# Object-Centric Multi-View Aggregation

Shubham Tulsiani[1]    Or Litany[2†]    Charles R. Qi[†]    He Wang[2†]    Leonidas J. Guibas[2†]

[1]Facebook AI    [2]Stanford University

## Abstract

*We present an approach for aggregating a sparse set of views of an object in order to compute a semi-implicit 3D representation in the form of a volumetric feature grid. Key to our approach is an object-centric canonical 3D coordinate system into which views can be lifted, without explicit camera pose estimation, and then combined – in a manner that can accommodate a variable number of views and is view order independent. We show that computing a symmetry-aware mapping from pixels to the canonical coordinate system allows us to better propagate information to unseen regions, as well as to robustly overcome pose ambiguities during inference. Our aggregate representation enables us to perform 3D inference tasks like volumetric reconstruction and novel view synthesis, and we use these tasks to demonstrate the benefits of our aggregation approach as compared to implicit or camera-centric alternatives.*

## 1. Introduction

A central problem in computer vision is to recover the 3D structure of the world from 2D views of it, namely images. Classical approaches such as Structure from Motion (SfM) operate in the setting where many views (say tens to hundreds) are available so that geometric correspondences between them can be used to infer camera pose and induce 3D structure [29, 15]. More recently there has been a flurry of activity on object reconstruction from a single image, exploiting machine learning and especially deep learning approaches to hallucinate information invisible in that view [7, 3, 4, 30], building on large 3D model repositories such as ShapeNet [2] for training. We aim to tackle the scenario in between these two extremes, and addresses situations where a few images (say one to four) are available *e.g.*, online product marketplaces.

In this work, we focus on the single-object setup: how to recover a '3D representation' of the underlying object given one or more images. The key questions for this task pertain to the form of this 3D representation, and how can one enable aggregation of the information across views to compute a unified representation. While classical methods pursue explicit representations such as point clouds and aggregate views via explicit correspondence inference, these choices are not easily applicable to our setup with a small number of input images with unknown camera poses. In contrast, learning based methods typically represent 3D implicitly [7, 3] *e.g.* via a single latent vector, and can be extended to implicitly aggregate images *e.g.*, via an LSTM [3]. However, this fully implicit representation and aggregation ignores the underlying geometric structure of the task. We instead pursue semi-implicit 3D representations [26], as these combine a 3D voxel grid with an implicit latent vector in each voxel coding the contents of that cell, and allow easily recovering explicit structure. We propose an aggregation mechanism to infer such a semi-latent representation given multiple images, and show that it allows us to perform 3D centric tasks *e.g.* shape inference or novel view synthesis.

Our key insight towards designing the aggregation mechanism is that for each pixel, we can 'lift' it to a canonical object-centric space, and then process the information from across the images in this shared canonical space. Unlike more traditional computer vision techniques that are camera-focused and need to know or estimate the camera pose for each view in a world coordinate system so as to properly integrate them, we instead accomplish view aggregation without explicit camera pose estimation, by directly lifting views into an object-centric space. This further allows exploitation of prior object knowledge that is difficult to incorporate in camera-centric approaches. For example, many objects possess important symmetries that can be inferred from knowledge of the object class and the observed view or views. With such knowledge, when we lift a given view into the canonical object-centric space in 3D, we can augment the observed regions by their inferred symmetric counterparts, effectively enhancing our understanding of the 3D structure of the object and allowing us to produce a complete 3D representation from fewer views. Critically, symmetry induced augmentation also solves the problem of pose ambiguities due to object symmetries or part occlusions (e.g., resolving 'front' or 'back' of a bottle) – by effectively generating in the
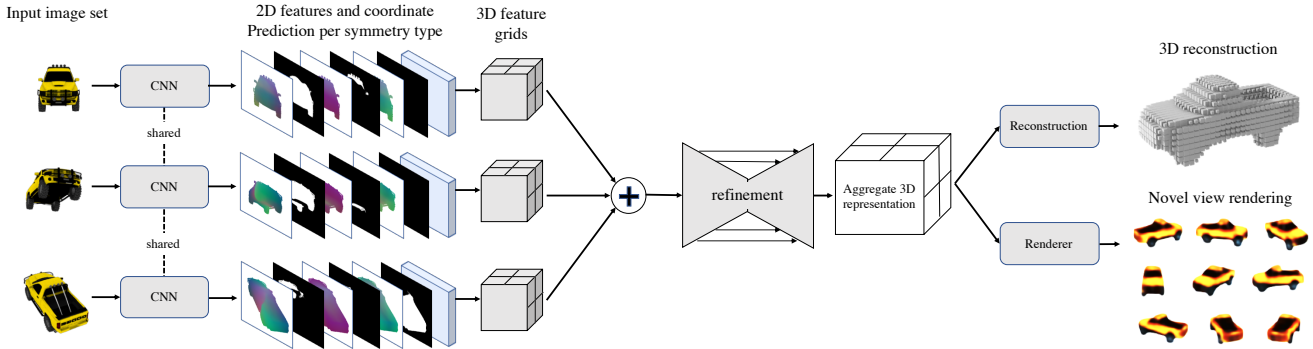
---

**Figure 1: Multi-view aggregation and its applications.** An input image set is first passed through a 2D CNN that outputs per-pixel features, and for each symmetry type a dense (pre-pixel) coordinate prediction (3-dim) and a confidence map (1-dim). Three pairs of predicted coordinates and confidence maps are shown as an illustration. The computed features are then lifted to 3D voxel locations prescribed by the coordinate prediction and weighted by the confidences. A 3D-UNet is then used to refine the averaged 3D features to yield an aggregate 3D representation, which is further used for (a) 3D reconstruction, and (b) novel-view synthesis.

object coordinate space the union of the lift from all possible valid poses, thus making such ambiguities immaterial.

