820	22.80	0.7221	1.1007	32.70	0.2548	0.2499	8.63	0.7622	0.6830	78.92	0.1778	0.2366
<b>OCH3R (ours)</b>	<b>94.82</b>	<b>0.0834</b>	<b>0.1039</b>	42.77	0.2900	0.8042	<b>69.96</b>	0.1933	<u>0.2044</u>	<b>61.36</b>	<b>0.2192</b>	<b>0.1812</b>	<b>98.43</b>	<b>0.0923</b>	<b>0.1066</b>

Table 2. **Monocular metric depth estimation results** on PACE, OMNI, YCB-V, HOPE, and NOCS real. Each dataset block reports  $\delta_1$  (in percentage), AbsRel, and RMSE. Bold indicates the best result, and underline indicates the second best.

Model	PACE			OMNI			YCB-V			HOPE			NOCS real		
	mIoU $\uparrow$	FB-IoU $\uparrow$	hit@5 $\uparrow$	mIoU $\uparrow$	FB-IoU $\uparrow$	hit@5 $\uparrow$	mIoU $\uparrow$	FB-IoU $\uparrow$	hit@5 $\uparrow$	mIoU $\uparrow$	FB-IoU $\uparrow$	hit@5 $\uparrow$	mIoU $\uparrow$	FB-IoU $\uparrow$	hit@5 $\uparrow$
LSeg [34]	0.26	17.23	8.52	2.43	4.55	16.31	2.63	12.41	23.05	0.30	12.09	20.69	5.48	8.73	79.82
OVSeg [36]	1.28	13.57	28.84	11.24	2.94	52.19	1.29	12.20	58.66	0.31	11.87	25.99	1.87	8.54	79.21
ODISE [69]	1.69	13.70	36.05	13.27	3.66	58.54	4.51	12.17	76.76	0.99	11.87	68.77	6.31	8.53	90.60
FC-CLIP [77]	<u>3.67</u>	13.56	61.62	<u>22.29</u>	3.61	71.11	3.69	12.15	69.96	1.85	11.91	71.80	4.72	9.40	94.68
CAT-Seg [9]	1.26	21.48	30.53	9.61	4.14	52.62	2.38	12.12	65.78	2.87	11.91	75.97	5.31	8.79	95.80
SAN [70]	1.70	<u>71.05</u>	31.87	13.75	<u>40.82</u>	53.09	4.25	<u>44.77</u>	72.83	<u>3.27</u>	<u>29.96</u>	63.94	5.22	<u>45.30</u>	94.56
MAFT+ [27]	3.66	13.76	<u>64.39</u>	<b>23.02</b>	3.61	<b>76.69</b>	<u>5.00</u>	12.16	<u>79.01</u>	3.21	11.87	<b>89.23</b>	<u>8.10</u>	8.53	<u>96.14</u>
<b>OCH3R (ours)</b>	<b>7.34</b>	<b>94.95</b>	<b>84.76</b>	18.43	<b>84.21</b>	<u>75.23</u>	<b>9.14</b>	<b>72.00</b>	<b>80.03</b>	<b>6.90</b>	<b>71.78</b>	<u>79.27</u>	<b>13.40</b>	<b>92.95</b>	<b>98.51</b>

Table 3. **Open-vocabulary semantic segmentation results** on PACE, OMNI, YCB-V, HOPE, and NOCS real. Each dataset block reports mIoU, FB-IoU, and hit@5 in percentages. Bold indicates the best result, and underline indicates the second best.

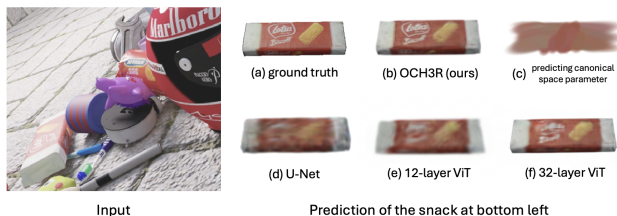


Figure 5. Qualitative ablations showing that (i) predicting Gaussians directly in canonical space collapses, and (ii) OCH3R’s formulation remains robust across model scales and architectures.

