# Multi-view Consistency as Supervisory Signal
# for Learning Shape and Pose Prediction

Shubham Tulsiani, Alexei A. Efros, Jitendra Malik
University of California, Berkeley
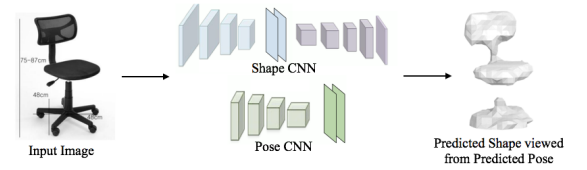{shubhtuls, efros, malik}@eecs.berkeley.edu

## Abstract

*We present a framework for learning single-view shape and pose prediction without using direct supervision for either. Our approach allows leveraging multi-view observations from unknown poses as supervisory signal during training. Our proposed training setup enforces geometric consistency between the independently predicted shape and pose from two views of the same instance. We consequently learn to predict shape in an emergent canonical (view-agnostic) frame along with a corresponding pose predictor. We show empirical and qualitative results using the ShapeNet dataset and observe encouragingly competitive performance to previous techniques which rely on stronger forms of supervision. We also demonstrate the applicability of our framework in a realistic setting which is beyond the scope of existing techniques: using a training dataset comprised of online product images where the underlying shape and pose are unknown.*

## 1. Introduction

Consider the flat, two-dimensional image of a chair in Figure 1(a). A human observer cannot help but perceive its 3D structure. Even though we may have never seen this particular chair before, we can readily infer, from this single image, its likely 3D shape and orientation. To make this inference, we must rely on our knowledge about the 3D structure of other, previously seen chairs. But how did we acquire this knowledge? And can we build computational systems that learn about 3D in a similar manner?

Humans are moving organisms: our ecological supervision [15] comprises of observing the world and the objects in it from different perspectives, and these multiple views inform us of the underlying geometry. This insight has been successfully leveraged by a long line of geometry-based reconstruction techniques. However these structure from motion or multi-view stereo methods work for specific instances and do not, unlike humans, generalize to predict the 3D shape of a novel instance given a single view. Some

---

Project website with code: https://shubhtuls.github.io/mvcSnP/



a) Testing: Our learned CNNs can infer the shape and pose from a **single input image**.
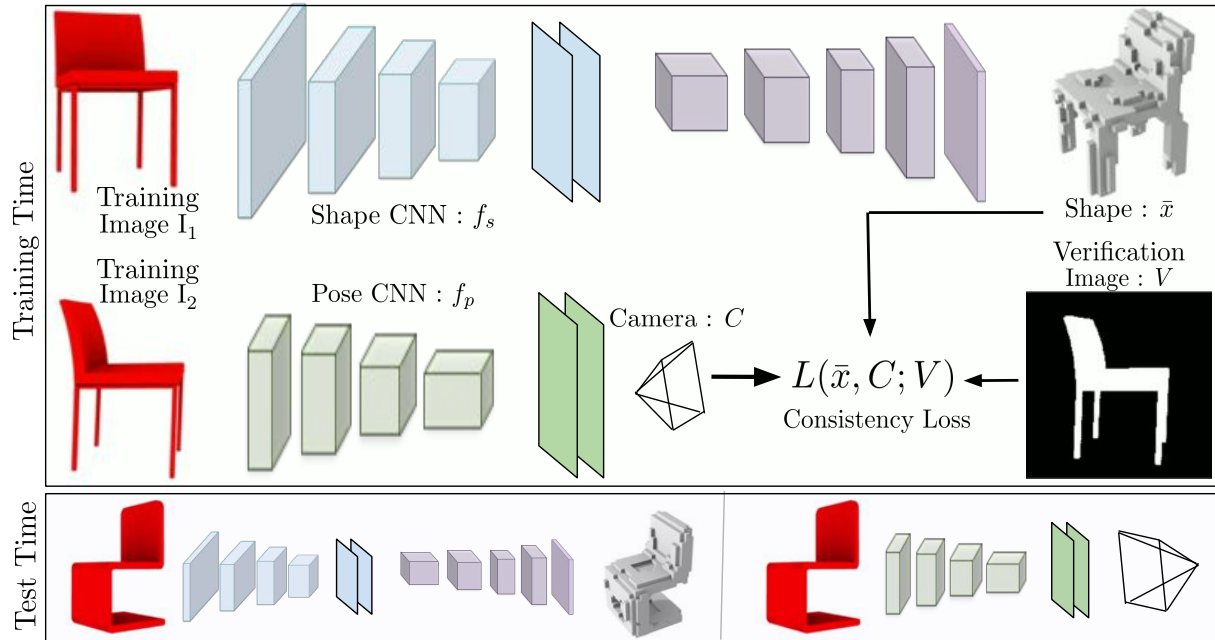


b) Training data: Multiple images of objects with **unknown shape** under **unknown camera pose** with associated (approximate) depth/mask