Our overall architecture consists of three parts: a) a network that lifts a single image into a semi-implicit representation of the seen object in a voxelized object-centric 3D space annotated with learned features per voxel, b) a symmetric view aggregator that can combine lifts from different views into a single semi-implicit representation, and c) task specific networks that then use this representation for different downstream tasks, such as volumetric reconstruction or novel view synthesis.

In summary, we present an approach to infer a 3D representation given a sparse set of images with unknown camera poses, and show that this representation can be used for downstream 3D tasks. Our key contributions are:

- By using an object-centric coordinate system we obviate the need for camera pose estimation in 3D inference and view aggregation.

- Our symmetry-aware view lifting into the object-centric space allows us to extract more information from each view, while also bypassing ambiguities.

- Our semi-implicit 3D representation carries a strong geometric inductive bias in its formulation, and combines latent local elements that that can be trained in an end-to-end fashion for multiple downstream tasks.

## 2. Related Work

The task of recovering the underlying 3D structure from multiple images is a classical one in the computer vision community. Early approaches can broadly be categorized as tackling Multi-view Stereopsis (MVS) [25, 5], where the aim is to recover 3D given camera poses, or Structure from Motion (SfM) [31, 9], where the camera pose is also unknown. A long line of work in both setups [1, 27, 6, 24], relying on

cues such as geometric constraints and photometric consistency, has yielded impressive results, but in scenarios where numerous views of a scene/object are available. However, due to the reliance on such cues, these learning-free methods are not well-suited for handling a sparse set of images with possibly little overlap. We instead aim to build a system that can retain the geometric inductive biases of these methods, while also leveraging data-driven priors via learning to tackle inference from sparse set of images with unknown camera poses.

A scenario where these data-driven priors have been successfully exploited by learning based approaches is that of 3D prediction from a single input image. Driven by the success of deep learning, several methods tackle the task of inferring volumetric [3, 7, 30, 34, 10], mesh [13, 21, 16, 33], point cloud [4, 17], or implicit 3D [8, 18] representations from an input image. These approaches learn deep neural network based models that during inference, directly predict the underlying 3D in a feed forward manner given a single image, without explicitly relying on the geometry of the task. Our goal in this work is to extend these approaches, in particular ones for volumetric inference, to leverage multiple input images, and we do so using a geometrically motivated aggregation mechanism.

There have been several recent attempts [14, 11, 35, 20, 23, 12, 19, 36] which, sharing our motivation for integrating geometric inductive biases in a learning framework, tackle the task of learning based multi-view stereopsis. While these demonstrate impressive results, they all crucially rely on the availability of the ground-truth camera poses during inference. We instead pursue the task of 3D reconstruction with only a set of images at inference, without the availability of known camera poses. In contrast to these approaches which rely on global camera estimates and camera-centric predictions *e.g.* depth, we propose to directly predict object-

centric coordinates. Closer to our setup is the work by Choy *et al*. [3] which similarly infers 3D from multiple images without known camera poses, but performs implicit feature aggregation, whereas we propose a more geometrically motivated mechanism. More recently, work from Sridhar et. al. [28] also investigated the task of aggregating multiple views, but in contrast to our approach, did not account for the ambiguities due to symmetry.

## 3. Approach

Given multiple images of an object with associated foreground masks, we aim to compute an aggregate representation that incorporates information from across the images, and then leverage this representation for tasks such as 3D reconstruction and novel view synthesis. Our key insight is that we can first 'lift' the information from the foreground pixels across different images of an instance to an object-centric 3D space, and then aggregate and process the features from across the images in this shared canonical space. This yields a canonicalized 3D representation for the underlying object, which can then be used towards various 3D related tasks.

Our proposed system, as depicted in Figure 1, is comprised of three stages: a) lifting pixels to an object-centric 3D space with symmetry augmentation, b) multi-image feature aggregation, and c) task-specific computation. We first present the lifting procedure in Section 3.1 where we learn a (probabilistic) mapping from pixels to coordinates in a normalized space, while allowing symmetries to overcome ambiguities. We then describe in Section 3.2 how these learned mappings enable copying pixel-wise features from across the images into a volumetric feature grid, where these are aggregated and refined via a learned function. Finally, we show in Section 3.3 that the aggregate object representation can be leveraged for various 3D-centric tasks.

### 3.1. Symmetry-Aware Object-Centric Coordinate Prediction

We aim to compute a mapping from pixels to a canonical 3D space to enable aggregation across images. We build on the insight by Wang *et al*. [32] that such a mapping can be defined using a normalized and aligned shape collection. We briefly summarize the canonical coordinate prediction [32] into an object-centric space, and describe in detail our symmetry-aware formulation that allows overcoming ambiguities in the task, and the corresponding training objective.