520 an unstructured blob, confirming that unconstrained pixel-  
521 to-Gaussian mapping is highly underdetermined. By in-  
522 stead predicting per-pixel offsets in camera space and using  
523 Canonical-Space Supervision to organize the final layout,  
524 OCH3R obtains a stable and expressive inductive bias that  
525 enables high-quality reconstruction.

526 **Model scale and architecture.** The core of our  
527 method, particularly the multi-task learning paradigm and  
528 Canonical-Space Supervision, is orthogonal to model archi-  
529 tecture, so in principle any dense predictor could be used.  
530 We validate this by testing a U-Net (matched in parameter  
531 count to the 32-layer ViT) and two smaller ViT variants. As  
532 shown in Fig. 5, while larger models yield sharper textures  
533 and cleaner geometry, all architectures successfully recover  
534 the object’s overall shape. This confirms that OCH3R’s  
535 pipeline transfers across dense predictors and scales well  
536 with model capacity.

## 7. Conclusion

537  
538 In this paper, we try to address the long standing problem  
539 of understanding a scene as a set of discrete, posed objects  
540 from a single RGB image. Rather than relying on multi  
541 stage pipelines that decompose the task into segmentation,  
542 per object reconstruction, and post hoc alignment, we pro-  
543 posed OCH3R, a unified, feed forward model that predicts  
544 all object instances, their category level SIM(3) poses, and  
545 high fidelity 3D Gaussians in a single pass. The key ingredi-  
546 ents are a transformer that produces dense, pixel aligned at-  
547 tributes (metric depth, CLIP based semantics, NOCS coor-  
548 dinates, and per pixel Gaussians), a simple inference proce-  
549 dure for instance discovery and pose estimation, and canoni-  
550 cal space supervision that trains amodally complete Gaus-  
551 sians without per image Gaussian labels.

552 To support this setting, we assembled a large scale  
553 dataset for holistic object centric 3D scene representation  
554 by aligning PACE, OmniDPose, HOPE, YCB-Video and  
555 NOCS into a unified benchmark with per instance masks,  
556 semantics, 6D poses, and 3D models. Across this bench-  
557 mark, OCH3R outperforms previous baselines on monocu-  
558 lar depth estimation, open vocabulary segmentation, and  
559 category level pose prediction, while producing more com-  
560 plete, editable 3D object reconstructions with feed forward  
561 inference that scales essentially independently of the num-  
562 ber of objects.