**Figure 1:** We learn to predict the shape and pose of an object from a single input view. Our framework can leverage training data of the form of multi-view observations of objects, and learn shape and pose prediction despite the lack of any direct supervision.

recent learning-based methods [8, 16] have attempted to address single-view 3D inference task, but this ability has come at a cost. These approaches rely on full 3D supervision and require known 3D shape for each training image. Not only is this form of supervision ecologically implausible, it is also practically tedious to acquire and difficult to scale. Instead, as depicted in Figure 1(b), our goal is to learn 3D prediction using the more naturally plausible multi-view supervision.

The broader goal of learning from data without explicit supervision is the focus of of considerable attention in the deep learning literature. Mechanisms that have been proposed include the use of information bottlenecks or proxy tasks such as prediction that encourage learning about the temporal or spatial structure. Similarly, in this paper, we rely on enforcing a geometric bottleneck for the task of explaining novel views and leverage the principle of multi-view consistency: a common geometry, observed from different perspectives can consistently explain multiple views of an instance. While some recent approaches [25, 30, 35] have utilized these principles to learn 3D shape prediction, they all crucially rely on object pose supervision during training.

**Figure 2:** Overview of our approach. During training, we use paired views of the same instance along with a depth/mask verification image from the second view. We predict shape from the first image and pose from the second, and enforce consistency between the shape, pose and the verification image. At test time, our learned models are used to infer the shape and pose from a single RGB input image.

Our proposed framework allows us to go a step further, and learn single-view shape and pose prediction using multi-view observations from *unknown* poses. Therefore, unlike previous methods which require either shape or pose supervision, we relax the requirement for *both* these forms of supervision.

Our approach, as summarized in Figure 2, learns shape and pose prediction by enforcing consistency between the predictions and available (novel view) observations. Concretely, given one image of an object instance, we predict a corresponding shape. In parallel, given a *different* image of the same instance, we independently predict a corresponding pose. Then, we enforce that the predicted shape (using the former image) should be 'consistent' with a depth/mask observation for the latter image when viewed from the predicted pose. As we discuss in Section 3, and demonstrate qualitatively and quantitatively demonstrate in Section 4, this allows us to learn single-view shape and pose prediction despite not having direct supervision for either.

## 2. Related Work

**Structure from Motion and Multi-view Instance Reconstruction.** Structure from motion (SfM) [31] based methods (e.g. [4, 28]) aim to recover the geometry, typically as sparse 3D point clouds, and the camera pose for each image. It was also shown that volumetric representations can be inferred by fusing multiple range images [9] or foreground masksl [3, 22, 24]. More closely related to our formulation, ray-potential based optimization methods [10, 23] can be used to infer discrete or probabilistic [32] volumetric

representations from multiple color images. This class of optimization techniques can be further extended to incorporate additional signals *e.g.* depth or semantics [21, 26, 27]. The goal of all these multi-view instance reconstruction methods is to infer the 3D structure of a specific scene/object given a large number of views of the *same instance*. Our method can be thought of as trying to minimize similar cost functions during training, but at test time, we can infer the pose and shape from a *single RGB image* – something that these classical techniques cannot do.

**Generative 3D Modeling without 3D Supervision.** Blanz and Vetter [2], using 3D supervision, captured the shapes of faces using a deformable model. Cashman and Fitzgibbon [5] subsequently demonstrated that similar generative models could be learned using only image based annotations. Kar *et al*. [19] extended these ideas to more general categories and automated test-time inference using off-the shelf recognition systems. However, these models are restricted to only capture deformations around a mean shape(s), thus limiting their expressiveness. Recently, Gadhela *et al*. [13] presented a more expressive generative model for shapes learned using a collection of silhouette images but did not examine applications for inference conditioned on image evidence. Eslami *et al*. [12] also learned a generative model with a corresponding inference module using only RGB images but only demonstrated 3D inference in scenarios where object shapes were known a priori. While the recent successes indicate that multi-view (or even single-view) ob-

servations can allow learning expressive generative models, their applications for single-view reconstruction have not been demonstrated conclusively. We instead propose to discriminatively train single-view shape and pose estimation systems using similar multi-view observations.

**Multi-view Supervision for Single-view Depth Prediction.** A recent direction pursued in the area of learning-based single-view depth prediction is to forego the need for direct supervision [11] and instead rely on multi-view observations for training [14, 17, 36]. Garg *et al*. [14] and Godard *et al*. [17] leverage stereo images as supervision to learn single image depth prediction. Zhou *et al*. [36] further relax the assumption of known relative pose between the multiple views, and learn single-view depth and ego-motion prediction models from monocular videos. Similarly, we leverage multiple views from unknown poses as supervisory signal but we pursue 3D instead of 2.5D predictions.