**Object-Centric Coordinate Prediction.** Each foreground pixel corresponds to a 3D point on the underlying object surface. Using a normalized shape collection, where all shapes are aligned and scaled to fit unit diagonal cubes, one can learn a mapping from pixels to their corresponding

coordinates in this space. Given an image $I$, we can learn to predict a pixel-wise mapping $C$, where for a pixel $u$, $C[u] \in \mathbb{R}^3$ corresponds to the canonical coordinate of the 3D point visible at that pixel. Using synthetically rendered images, it is easy to obtain the ground-truth mapping $\hat{C}$ and learn a parametric prediction function $f_\theta(I) \equiv C$ using such supervision. While Wang *et al*. [32] leveraged such a mapping for the task of pose estimation from a single image, we observe that it can also allow us to aggregate information across several images of an instance via lifting pixels (and associated features) into this canonical space.

**Overcoming Ambiguities via Allowing Symmetries.** Unfortunately, inferring this mapping from pixels to their canonical coordinates for generic objects is an inherently ambiguous task. Consider a pixel on the leg of a symmetrical square table with four legs. While this pixel does have a unique canonical coordinate (as defined by the the position of the corresponding 3D point in the normalized object model), it is not possible to infer this coordinate given an image alone. Such ambiguities are common due to local and global symmetries across objects.

Our insight is that instead of predicting a single canonical coordinate for each pixel, we can instead learn to predict a *symmetry-aware distribution*. This allows us to both: a) overcome mean prediction effects when dealing with ambiguities, b) propagate information from a pixel to multiple 3D locations depending on the inferred symmetries. We therefore make predictions of coordinates and probabilities for multiple symmetry types. We assume a set of possible symmetry types $\mathcal{G}$, with each $g \in \mathcal{G}$ indicating either a rotational or reflection symmetry. We overload notation, and also use $g$ to denote a function that, given an input point $x \in \mathbb{R}^3$, generates the corresponding set of points under the symmetry type $g$, *i.e.* $g(x)$ is the closure set for $x$ under $g$. Concretely, we consider 5 global symmetry types: identity, reflection along $y$ axis, and $2, 4$, or continuous rotational symmetry along z-axis. The first type corresponds to no symmetry prediction, while the others correspond to some commonly occurring global symmetries across objects.

Given an input image $I$, we predict $f_\theta(I) \equiv \{(C^g, P^g)\}$, where for a pixel $u$, $P^g[u]$ denotes the probability that the underlying 3D point belongs to the symmetry type $g$ (the per-pixel probabilities sum to 1), and $C^g[u]$ indicates its object-centric canonical coordinate if it does. Note that this is equivalent to predicting the set of points $g(C^g[u])$ with probability $P^g[u]$ at the pixel $u$.

**Training Objective.** We can learn this per-symmetry type canonical coordinate prediction using supervision in the form of the ground-truth coordinates $\hat{C}$. Note that this is similar supervision as in the case of learning canonical prediction without allowing symmetries, hence *we do not need to rely*

*on explicit supervision for the per-type coordinate or probability predictions.* Our training objective comprises two terms that work together to (a) encourage the predicted symmetry type at each pixel to contain the correct coordinate, while (b) penalizing spurious predictions.

For a pixel $u$, we penalize the distance between its associated canonical coordinate $\hat{C}[u]$ and the closest point in each symmetry type predicted, and weight this loss by the corresponding probability.

$$L_c = \sum_u \sum_g P^g[u] \min_{x \in g(C^g[u])} \|x - \hat{C}[u]\|. \quad (1)$$

This loss enables overcoming ambiguities inherent in the task due to symmetries, as instead of requiring the prediction of the true canonical coordinate, it only requires predicting a possible canonical coordinate in the respective symmetry type.

However, this encourages over-predicting symmetry types as the additional points can only reduce the loss. We therefore introduce a second loss that reduces spurious predictions by penalizing them for inducing points that do not exist in the underlying 3D shape. Denoting by $\mathcal{D}(S, x)$ the distance from a point $x$ to its closet point on a shape $S$, the additional objective is:

$$L_s = \sum_u \sum_g P^g[u] \max_{x \in g(C^g[u])} \mathcal{D}(S, x). \quad (2)$$

We use a UNet [22] based CNN as the parametrized predictor $f_\theta$ to learn the symmetry-aware canonical coordinate prediction. As an implementation detail, while it is easy to analytically compute the objective in Eq. 1 for the symmetry types we consider, we use a finite number of random samples for symmetry types that induce sets with infinite cardinality for the objective in Eq. 2.

### 3.2. Multi-Image Feature Aggregation

Given the learned (probabilistic) mapping from pixels across the images to a canonical 3D space, we can compute an aggregate representation in the form of a volumetric feature grid. We do so via first computing 2D per-pixel features across the input images and lift these features to a shared 3D grid using the inferred embeddings. We can then combine and further process the features from the different images in this volumetric space to obtain an aggregate representation that encompasses the information from across the input images. Concretely, given $K$ images $\{I_k\}$ of an instance, we first compute corresponding 2D features $\{F_k\}$. We then lift these features using the predicted canonical coordinates to a per-image volumetric feature grid $V_k$, and then aggregate these to obtain a volumetric feature representation $V$.