563

## References

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

- [1] Andreea Ardelean, Mert Özer, and Bernhard Egger. Generalizable 3d scene reconstruction via divide and conquer from a single view. In *International Conference on 3D Vision (3DV)*, pages 616–626, 2025. 3, 6, 7
- [2] Shariq Farooq Bhat, Reiner Birkel, Diana Wofk, Peter Wonka, and Matthias Müller. Zoedepth: Zero-shot transfer by combining relative and metric depth, 2023. 8
- [3] Aleksei Bochkovskii, Amaël Delaunoy, Hugo Germain, Marcel Santos, Yichao Zhou, Stephan R. Richter, and Vladlen Koltun. Depth pro: Sharp monocular metric depth in less than a second, 2025. 5, 6, 8
- [4] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, Pete Florence, Chuyuan Fu, Montse Gonzalez Arenas, Keerthana Gopalakrishnan, Kehang Han, Karol Hausman, Alexander Herzog, Jasmine Hsu, Brian Ichter, Alex Irpan, Nikhil Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Lisa Lee, Tsang-Wei Edward Lee, Sergey Levine, Yao Lu, Henryk Michalewski, Igor Mordatch, Karl Pertsch, Kanishk Rao, Krista Reymann, Michael Ryoo, Grecia Salazar, Pannag Sanketi, Pierre Sermanet, Jaspier Singh, Anikait Singh, Radu Soricut, Huong Tran, Vincent Vanhoucke, Quan Vuong, Ayzan Wahid, Stefan Welker, Paul Wohlhart, Jialin Wu, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. Rt-2: Vision-language-action models transfer web knowledge to robotic control, 2023. 6
- [5] David Charatan, Sizhe Li, Andrea Tagliasacchi, and Vincent Sitzmann. pixelsplat: 3d gaussian splats from image pairs for scalable generalizable 3d reconstruction, 2024. 2, 3
- [6] Anpei Chen, Haofei Xu, Stefano Esposito, Siyu Tang, and Andreas Geiger. Lara: Efficient large-baseline radiance fields, 2024. 3
- [7] Yamei Chen, Yan Di, Guangyao Zhai, Fabian Manhardt, Chenyangguang Zhang, Ruida Zhang, Federico Tombari, Nassir Navab, and Benjamin Busam. Secondpose: Se(3)-consistent dual-stream feature fusion for category-level pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9959–9969, 2024. 7
- [8] Yuedong Chen, Haofei Xu, Chuanxia Zheng, Bohan Zhuang, Marc Pollefeys, Andreas Geiger, Tat-Jen Cham, and Jianfei Cai. *MVSplat: Efficient 3D Gaussian Splatting from Sparse Multi-view Images*, page 370–386. Springer Nature Switzerland, 2024. 2, 3
- [9] Seokju Cho, Heeseong Shin, Sunghwan Hong, Anurag Arnab, Paul Hongsuck Seo, and Seungryong Kim. Cat-seg: Cost aggregation for open-vocabulary semantic segmentation, 2024. 8
- [10] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes, 2017. 6
- [11] Tianyuan Dai, Josiah Wong, Yunfan Jiang, Chen Wang, Cem Gokmen, Ruohan Zhang, Jiajun Wu, and Li Fei-Fei. Automated creation of digital cousins for robust policy learning. In *Conference on Robot Learning (CoRL)*, 2024. 6, 7
- [12] Laura Downs, Anthony Francis, Nate Koenig, Brandon Kinman, Ryan Hickman, Krista Reymann, Thomas B. McHugh, and Vincent Vanhoucke. Google scanned objects: A high-quality dataset of 3d scanned household items, 2022. 2, 6
- [13] Ruofei Du, Eric Turner, Maksym Dzitsiuk, Luca Prasso, Ivo Duarte, Jason Dourgarian, Joao Afonso, Jose Pascoal, Josh Gladstone, Nuno Cruces, Shahram Izadi, Adarsh Kowdle, Konstantine Tsotsos, and David Kim. Depthlab: Real-time 3d interaction with depth maps for mobile augmented reality. In *Proceedings of the 33rd Annual ACM Symposium on User Interface Software and Technology*, page 829–843, New York, NY, USA, 2020. Association for Computing Machinery. 6
- [14] Hao-Shu Fang, Chenxi Wang, Minghao Gou, and Cewu Lu. Graspnet-1billion: A large-scale benchmark for general object grasping. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 6
- [15] Martin A. Fischler and Robert C. Bolles. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Commun. ACM*, 24(6):381–395, 1981. 4
- [16] Yiming Gao, Yan-Pei Cao, and Ying Shan. Surfelfnerf: Neural surfel radiance fields for online photorealistic reconstruction of indoor scenes, 2023. 2
- [17] Georgia Gkioxari, Jitendra Malik, and Justin Johnson. Mesh r-cnn, 2020. 2
- [18] Thibault Groueix, Matthew Fisher, Vladimir G. Kim, Bryan C. Russell, and Mathieu Aubry. Atlasnet: A papier-mâché approach to learning 3d surface generation, 2018. 2
- [19] Can Gümeli, Angela Dai, and Matthias Nießner. Roca: Robust cad model retrieval and alignment from a single image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4022–4031, 2022. 3
- [20] Aleksander Holynski and Johannes Kopf. Fast depth densification for occlusion-aware augmented reality. *ACM Transactions on Graphics (TOG)*, 37(6), 2018. 6
- [21] Yicong Hong, Kai Zhang, Jiuxiang Gu, Sai Bi, Yang Zhou, Difan Liu, Feng Liu, Kalyan Sunkavalli, Trung Bui, and Hao Tan. Lrm: Large reconstruction model for single image to 3d, 2024. 2
- [22] Mu Hu, Wei Yin, Chi Zhang, Zhipeng Cai, Xiaoxiao Long, Hao Chen, Kaixuan Wang, Gang Yu, Chunhua Shen, and Shaojie Shen. Metric3d v2: A versatile monocular geometric foundation model for zero-shot metric depth and surface normal estimation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(12):10579–10596, 2024. 6, 8
- [23] Zehuan Huang, Yuan-Chen Guo, Xingqiao An, Yunhan Yang, Yangguang Li, Zi-Xin Zou, Ding Liang, Xihui Liu, Yan-Pei Cao, and Lu Sheng. Midi: Multi-instance diffusion for single image to 3d scene generation, 2025. 3
- [24] Muhammad Zubair Irshad, Thomas Kollar, Michael Laskey, Kevin Stone, and Zsolt Kira. Centersnap: Single-shot multi-object 3d shape reconstruction and categorical 6d pose and size estimation, 2022. 2, 3