**Multi-view Supervised Single-view Reconstruction.** Initial CNN-based methods [8, 16, 34] predicted voxel occupancy representations from a single input image but required full 3D supervision during training. Recent approaches have advocated using alternate forms of supervision. Zhu *et al*. [37] showed that systems trained using synthetic shape and pose supervision could be adapted to real data using only image based annotation. Their pre-training, however, crucially relied on direct shape and pose supervision. Towards relaxing the need of any shape supervision, some recent methods demonstrated the feasibility of using multi-view foreground masks [18, 25, 35] or more general forms of observation *e.g*. depth, color, masks, semantics *etc*. [30] as supervisory signal. Our work adheres to this ideology of using more natural forms of supervision for learning 3D prediction and we take a step further in this direction. The previous multi-view supervised approaches [18, 25, 30, 35] required known camera poses for the multiple views used during training and our work relaxes this requirement.

## 3. Approach

We aim to learn shape and pose prediction systems, denoted as $f_s$ and $f_p$ respectively, which can infer the corresponding property for the underlying object from a single image. However, instead of direct supervision, the supervision available is of the form of multi-view observations from unknown poses. We first formally define our problem setup by describing the representations inferred and training data leveraged and then discuss our approach.

**Training Data.** We require a sparse set of multi-view observations for multiple instances of the same object category. Formally, denoting by $\mathcal{N}(i)$ the set of natural numbers up to $i$, we assume a dataset of the form $\{\{(I_v^i, V_v^i) \mid v \in \mathcal{N}(N_i)\} \mid i \in \mathcal{N}(N)\}$. This corresponds to $N$ object instances, with $N_i$ views available for the $i^{th}$ instance. Associated with each image $I_i^v$, there is also a depth/mask image

$V_i^v$ that is used for consistency verification during training. Note that there is no direct pose or shape supervision used – only multi-view observations with identity supervision.

**Shape and Pose Parametrization.** The (predicted) shape representation $\bar{x}$ is parametrized as occupancy probabilities of cells in a 3D grid. The pose of the object, parametrized as a translation $t$ and rotation $R$, corresponds to the camera extrinsic matrix. While we assume known camera intrinsics for our experiments, our framework can also be extended to predict these.

### 3.1. Geometric Consistency as Supervision

Multiple images of the same instance are simply renderings of a common geometry from diverse viewpoints. Therefore, to correctly 'explain' multiple observations of an instance, we need the correct geometry (shape) of the instance and the corresponding viewpoints (pose) for each image. Our approach, which is depicted in Figure 2, builds on this insight and proposes to predict *both*, shape and pose s.t. the available multi-view observations can be explained.

Concretely, during training, we use one image of an instance to predict the instance shape. In parallel, we use a *different* image of the same instance to predict pose. Then, we enforce that the predicted shape, when viewed according to the predicted pose, should be consistent with a depth/mask image from the latter view. We therefore use the notion of *consistency* as a form of meta-supervision *i.e*. while the ground-truth shape and pose are unknown, we know that they should be consistent with the available verification image. After the training stage, our learned models can infer shape and pose from a single view of a novel instance.

A crucial aspect of the designed training setup is that the shape and pose estimates are *independently* obtained from *different* images of the same instance. This enforces that the optimal solution corresponds to predicting the correct shape and pose. Another interesting property is that the shape is predicted in an emergent canonical, view-independent frame, and the predicted pose is with respect to this frame.

**Correctness of Optimal Shape and Pose.** We consider Figure 2 and first examine the shape prediction CNN $f_s$. It predicts a shape $f_s(I_1)$ given some input image. This shape is verified against $V$ from a different view which is unknown to $f_s$. The optimal predicted shape should therefore be consistent with *all* possible novel views of this instance, and therefore correspond to the true shape (upto some inherent ambiguities *e.g*. concavities in case of mask supervision). Similarly, the pose prediction CNN $f_p$ is required to infer a viewpoint under which the predicted geometry can explain the verification image $V$. As $V$ is chosen to be from the same viewpoint as the image $I_2$, the pose CNN should predict the correct viewpoint corresponding to its input image ($I_2$).