**2D Feature Extraction.** We want a 2D encoder that can capture global context while preserving the low-level details.

Towards this, we extend the UNet [22] $f_\theta$ presented in Section 3.1 to additionally output per-pixel features $F$ given input image $I$.

**Probabilistic Feature Splatting.** The predicted canonical coordinates associated with image $I_k$, $\{(C_k^g, P_k^g)\} \equiv f_\theta(I_k)$ allow us to lift the associated 2D features $F_k$ to a volumetric feature grid $V_k$. Intuitively, we start with an empty 3D feature grid and for each pixel, we add the associated 2D feature (weighted by the corresponding symmetry type probability) at the 3D location(s) implied by the canonical coordinate prediction. Note that a single pixel may lead to placing features at multiple 3D locations based on the underlying symmetry type.

We denote by $\mathcal{V}(x, f)$ a 3D feature grid obtained by placing a feature $f$ at the 3D coordinate $x$ in an initially empty grid. Note that this grid is empty at all locations except up to 8 cells immediately around the coordinate $x$ (see appendix for details). Using this notation, we can define our lifted feature grid as the probability weighted combination of all the 3D grids obtained for each pixel.

$$V_k = \sum_u \sum_g \sum_{x \in g(C_k^g[u])} P_k^g[u] \, \mathcal{V}(x, F_k[u]). \quad (3)$$

This procedure allows us to copy features from pixels to possibly multiple locations as implied by the symmetry predictions. Further, for use in aggregating feature grids across images, we also compute a 'weight' grid that records the total number of pixels that contribute to each cell (weighted by probabilities).

$$W_k = \sum_u \sum_g \sum_{x \in g(C_k^g[u])} P_k^g[u] \, \mathcal{V}(x, 1). \quad (4)$$

**Averaging and Refinement.** Having obtained feature and weight grids $\{(V_k, W_k)\}$ for each input image $I_k$, we can now construct a sum weight and an average feature across the images:

$$\bar{W} = \sum_k W_k \;\; ; \;\; \bar{V} = \frac{\sum_k V_k}{\bar{W}}, \quad (5)$$

Where the division is understood as a voxel-wise operation. Finally, we use a 3D UNet based CNN $h_\psi$ to process the features and yield a final aggregate representation $V$ that incorporates information from all the views. This additional processing allows us to perform 3D reasoning using the lifted 2D features, and can implicitly perform noise filtering *etc.* and also propagate information to regions of the 3D volume without any direct image evidence.

$$V = h_\psi([\bar{V}; \bar{W}]). \quad (6)$$

## 3.3. Learning 3D-Centric Tasks

Our aggregate feature representation $V$ is a volumetric feature grid that integrates information from multiple images of the object. We leverage this representation for learning 3D tasks, and show how we can train our system using tasks like volumetric prediction and novel view synthesis as supervision.

**Volumetric 3D Prediction.**  A task pursued by previous multi-image prediction systems [3] is that of 3D reconstruction. We show that our aggregate representation can also be used for this task, and that our approach improves over fully implicit aggregation mechanisms. We use a lightweight 2-layer 3D CNN to predict voxel occupancy $O$ from our volumetric representation and use a cross-entropy loss $L_{vol}$ between the ground-truth and predicted occupancies.

**Novel View Synthesis.**  To demonstrate that the aggregate representation can capture the *appearance* of 3D objects, we synthesize novel views from it.

We adopt the pipeline from [26] and train a renderer $\mathcal{R}$ that yields an image $\mathcal{R}(V, \pi)$ given an input feature grid $V$ and camera viewpoint $\pi$.

While [26] used this renderer in conjunction with an optimized feature grid that required hundreds of input views with known camera poses, we adopt it to our setting where the 3D representation is predicted from only few images with unknown camera poses. As supervision for training, we assume novel views of the object $\{I_{k'}, \pi_{k'}\}$ where the image $I_{k'}$ has a corresponding camera viewpoint $\pi_{k'}$. Our view synthesis loss is defined as: $L_{vs} = \sum_{k'} \|\mathcal{R}(V, \pi_{k'}) - I_{k'}\|$. Note that while the novel views used for supervising the rendering have known camera viewpoints, the images used to compute the aggregate representation do not.

**Overall Training Objective.**  Our training objective comprises terms for learning the canonical coordinate prediction $(L_s, L_c)$ as well as task-specific reconstruction and rendering objectives $(L_{vol}, L_{vs})$. We weight the loss terms to (approximately) equalize their contribution to the total loss. Additionally, we find it beneficial to decouple the learning of the coordinate prediction from the downstream tasks, *i.e.* the task-specific losses only influence the learned feature representations. We learn a common model across all categories, trained jointly for both tasks (reconstruction and view synthesis). We will publicly release our implementation for reproducibility.

## 4. Experiments

### 4.1. Training Setup

**Dataset.**  We use the ShapeNet dataset [2] for empirical validation of our approach. We use models from 13 object categories (similar to previous approaches [3, 14]), and use random train/val/test splits with (0.7, 0.1, 0.2) fraction of the models. We render each instance from 10 randomly sampled camera viewpoints with azimuth $\in [0, 360)$, elevation $\in [-20, 40]$ degrees, and an additional random camera translation $\in [-0.1, 0.1]$ units in each dimension (where ShapeNet models lie in a unit diameter ball). We train our model (and all baselines) using $K = 4$ input images during training, and use $K' = 5$ images as supervision for learning novel view synthesis. We use voxelized representations of the models with a grid size 32 for training and evaluating the volumetric reconstruction.