- 677 [25] Muhammad Zubair Irshad, Sergey Zakharov, Rares Ambrus, 734  
678 Thomas Kollar, Zsolt Kira, and Adrien Gaidon. Shapo: Im- 735  
679 plicit representations for multi-object shape, appearance, and 736  
680 pose optimization, 2022. 2, 3 737
- 681 [26] Hamid Izadinia, Qi Shan, and Steven M. Seitz. Im2cad, 738  
682 2017. 2 739
- 683 [27] Siyu Jiao, Hongguang Zhu, Jiannan Huang, Yao Zhao, Yun- 740  
684 chao Wei, and Humphrey Shi. Collaborative vision-text rep- 741  
685 resentation optimizing for open-vocabulary segmentation, 742  
686 2024. 7, 8 743
- 687 [28] Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task 744  
688 learning using uncertainty to weigh losses for scene geome- 745  
689 try and semantics, 2018. 5 746
- 690 [29] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, 747  
691 and George Drettakis. 3d gaussian splatting for real-time 748  
692 radiance field rendering. *ACM Trans. Graph.*, 42(4):139–1, 749  
693 2023. 2, 5 750
- 694 [30] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, 751  
695 and George Drettakis. 3d gaussian splatting for real-time 752  
696 radiance field rendering, 2023. 2, 3 753
- 697 [31] Philipp Krähenbühl and Vladlen Koltun. Efficient inference 754  
698 in fully connected crfs with gaussian edge potentials. *Ad- 755  
699 vances in neural information processing systems*, 24, 2011. 756  
700 4 757
- 701 [32] Weicheng Kuo, Anelia Angelova, Tsung-Yi Lin, and Angela 758  
702 Dai. Mask2cad: 3d shape prediction by learning to segment 759  
703 and retrieve, 2020. 3 760
- 704 [33] Lubor Ladicky, Jianbo Shi, and Marc Pollefeys. Pulling 761  
705 things out of perspective. In *Proceedings of the IEEE Confer- 762  
706 ence on Computer Vision and Pattern Recognition (CVPR)*, 763  
707 2014. 6 764
- 708 [34] Boyi Li, Kilian Q. Weinberger, Serge Belongie, Vladlen 765  
709 Koltun, and René Ranftl. Language-driven semantic seg- 766  
710 mentation, 2022. 3, 8 767
- 711 [35] Jiaxin Li, Zijian Feng, Qi She, Henghui Ding, Changhu 768  
712 Wang, and Gim Hee Lee. Mine: Towards continuous depth 769  
713 mpi with nerf for novel view synthesis, 2021. 2 770
- 714 [36] Feng Liang, Bichen Wu, Xiaoliang Dai, Kunpeng Li, Yinan 771  
715 Zhao, Hang Zhang, Peizhao Zhang, Peter Vajda, and Diana 772  
716 Marculescu. Open-vocabulary semantic segmentation with 773  
717 mask-adapted clip, 2023. 8 774
- 718 [37] Xiao Lin, Wenfei Yang, Yuan Gao, and Tianzhu Zhang. 775  
719 Instance-adaptive and geometric-aware keypoint learning for 776  
720 category-level 6d object pose estimation. In *Proceedings of 777  
721 the IEEE/CVF Conference on Computer Vision and Pattern 778  
722 Recognition*, pages 21040–21049, 2024. 7 779
- 723 [38] Jian Liu, Wei Sun, Hui Yang, Jin Zheng, Zichen Geng, Hos- 780  
724 sein Rahmani, and Ajmal Mian. Monodiff9d: Monocular 781  
725 category-level 9d object pose estimation via diffusion model. 782  
726 In *IEEE International Conference on Robotics and Automa- 783  
727 tion (ICRA)*, 2025. 6, 7 784
- 728 [39] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao 785  
729 Zhang, Jie Yang, Qing Jiang, Chunyuan Li, Jianwei Yang, 786  
730 Hang Su, et al. Grounding dino: Marrying dino with 787  
731 grounded pre-training for open-set object detection. In *Euro- 788  
732 pean conference on computer vision*, pages 38–55. Springer, 789  