**Emergent Canonical Frame.** Under our proposed setup, the predicted pose $f_p(I_2)$ is agnostic to the image $I_1$. How-

ever, to explain the verification image $V$, the pose CNN is required to predict a pose w.r.t the inferred shape $f_s(I_1)$. So how can $f_p$ infer pose w.r.t $f_s(I_1)$ when it does not even have access to $I_1$? The resolution to this is that the shape prediction CNN $f_s$ automatically learns to predict shape in some (arbitrary) view-agnostic canonical frame (e.g. 'front' of chairs may always face towards the X axis), and the pose CNN $f_p$ learns to predict pose w.r.t this frame. Therefore, even though it is not explicitly enforced, our approach of independently inferring shape and pose makes the learnt CNNs automatically adhere to some emergent canonical frame.

Towards implementing our framework, we require a consistency loss $L(\bar{x}, C; V)$ which measures whether the (predicted) shape $\bar{x}$ and camera pose $C$ can geometrically explain a depth/mask image $V$. We present a formulation for this loss in Section 3.2 and then describe the training process in Section 3.3. We finally describe some modifications required to make the training more robust.

### 3.2. Pose-differentiable Consistency Loss

We formulate a view consistency loss $L(\bar{x}, C; V)$ that measures the inconsistency between a shape $\bar{x}$ viewed according to camera $C$ and a depth/mask image $V$. Our formulation builds upon previously proposed differentiable ray consistency formulation [30]. However, unlike the previous formulation, our proposed view consistency loss is differentiable w.r.t pose (a crucial requirement for usage in our learning framework). Here, we very briefly recall the previous formulation and mainly highlight our proposed extension. A more detailed and complete formulation of the view consistency loss can be found in the appendix.

**Differentiable Ray Consistency [30].** The view consistency loss formulated by Tulsiani *et al*. [30] could be decomposed into per-pixel (or ray) based loss terms where $L_p(\bar{x}, C; v_p)$ denotes the consistency of the shape and camera with the observation $v_p$ at pixel $p$. The per-pixel loss is defined as the *expected event cost*:

$$L_p(\bar{x}, C; v_p) \ = \ \sum_{i=1}^{N} q_p(i)\psi_p(i) \tag{1}$$

Here, $\psi_p(i)$ denotes the cost for each event, determined by $v_p$, and $q_p(i)$ indicates the *event probability i.e.* the likelihood of the ray stopping at the $i^{th}$ voxel in its path. The event probability, $q_p(i)$ is in turn instantiated using the probabilities $\{x_p^i\}$ - where $x_p^i$ denotes the occupancy probability of the $i^{th}$ voxel in the ray's path. See appendix for details.

**Sampling Occupancies along a Ray.** The loss function as defined above is differentiable w.r.t shape $\bar{x}$, but not the camera parameters. This is because the quantity $\{x_p^i\}$ is not a differentiable function of the camera (since the ordering of voxels on a ray's path is a discrete function). Our insight is that instead of looking up *voxels* on the ray's path, we can consider *samples* along its path. Thus, our formulation

is similar to that proposed by Tulsiani *et al*. [30], with the difference that the variable $\{x_p^i\}$ is redefined to correspond to the occupancy at the $i^{th}$ point sample along the ray.

Concretely, we sample points at a fixed set of $N = 80$ depth values $\{d_i | 1 \le i \le N\}$ along each ray. To determine $x_i^p$, we look at the 3D coordinate of the corresponding point (determined using camera parameters), and trilinearly sample the shape $\bar{x}$ to determine the occupancy at this point.

$$l_i \equiv (\frac{u - u_0}{f_u}d_i, \frac{v - v_0}{f_v}d_i, d_i) \tag{2}$$

$$x_i^p = \mathcal{T}(\bar{x}, R \times (l_i + t)) \tag{3}$$

As the trilinear sampling function $\mathcal{T}$ is differentiable w.r.t its arguments, the sampled occupancy $x_i^p$ is differentiable w.r.t the shape $\bar{x}$ and the camera $C$. We note that Yan *et al*. [35] also used a similar sampling trick but their formulation is restricted to specifically using mask verification images and is additionally not leveraged for learning about pose.

### 3.3. Learning

**Training Objective.** To train the shape and pose predictors, we leverage the view consistency loss previously defined (Section 3.2) and train $f_s, f_p$ jointly to minimize $L_{data} = \sum_{i=1}^{N} \sum_{u=1}^{N_i} \sum_{v=1}^{N_i} L(f_s(I_u^i), f_p(I_v^i); V_v^i)$. Therefore, the shape predicted using every image $f_s(I_u^i)$ should be consistent with *all* available verification images of the same instance ($\{V_v^i\}$) when viewed from the corresponding (predicted) poses ($\{f_p(I_v^i)\}$). As detailed earlier, the independent prediction of shape and pose from different images ensures that the CNNs learn to infer the correct shape and pose under some emergent canonical frame.