**Baselines.**  We compare our method with several approaches that aggregate multiple images:

*a) 3D-R2N2*: an implicit LSTM based aggregation. To enable view synthesis we extend the model originally presented in [3] with a learned decoder to upsample the hidden state to a spatial resolution (similar to our representation $V$), followed by a renderer as described in Section 3.3.

*b) Depth and Pose based Aggregation (DnP)*: Instead of our symmetry-aware object-centric coordinates from 2D images, one can predict a camera-centric per-pixel depth and the global camera pose. We therefore present a baseline which replaces our canonical coordinate prediction in Section 3.1 with depth and pose (trained with corresponding supervision), while keeping all other aspects unchanged.

*c) Ours (w/o symmetry)*: To highlight the importance of allowing possible symmetries, we present a baseline which does not leverage symmetry, and instead predicts a single coordinate per pixel.

### 4.2. Evaluation

**Volumetric Reconstruction.**  We measure the performance of methods using the intersection over union (IoU) between the predicted and ground-truth $32^3$ volumes. We report the mean IoU score across the 13 categories. As all approaches predict a continuous probability, we report performance for each method using the corresponding optimal binarization threshold (typically around $0.4$). The performance of various approaches in the setting with 4 input views is reported in Table 1.

We consistently improve over the alternatives of leveraging implicit aggregation or camera-centric prediction. We also clearly see the benefits of incorporating symmetry, in particular for classes such as bench, lamp, and table. While all methods were trained with exactly 4 input views, we also test their performance with fewer/more views at inference
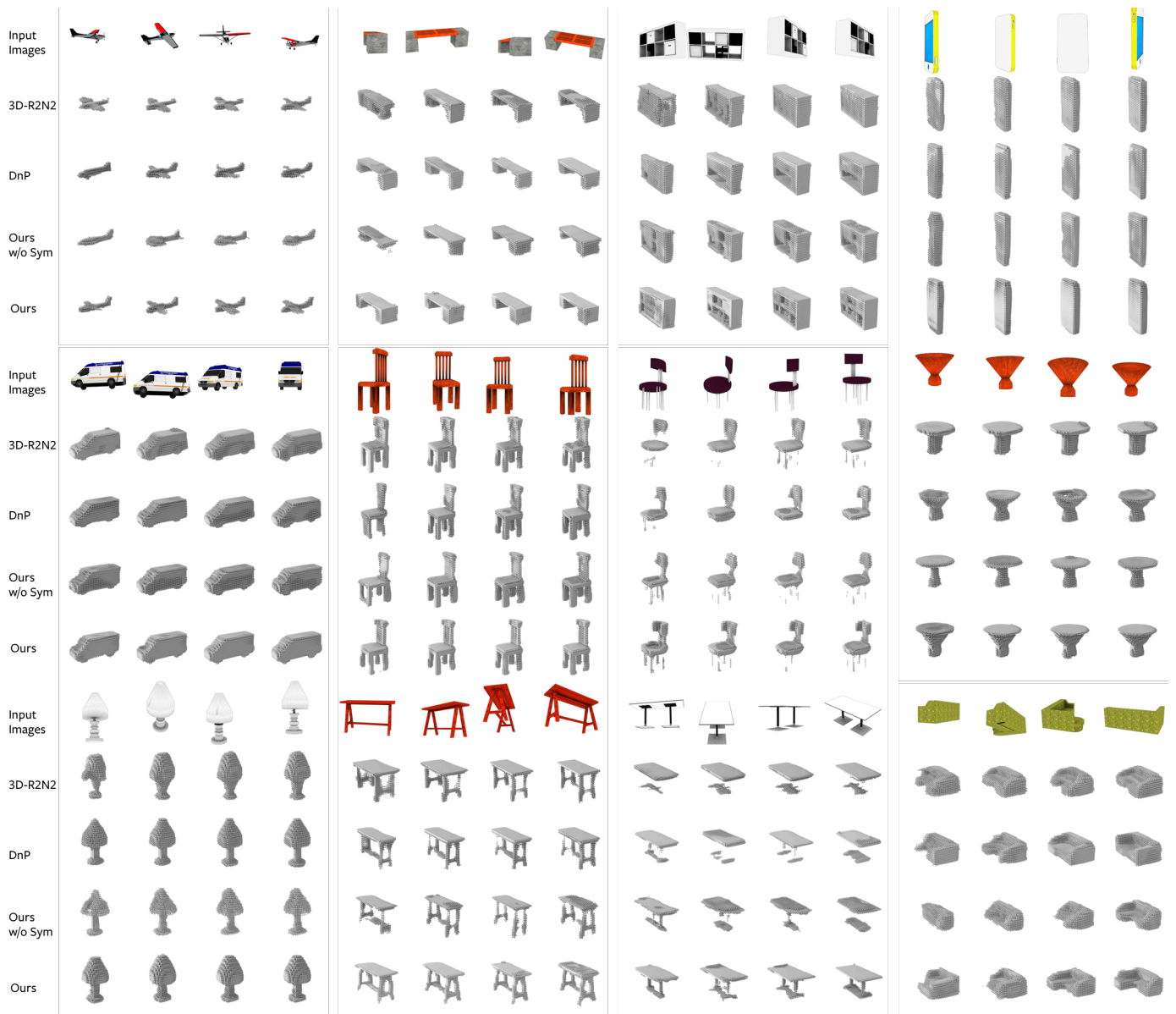
**Figure 2: Volumetric prediction**. We show results for 8 example objects. For each example, the top row shows the 4 input views. The second to fifth rows show reconstructed shapes from different methods. The columns correspond to results when using the initial 1, 2, 3, or all 4 views as input.