733 2024. 6 790
- [40] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, 734  
Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: 735  
Representing scenes as neural radiance fields for view syn- 736  
thesis, 2020. 3 737
- [41] Yinyu Nie, Xiaoguang Han, Shihui Guo, Yujian Zheng, Jian 738  
Chang, and Jian Jun Zhang. Total3dunderstanding: Joint lay- 739  
out, object pose and mesh reconstruction for indoor scenes 740  
from a single image, 2020. 2 741
- [42] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy 742  
Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, 743  
Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mah- 744  
moud Assran, Nicolas Ballas, Wojciech Galuba, Russell 745  
Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael 746  
Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Je- 747  
gou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr 748  
Bojanowski. Dinov2: Learning robust visual features with- 749  
out supervision, 2024. 5 750
- [43] Luigi Piccinelli, Christos Sakaridis, Yung-Hsu Yang, Mat- 751  
tia Segu, Siyuan Li, Wim Abbeels, and Luc Van Gool. 752  
Unidepthv2: Universal monocular metric depth estimation 753  
made simpler, 2025. 8 754
- [44] Stefan Popov, Pablo Bauszat, and Vittorio Ferrari. Corenet: 755  
Coherent 3d scene reconstruction from a single rgb image, 756  
2020. 2 757
- [45] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya 758  
Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, 759  
Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen 760  
Krueger, and Ilya Sutskever. Learning transferable visual 761  
models from natural language supervision, 2021. 2, 3 762
- [46] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vi- 763  
sion transformers for dense prediction, 2021. 5 764
- [47] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang 765  
Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman 766  
Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2: 767  
Segment anything in images and videos. *arXiv preprint 768  
arXiv:2408.00714*, 2024. 6 769
- [48] Mike Roberts, Jason Ramapuram, Anurag Ranjan, Atulit 770  
Kumar, Miguel Angel Bautista, Nathan Paczan, Russ Webb, 771  
and Joshua M. Susskind. Hypersim: A photorealistic syn- 772  
thetic dataset for holistic indoor scene understanding, 2021. 773  
2, 6 774
- [49] Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Perceiver- 775  
actor: A multi-task transformer for robotic manipulation, 776  
2022. 6 777
- [50] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob 778  
Fergus. Indoor segmentation and support inference from 779  
rgbd images. In *Proceedings of the 12th European Confer- 780  
ence on Computer Vision (ECCV)*, pages 746–760, Berlin, 781  
Heidelberg, 2012. Springer. 6 782
- [51] Shuran Song, Samuel P. Lichtenberg, and Jianxiong Xiao. 783  
Sun rgb-d: A rgb-d scene understanding benchmark suite. 784  
In *Proceedings of the IEEE Conference on Computer Vision 785  
and Pattern Recognition (CVPR)*, pages 567–576, 2015. 6 786
- [52] Stanislaw Szymanowicz, Christian Rupprecht, and Andrea 787  
Vedaldi. Splatter image: Ultra-fast single-view 3d recon- 788  
struction. In *Proceedings of the IEEE/CVF conference 789  
on computer vision and pattern recognition*, pages 10208– 790  
10217, 2024. 2, 3 791