**Architecture and Optimization Details.** We use a minibatch size of 8 images $I_u^i$ for which shape is predicted. For each of these images, we randomly sample at least 2, and upto 3 if available, out of $N_i$, views $I_v^i$ of the same instance *i.e.* the mini-batch size for the pose prediction CNN is between 16 and 24. We use extremely simple CNN architectures (depicted in Figure 2) corresponding to $f_s$ and $f_p$. Note that both these CNNs are initialized randomly (without any pre-training) and trained using ADAM [20].

*Shape Prediction.* Our shape prediction CNN has an encoder-decoder structure similar to the one used by Tulsiani *et al*. [30]. The input to the CNN is an RGB image of size $64 \times 64$ and the outputs are corresponding voxel occupancy probabilities for a $32 \times 32 \times 32$ grid.

*Pose Prediction.* Our pose prediction CNN $f_p$ has a similar encoder to $f_s$, but outputs the predicted pose via fully connected layers. The rotation aspect of the pose is parametrized using two euler angles (azimuth, elevation) and the predicted translation $\in \mathbb{R}^3$. However, for some analysis experiments, we also assume that the object is at a known location w.r.t

the camera and only predict the camera rotation. While in this work we assume known intrinsic parameters, the pose prediction CNN could in principle be extended to infer these.

### 3.4. Overcoming Local Minima

We observed that our training is susceptible to local minima, in particular for the pose prediction CNN $f_p$. This is not too surprising since we have to learn both shape and pose from scratch, and erroneous estimates for one could confound the learning for the other, particularly in the in the initial stages We observe that the $f_p$ learns to predict only a small range of poses and *e.g.* instead of predicting back-facing chairs, it confuses them with front-facing chairs. To avoid such local minima, we introduce two changes to the setup previously described.

**Incorporating a Pose Prior.** We encourage the distribution of the predicted poses to be similar to a prior distribution (uniform azimuth $\in [0, 360)$, elevation $\in [-20, 40)$ degrees). We do so by adding an adversarial loss for the predictions of $f_p$ where the 'real' samples are drawn from the prior distribution and 'generated' samples are those predicted by $f_p$. We empirically show that our training is robust to the exact prior and that it can be different from the true distribution.

**Allowing Diverse Predictions.** While the adversarial loss encourages diverse predictions, we also need some architectural changes to easily capture these. Instead of directly regressing to a single pose estimate in the last layer, we predict $N_p = 8$ estimates and additionally predict a probability distribution over these. We then sample a pose according to the predicted distribution. We use Reinforce [33] to obtain gradients for the probability predictions.

## 4. Experiments

We consider two different scenarios where we can learn single-view shape and pose prediction using multi-view observations from unknown poses. We first examine the ShapeNet dataset where we can synthetically generate images and compare our approach against previous techniques which rely on stronger forms of supervision. We then consider a realistic setting where the existing approaches, all of which require either shape or pose supervision, cannot be applied due to lack of any such annotation. Unlike these existing methods, we show that our approach can learn using an online product dataset where multiple images on objects are collected from product websites *e.g.* eBay.

### 4.1. Empirical Analysis using ShapeNet

#### 4.1.1 Experimental Setup

**Dataset.** We use the ShapeNet dataset [6] to empirically validate our approach. We evaluate on three representative object categories with a large number of models : airplanes, cars, and chairs. We create random train/val/test splits with $(0.7, 0.1, 0.2)$ fraction of the models respectively. For each

training model, we use $N_i = 5$ images available from different (unknown) views with corresponding depth/mask observations. The images are rendered using blender and correspond to a viewpoint from a randomly chosen azimuth $\in [0, 360)$ degrees and elevation $\in [-20, 40]$ degrees. We additionally use random lighting variations during rendering.