| Classes | aero | bench | cabinet | car | chair | display | lamp | speaker | rifle | sofa | table | phone | vessel | mean |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **3D-R2N2** [3] | 58.0 | 55.0 | 71.5 | 79.3 | 57.2 | 51.4 | 46.4 | 65.2 | 61.4 | 67.7 | 60.0 | 70.4 | 59.5 | 61.8 |
| **DnP** | 57.6 | 56.1 | 70.7 | 80.1 | 61.8 | **54.1** | 46.4 | 66.2 | 65.2 | 70.5 | 60.1 | 69.8 | 58.8 | 62.9 |
| **Ours (w/o sym)** | 55.8 | 52.0 | 69.5 | 78.7 | 58.2 | 49.5 | 42.3 | 64.1 | 63.5 | 69.2 | 55.4 | 65.8 | 56.5 | 60.0 |
| **Ours** | **58.6** | **59.4** | **74.0** | **80.3** | **62.1** | 53.0 | **49.0** | **66.6** | **66.0** | **72.6** | **65.0** | **71.6** | **60.6** | **64.5** |

**Table 1:** Mean voxel IoU for 3D shape reconstruction with 4 input views.

**Figure 3: Novel view synthesis**. We show novel views synthesized for 4 example objects. The views are generated through a neural renderer from the aggregate representation. Best viewed in color with zoom in.

| Classes | aero | bench | cabinet | car | chair | display | lamp | speaker | rifle | sofa | table | phone | vessel | mean |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **3D-R2N2** [3] | 1.54 | 3.40 | 4.43 | 4.03 | 4.36 | 4.76 | 2.58 | 5.92 | 1.89 | 3.80 | 3.85 | 4.57 | 2.52 | 3.67 |
| **DnP** | 1.42 | 3.12 | 4.13 | 3.47 | 3.74 | 4.44 | 2.40 | 5.47 | 1.61 | 3.21 | 3.68 | 3.86 | 2.40 | 3.30 |
| **Ours (w/o sym)** | 1.44 | 3.30 | 3.77 | 3.43 | 3.86 | 4.40 | 2.63 | 5.44 | 1.65 | 3.15 | 3.94 | **3.79** | 2.44 | 3.33 |
| **Ours** | **1.41** | **2.87** | **3.61** | **3.19** | **3.63** | **4.25** | **2.28** | **5.27** | **1.58** | **2.98** | **3.13** | 4.06 | **2.29** | **3.12** |

**Table 2: Novel view synthesis.** Mean L1 error (scaled by 100) across classes when using 4 input images for inference.

and visualize the mean IoU in Fig 4a. Our approach out-performs the baselines over the spectrum, and performance consistently increases with additional views. Although all methods were similarly trained using 4 input views, we observe a more significant improvement over baselines when using smaller number of input views, indicating that the symmetries allow us to better leverage the information. We visualize in Figure 2 the predictions with varying number of input views.

**Novel-view Synthesis.** We evaluate the performance for the task of view synthesis using L1 error between predicted and ground-truth images (using 5 novel views per instance). We report the category-wise mean error for various approaches in the setting with 4 input views in Table 2, and highlight the mean error across classes with varying input views in Fig 4b. We also visualize sample results in the setting with 4 input views in Figure 3.

We notice a similar trend as in the case of volumetric prediction – our method improves over the baselines, and error reduces with additional views. In particular, the implicit aggregation method [3] has a spatially low resolution
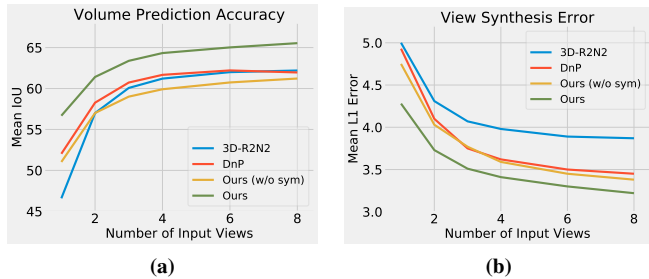
7

**Figure 4:** (a) Mean voxel IoU and (b) mean image L1 error vs number of input views.
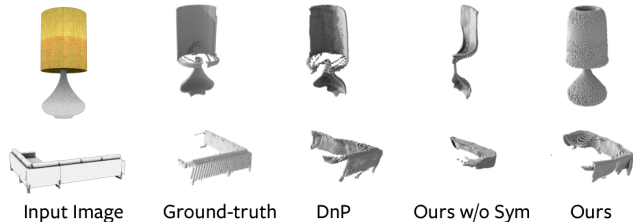