- 792 [53] Stanislaw Szymanowicz, Eldar Insafutdinov, Chuanxia  
793 Zheng, Dylan Campbell, Joao F Henriques, Christian Rup-  
794 precht, and Andrea Vedaldi. Flash3d: Feed-forward general-  
795 isable 3d scene reconstruction from a single image. In *2025*  
796 *International Conference on 3D Vision (3DV)*, pages 670–  
797 681. IEEE, 2025. 2
- 798 [54] Stanislaw Szymanowicz, Eldar Insafutdinov, Chuanxia  
799 Zheng, Dylan Campbell, João F. Henriques, Christian Rup-  
800 precht, and Andrea Vedaldi. Flash3d: Feed-forward general-  
801 isable 3d scene reconstruction from a single image, 2025. 2,  
802 3
- 803 [55] Jiaxiang Tang, Zhaoxi Chen, Xiaokang Chen, Tengfei Wang,  
804 Gang Zeng, and Ziwei Liu. Lgm: Large multi-view gaussian  
805 model for high-resolution 3d content creation, 2024. 2, 3
- 806 [56] Jiapeng Tang, Yinyu Nie, Lev Markhasin, Angela Dai, Justus  
807 Thies, and Matthias Nießner. Diffuscene: Denoising diffu-  
808 sion models for generative indoor scene synthesis, 2024. 4
- 809 [57] Maxim Tatarchenko, Alexey Dosovitskiy, and Thomas Brox.  
810 Octree generating networks: Efficient convolutional archi-  
811 tectures for high-resolution 3d outputs, 2017. 2
- 812 [58] Shubham Tulsiani, Tinghui Zhou, Alexei A. Efros, and Ji-  
813 tendra Malik. Multi-view supervision for single-view recon-  
814 struction via differentiable ray consistency, 2017. 2
- 815 [59] Shubham Tulsiani, Richard Tucker, and Noah Snavely.  
816 Layer-structured 3d scene inference via view synthesis,  
817 2018. 2
- 818 [60] Stephen Tyree, Jonathan Tremblay, Thang To, Jia Cheng,  
819 Terry Mosier, Jeffrey Smith, and Stan Birchfield. 6-dof pose  
820 estimation of household objects for robotic manipulation:  
821 An accessible dataset and benchmark, 2022. 6
- 822 [61] S. Umeyama. Least-squares estimation of transformation pa-  
823 rameters between two point patterns. *IEEE Transactions on*  
824 *Pattern Analysis and Machine Intelligence*, 13(4):376–380,  
825 1991. 4
- 826 [62] He Wang, Srinath Sridhar, Jingwei Huang, Julien Valentin,  
827 Shuran Song, and Leonidas J Guibas. Normalized object  
828 coordinate space for category-level 6d object pose and size  
829 estimation. In *Proceedings of the IEEE/CVF conference on*  
830 *computer vision and pattern recognition*, pages 2642–2651,  
831 2019. 2, 4
- 832 [63] He Wang, Srinath Sridhar, Jingwei Huang, Julien Valentin,  
833 Shuran Song, and Leonidas J. Guibas. Normalized object  
834 coordinate space for category-level 6d object pose and size  
835 estimation, 2019. 3, 6
- 836 [64] Jianyuan Wang, Minghao Chen, Nikita Karaev, Andrea  
837 Vedaldi, Christian Rupprecht, and David Novotny. Vggt: Vi-  
838 sual geometry grounded transformer, 2025. 3, 5, 8
- 839 [65] Yunsong Wang, Tianxin Huang, Hanlin Chen, and Gim Hee  
840 Lee. Freesplat: Generalizable 3d gaussian splatting towards  
841 free-view synthesis of indoor scenes, 2024. 2, 3
- 842 [66] Bowen Wen, Wei Yang, Jan Kautz, and Stan Birchfield.  
843 Foundationpose: Unified 6d pose estimation and tracking of  
844 novel objects, 2024. 6
- 845 [67] Yu Xiang, Tanner Schmidt, Venkatraman Narayanan, and  
846 Dieter Fox. Posecnn: A convolutional neural network for 6d  
847 object pose estimation in cluttered scenes. In *Proceedings of*  
848 *Robotics: Science and Systems (RSS)*, 2018. 6
- [68] Haoifei Xu, Songyou Peng, Fangjinhua Wang, Hermann  