We also render the training objects under two settings - a) origin centred, or b) randomly translated around the origin. As the camera is always at a fixed distance away from the origin, the first setting corresponds to training with a known camera translation, but unknown rotation. The second corresponds to training with both translation and rotation unknown. To have a common test set across various control setting (and compare to [30]), we use the origin centered renderings for our validation and test sets. We note that these rendering settings are rather challenging and correspond to significantly more variation than commonly examined by previous multi-view supervised methods which examine settings with fixed translation [30], and sometimes only consider 24 [35] or even 8 [13] possible discrete views.
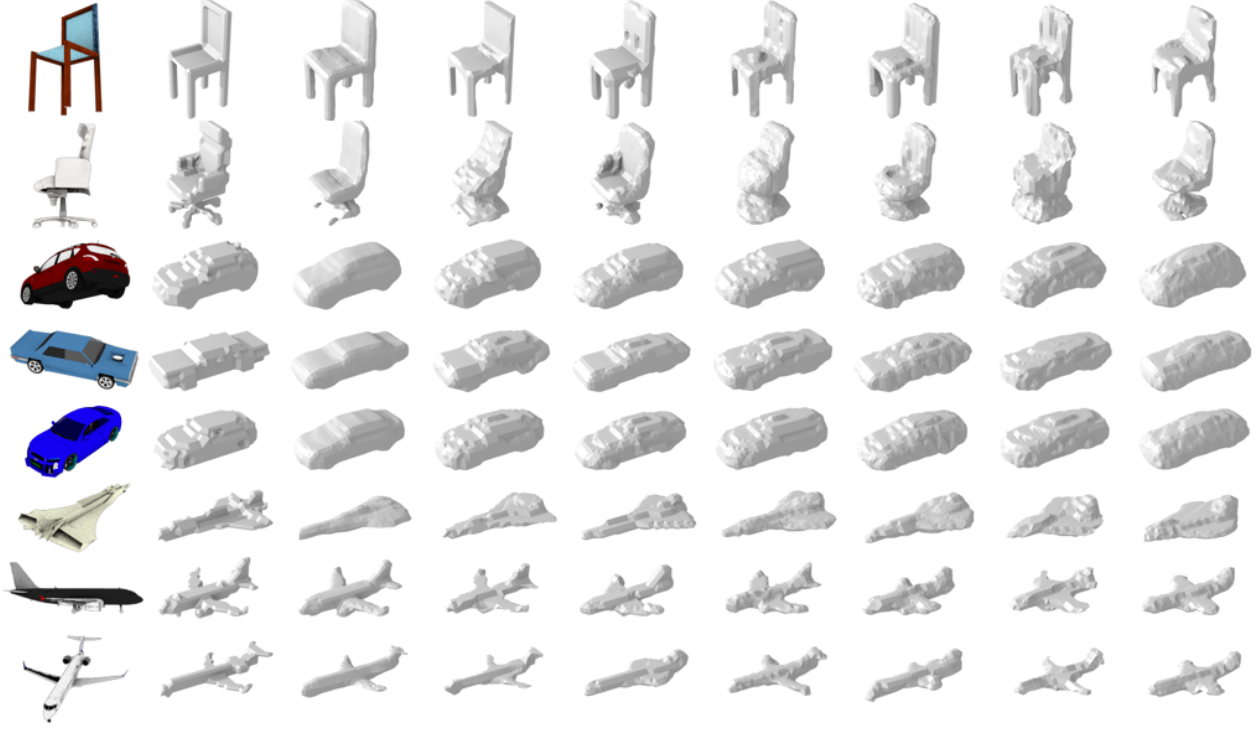
**Control Settings.** In addition to reporting the performance in the scenario where pose and shape supervision is unavailable, we also examine the settings where stronger supervision *e.g.* shape or pose can be used. These experiments serve to highlight the upper bound performance. In all the experiments, we train a separate model per object category. The various settings studied are :

*3D Supervision.* To mimic the setup used by 3D supervised approaches [8, 16], we assume known ground-truth 3D models for each training image and train the shape CNN using a cross-entropy loss.
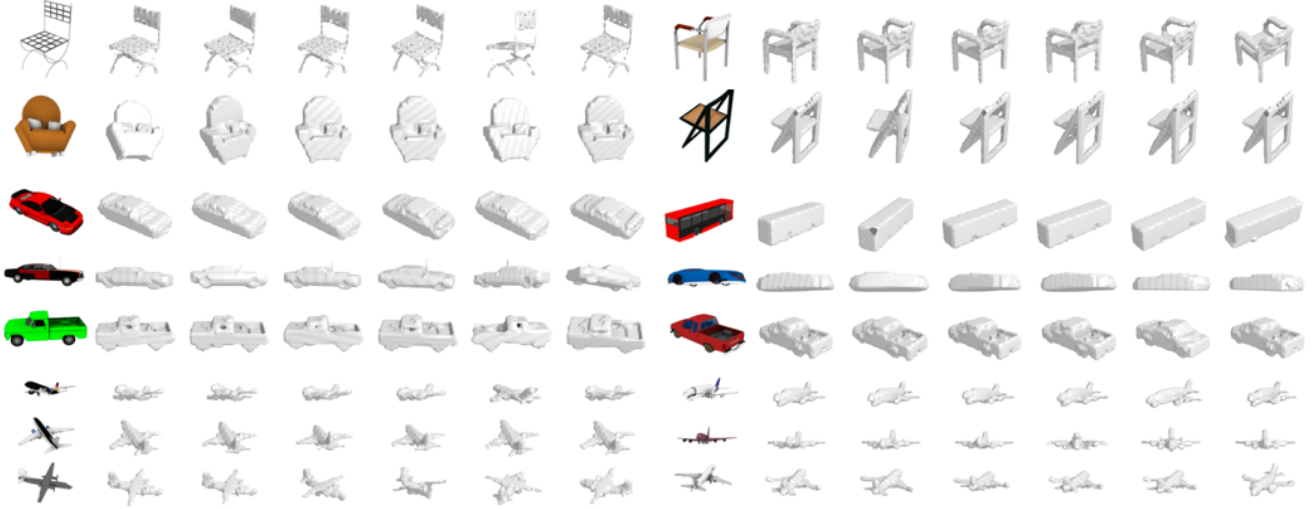
*Multi-view with Ground-truth Pose.* In this supervision setting used by previous multi-view supervised approaches, pose (but not shape) supervision is available for the multiple observations. We use our loss function but train the shape prediction CNN $f_s$ using the ground-truth pose instead of predicted poses. We separately train the pose prediction CNN $f_p$ using squared L2 loss in quaternion space (after accounting for antipodal symmetry of quaternions).