**Figure 5:** Predicted object-centric coordinates (visualized as point clouds) from a **single** image.
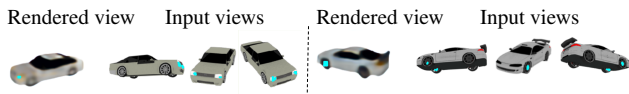


**Figure 6:** Back-propagating the gradient from a region (painted blue) in the rendered view (left) to the original images (right) reveals correspondence and symmetry.
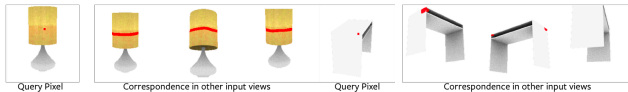


**Figure 7:** Given a pixel in an input view (left, highlighted in red), we visualize the pixels with the most similar symmetry-aware object coordinates (also highlighted in red).
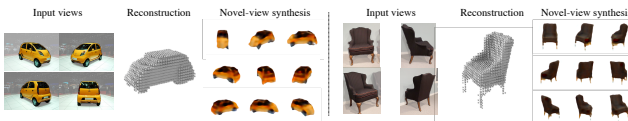


**Figure 8:** Multi-view aggregation on real images. Given 4 real images of an object in different poses (with object masks), our network is able to reconstruct the object shape and synthesize novel views.

aggregate feature, which prevents it from capturing appearance details well. While we note that all approaches produce slightly blurry results, and could be improved with additions *e.g.* adversarial losses, better renderer *etc.*, our goal here is to highlight the benefits our our aggregation method in comparison to alternates.

**Representation Analysis.** We inspect the structure learned by our model by back-tracing a pixel in a rendered image to its source set of input views. Specifically, we choose a region in the rendered image and compute the gradient of its sum with respect to the input images. We then highlight the source pixels with the most significant gradient magnitude. The result for several such regions is shown in Figure 6. We observe that the network relies on the information at the corresponding (or symmetric) pixels in the input images to render the target image *e.g.* the pixels on *both* front wheels in the source images are most influential for rendering the front *left* wheel in a novel views.

We also visualize object-centric coordinates obtained from a single image in Figure 5 (by selecting for each pixel the most likely predicted symmetry type), and compare our predictions to those obtained without symmetry inference, or via predicted per-pixel depth and global camera pose. We notice that our predictions are better aligned, and predict additional points for a symmetric object.

We additionally visualize in Figure 7 the corresponding pixels in other views that have similar symmetry-aware coordinate prediction to a given query pixel, and observe that the correspondences do respect the symmetry.

**Qualitative Results on Real Images.** There are several real-world scenarios that correspond to our inference setups, namely multiple images of an object without access to camera pose. These include for example an object on display on a turntable, images of products online *etc*. We show some qualitative results of our learned network on such data in Figure 8 using segmented images of a rotating car, and chair images from an online seller.

## 5. Discussion

We have presented an approach for aggregating multiple images of an object instance via predicting symmetry-aware object-centric coordinates, and have demonstrated that this aggregate representation can be leveraged for certain 3D tasks. While this has allowed us to improve over implicit, or camera-centric prediction based aggregation, our approach also has certain shortcomings. Classical SfM methods 'lift' pixels to 3D via reasoning across multiple images, as they rely on correspondence across images to do so. Our approach instead does this independently per image, and while the refinement could correct certain errors, the lifting itself could be improved via multi-image reasoning. Further, while the classical reconstruction methods are inherently 'unsupervised', our reliance on learning requires the use of supervisory data and it would be a desirable direction to lighten this burden. Lastly, our approach has tackled an object reconstruction setting using normalized coordinates, and it would be interesting to formulate extensions that could handle generic scenes.

# References

[1] Matthew Brown and David G Lowe. Unsupervised 3d object recognition and reconstruction in unordered datasets. In *3DIM*, 2005. 2

[2] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015. 1, 5

[3] Christopher B Choy, Danfei Xu, JunYoung Gwak, Kevin Chen, and Silvio Savarese. 3dr2n2: A unified approach for single and multi-view 3d object reconstruction. In *ECCV*, 2016. 1, 2, 3, 5, 6, 7

[4] Haoqiang Fan, Hao Su, and Leonidas J Guibas. A point set generation network for 3d object reconstruction from a single image. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 605–613, 2017. 1, 2

[5] Yasutaka Furukawa, Carlos Hernández, et al. Multi-view stereo: A tutorial. *Foundations and Trends® in Computer Graphics and Vision*, 9(1-2):1–148, 2015. 2

[6] Yasutaka Furukawa and Jean Ponce. Accurate, dense, and robust multiview stereopsis. *IEEE transactions on pattern analysis and machine intelligence*, 32(8):1362–1376, 2009. 2

[7] R. Girdhar, D.F. Fouhey, M. Rodriguez, and A. Gupta. Learning a predictable and generative vector representation for objects. 2016. 1, 2

[8] Thibault Groueix, Matthew Fisher, Vladimir G Kim, Bryan C Russell, and Mathieu Aubry. A papier-mâché approach to learning 3d surface generation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018. 2