Blum, Daniel Barath, Andreas Geiger, and Marc Pollefeys.  
Depthsplat: Connecting gaussian splatting and depth, 2025.  
2
- [69] Jiarui Xu, Sifei Liu, Arash Vahdat, Wonmin Byeon, Xiao-  
long Wang, and Shalini De Mello. Open-vocabulary panop-  
tic segmentation with text-to-image diffusion models, 2023.  
8
- [70] Mengde Xu, Zheng Zhang, Fangyun Wei, Han Hu, and Xi-  
ang Bai. Side adapter network for open-vocabulary semantic  
segmentation, 2023. 7, 8
- [71] Yinghao Xu, Zifan Shi, Wang Yifan, Hansheng Chen,  
Ceyuan Yang, Sida Peng, Yujun Shen, and Gordon Wet-  
zstein. Grm: Large gaussian reconstruction model for ef-  
ficient 3d reconstruction and generation, 2024. 2, 3
- [72] Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiao-  
gang Xu, Jiashi Feng, and Hengshuang Zhao. Depth any-  
thing v2, 2024. 6, 8
- [73] Kaixin Yao, Longwen Zhang, Xinhao Yan, Yan Zeng, Qix-  
uan Zhang, Wei Yang, Lan Xu, Jiayuan Gu, and Jingyi Yu.  
Cast: Component-aligned 3d scene reconstruction from an  
rgb image, 2025. 2, 3, 4, 6
- [74] Botao Ye, Sifei Liu, Haoifei Xu, Xueting Li, Marc Pollefeys,  
Ming-Hsuan Yang, and Songyou Peng. No pose, no problem:  
Surprisingly simple 3d gaussian splats from sparse unposed  
images, 2024. 2, 3, 5
- [75] Yang You, Kai Xiong, Zhening Yang, Zhengxiang Huang,  
Junwei Zhou, Ruoxi Shi, Zhou Fang, Adam W. Harley,  
Leonidas Guibas, and Cewu Lu. Pace: A large-scale dataset  
with pose annotations in cluttered environments, 2024. 2, 6
- [76] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa.  
pixelnerf: Neural radiance fields from one or few images,  
2021. 2
- [77] Qihang Yu, Ju He, Xueqing Deng, Xiaohui Shen, and Liang-  
Chieh Chen. Convolutions die hard: Open-vocabulary seg-  
mentation with single frozen convolutional clip, 2023. 8
- [78] Weihao Yuan, Xiaodong Gu, Zuozhuo Dai, Siyu Zhu, and  
Ping Tan. New crfs: Neural window fully-connected crfs for  
monocular depth estimation, 2022. 8
- [79] Cheng Zhang, Zhaopeng Cui, Yinda Zhang, Bing Zeng,  
Marc Pollefeys, and Shuaicheng Liu. Holistic 3d scene un-  
derstanding from a single image with implicit representation,  
2021. 2
- [80] Chubin Zhang, Hongliang Song, Yi Wei, Yu Chen, Jiwen  
Lu, and Yansong Tang. Geolrm: Geometry-aware large re-  
construction model for high-quality 3d gaussian generation,  
2024. 2
- [81] Jiyao Zhang, Weiyao Huang, Bo Peng, Mingdong Wu, Fei  
Hu, Zijian Chen, Bo Zhao, and Hao Dong. Omni6dpose: A  
benchmark and model for universal 6d object pose estima-  
tion and tracking, 2024. 2, 6
- [82] Kai Zhang, Sai Bi, Hao Tan, Yuanbo Xiangli, Nanxuan Zhao,  
Kalyan Sunkavalli, and Zexiang Xu. Gs-lrm: Large recon-  
struction model for 3d gaussian splatting, 2024. 2, 3
- [83] Qingcheng Zhao, Xiang Zhang, Haiyang Xu, Zeyuan Chen,  
Jianwen Xie, Yuan Gao, and Zhuowen Tu. Depr: Depth  
guided single-view scene reconstruction with instance-level  
diffusion, 2025. 3, 4

- 907 [84] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela  
908 Barriuso, and Antonio Torralba. Scene parsing through  
909 ade20k dataset. In *Proceedings of the IEEE Conference on*  
910 *Computer Vision and Pattern Recognition*, 2017. 6
- 911 [85] Bolei Zhou, Hang Zhao, Xavier Puig, Tete Xiao, Sanja Fi-  
912 dler, Adela Barriuso, and Antonio Torralba. Semantic under-  
913 standing of scenes through the ade20k dataset. *International*  
914 *Journal of Computer Vision*, 127(3):302–321, 2019. 6