*Multi-view without Pose Supervision.* This represents our target setting with the weakest form of supervision available. We train the shape and pose prediction CNNs jointly using our proposed loss. Further, we consider two variants of this setting - one where camera translation is known, one where both camera translation and rotation are unknown.
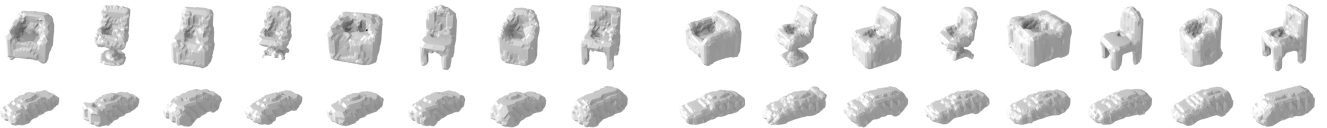
**Evaluation Metrics.** We report the results using predictions for 2 images per test model. For evaluating the shape prediction CNN, we report the mean intersection over union (IoU) between the ground-truth and predicted shapes. Since different CNNs can be calibrated differently, we search for the optimal threshold (per CNN on the validation set) to binarize the predictions. To evaluate the rotation prediction, we measure the angular distance between the predicted and

**Figure 3:** Shape predictions on the validation set using a single RGB input image. We visualize the voxel occupancies by rendering the corresponding mesh (obtained via marching cubes) from a canonical pose. Left to Right: a) Input Image b) Ground-truth c) 3D Supervised Prediction d,e) Multi-view & Pose Supervision (Mask, Depth) f,g) Mult-view w/o Rotation Supervision (Mask, Depth), and h,i) Mult-view w/o Rotation and Translation Supervision (Mask, Depth)



**Figure 4:** Rotation predictions on a random subset of the validation images. For visualization, we render the ground-truth voxel occupancies using the corresponding rotation. Left to Right: a) Input Image b) Ground-truth Rotation c) GT Supervised Prediction d,e) Multi-view w/o Rot Supervision (Mask, Depth), and f,g) Multi-view w/o Rot and Trans Supervision (Mask, Depth)
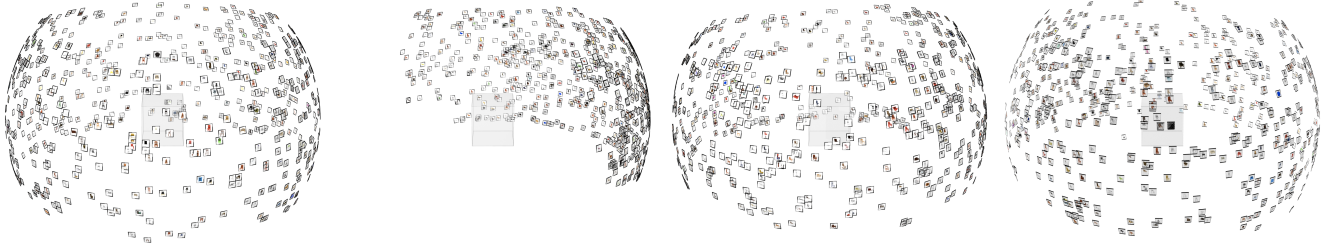


**Figure 5:** Visualization of 8 random predicted shapes from azimuth $= 60^\circ$, elevation $= 30^\circ$. Left: Original predictions from the shape CNN. Right: Shape predictions transformed according to the optimal rotation.