[9] Klaus Häming and Gabriele Peters. The structure-from-motion reconstruction pipeline–a survey with focus on short image sequences. *Kybernetika*, 46(5):926–937, 2010. 2

[10] Christian Häne, Shubham Tulsiani, and Jitendra Malik. Hierarchical surface prediction for 3d object reconstruction. In *3DV*, 2017. 2

[11] Po-Han Huang, Kevin Matzen, Johannes Kopf, Narendra Ahuja, and Jia-Bin Huang. Deepmvs: Learning multi-view stereopsis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2821–2830, 2018. 2

[12] Zeng Huang, Yuanlu Xu, Christoph Lassner, Hao Li, and Tony Tung. Arch: Animatable reconstruction of clothed humans. *arXiv preprint arXiv:2004.04572*, 2020. 2

[13] Angjoo Kanazawa, Shubham Tulsiani, Alexei A Efros, and Jitendra Malik. Learning category-specific mesh reconstruction from image collections. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 371–386, 2018. 2

[14] Abhishek Kar, Christian Häne, and Jitendra Malik. Learning a multi-view stereo machine. In *NeurIPS*, 2017. 2, 5

[15] Jan J Koenderink and Andrea J Van Doorn. Affine structure from motion. *JOSA A*, 8(2):377–385, 1991. 1

[16] Yiyi Liao, Simon Donné, and Andreas Geiger. Deep marching cubes: Learning explicit surface representations. In *CVPR*, 2018. 2

[17] Chen-Hsuan Lin, Chen Kong, and Simon Lucey. Learning efficient point cloud generation for dense 3d object reconstruction. In *AAAI*, 2018. 2

[18] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. *arXiv preprint arXiv:1812.03828*, 2018. 2

[19] Kyle Olszewski, Sergey Tulyakov, Oliver Woodford, Hao Li, and Linjie Luo. Transformable bottleneck networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 7648–7657, 2019. 2

[20] Despoina Paschalidou, Osman Ulusoy, Carolin Schmitt, Luc Van Gool, and Andreas Geiger. Raynet: Learning volumetric 3d reconstruction with ray potentials. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018. 2

[21] Jhony K Pontes, Chen Kong, Sridha Sridharan, Simon Lucey, Anders Eriksson, and Clinton Fookes. Image2mesh: A learning framework for single image 3d reconstruction. *arXiv preprint arXiv:1711.10669*, 2017. 2

[22] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*, page 234–241, 2015. 4

[23] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2304–2314, 2019. 2

[24] Johannes L Schönberger, Enliang Zheng, Jan-Michael Frahm, and Marc Pollefeys. Pixelwise view selection for unstructured multi-view stereo. In *European Conference on Computer Vision*, pages 501–518. Springer, 2016. 2

[25] Steven M Seitz, Brian Curless, James Diebel, Daniel Scharstein, and Richard Szeliski. A comparison and evaluation of multi-view stereo reconstruction algorithms. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, volume 1, pages 519–528. IEEE, 2006. 2

[26] Vincent Sitzmann, Justus Thies, Felix Heide, Matthias Nießner, Gordon Wetzstein, and Michael Zollhöfer. Deepvoxels: Learning persistent 3d feature embeddings. In *CVPR*, 2019. 1, 5

[27] Noah Snavely, Steven M Seitz, and Richard Szeliski. Photo tourism: exploring photo collections in 3d. In *ACM transactions on graphics (TOG)*, 2006. 2

[28] Srinath Sridhar, Davis Rempe, Julien Valentin, Bouaziz Sofien, and Leonidas J Guibas. Multiview aggregation for learning category-specific shape reconstruction. In *Advances in Neural Information Processing Systems*, pages 2348–2359, 2019. 3

[29] Carlo Tomasi and Takeo Kanade. Shape and motion from image streams under orthography: a factorization method. *International Journal of Computer Vision*, 9(2):137–154, 1992. 1

[30] Shubham Tulsiani, Tinghui Zhou, Alexei A Efros, and Jitendra Malik. Multi-view supervision for single-view reconstruction via differentiable ray consistency. In *Proceedings of the*

*IEEE Conference on Computer Vision and Pattern Recognition*, pages 2626–2634, 2017. 1, 2

[31] Shimon Ullman. The interpretation of structure from motion. *Proceedings of the Royal Society of London B: Biological Sciences*, 1979. 2

[32] He Wang, Srinath Sridhar, Jingwei Huang, Julien Valentin, Shuran Song, and Leonidas J Guibas. Normalized object coordinate space for category-level 6d object pose and size estimation. In *CVPR*, 2019. 3

[33] Nanyang Wang, Yinda Zhang, Zhuwen Li, Yanwei Fu, Wei Liu, and Yu-Gang Jiang. Pixel2mesh: Generating 3d mesh models from single rgb images. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 52–67, 2018. 2

[34] Jiajun Wu, Yifan Wang, Tianfan Xue, Xingyuan Sun, William T Freeman, and Joshua B Tenenbaum. MarrNet: 3D Shape Reconstruction via 2.5D Sketches. 2017. 2

[35] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 767–783, 2018. 2

[36] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Ronen Basri, and Yaron Lipman. Multiview neural surface reconstruction with implicit lighting and material, 2020. 2