6

| Training Data | 3D | Multi-view & GT Pose | | Multi-view w/o Rot | | Multi-view w/o Rot & Trans | | Training Data | GT Pose | | MV w/o Rot | | | | MV w/o Rot & Trans | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | Mask | | Depth | | Mask | | Depth | |
| class | | Mask | Depth | Mask | Depth | Mask | Depth | class | Acc | Err | Acc | Err | Acc | Err | Acc | Err | Acc | Err |
| aero | 0.57 | 0.55 | 0.43 | 0.52 | 0.44 | 0.38 | 0.37 | aero | 0.79 | 10.7 | 0.69 | 14.3 | 0.60 | 21.7 | 0.53 | 26.9 | 0.63 | 12.3 |
| car | 0.79 | 0.75 | 0.69 | 0.74 | 0.71 | 0.48 | 0.68 | car | 0.90 | 7.4 | 0.87 | 5.2 | 0.85 | 4.9 | 0.53 | 24.8 | 0.56 | 20.6 |
| chair | 0.49 | 0.42 | 0.45 | 0.40 | 0.43 | 0.35 | 0.37 | chair | 0.85 | 11.2 | 0.81 | 7.8 | 0.83 | 8.6 | 0.55 | 24.0 | 0.62 | 19.1 |
| mean | 0.62 | 0.57 | 0.52 | 0.55 | 0.53 | 0.40 | 0.47 | mean | 0.85 | 10.0 | 0.79 | 9.0 | 0.76 | 11.7 | 0.54 | 25.1 | 0.61 | 17.4 |

**Table 1:** Analysis of the performance for single-view shape (Left) and pose (Right) prediction. a) Shape Accuracy: Mean IoU on the test set using various supervision settings. b) Pose Accuracy/Error: Acc $\frac{\pi}{6}$ and Med-Err across different supervision settings.



**Figure 6:** Visualization of the predicted pose distribution under various training settings. Each small image is placed at the (predicted/known) location of the corresponding camera. The reference grid in the centre depicts the space in which shape is predicted. Left to Right : a) Ground-truth poses b) No pose prior c) True pose prior d) Incorrect pose prior, discarded midway through training. See text for details.

ground-truth rotation (in degrees) and report two metrics : a) Fraction of instances with error less than 30 degrees (Acc $\frac{\pi}{6}$), and b) Median Angular Error (Med-Err).

#### 4.1.2 Results

**Prediction Frame Alignment.** The ShapeNet models are all aligned in a canonical frame where X and Y axes represent lateral and upward directions. The shape and pose prediction CNNs learned using our approach are not constrained to adhere to this frame and in practice, learn to predict shape and pose w.r.t some arbitrary frame.

However, to evaluate these predictions, we compute an optimal rotation to best align the predictions to the canonical ShapeNet frame. We use 8 *random* images per category (the first validation mini-batch) alongwith the ground-truth 3D voxelizations and search for a rotation that maximizes the voxel overlap between the ground-truth and the rotated predicted shapes. We visualize the prediction frame alignment for car and chair CNNs trained using multi-view observations w/o pose via depth verification images in Figure 5. Note that the prediction frames across classes vary arbitrarily. After the alignment process, the predictions for both categories are in the canonical ShapeNet frame.

**Role of a Pose prior.** While the empirical results reported below correspond to using the correct pose prior, we first show that the primary benefit of this prior is that it encourages the CNN to predict diverse poses and avoid local minima, and that even an approximate prior is sufficient.

To further support this point, we conducted an experiment where we used an incorrect pose prior (elevation uniform $\in [-40, 80]$ instead of $\in [-20, 40]$) and removed the prior loss midway through training. We observed that this network also trained successfully, indicating that we do not require the true pose prior, rather only an approximate one. Figure 6 visualizes the pose distributions inferred under various settings. While using no prior results in a local optima, using the approximate prior (or the correct prior) does not.

**Single-view Shape Prediction.** Our results and the performance under various control settings with stronger supervision is reported in Table 1 and visualized in Figure 3. In general, we observe that the results using our approach are encouragingly close to those obtained using much stronger forms of supervision. This clearly indicates that our approach is able to learn single-view shape prediction despite the lack of either shape or pose information during training. As expected, we also observe that we cannot learn about concavities in chairs via consistency against mask validation images, though we can do so using depth images. e observe a noticeable performance drop in case of mask supervision with unknown translation, as this settings results in scale ambiguities which our evaluation does not account for *e.g.* we learn to predict larger cars, but further away, and this results in a low empirical score.

**Single-view Pose Estimation.** The results of our approach are reported in Table 1 and visualized in Figure 4. We observe a similar trend for the task of pose prediction – that

**Figure 7:** Visualization of predictions using the Stanford Online Product Dataset. (Top) Input image. (Middle) Predicted shape in the emergent canonical pose. (Bottom) Predicted shape rotated according to the predicted pose.

our approach performs comparably to directly supervised learning using ground-truth pose supervision. Interestingly, we often get lower median errors than the supervised setting. We attribute this to the different topologies of the loss functions. The squared L2 loss used in the supervised setting yields small gradients if the pose is almost correct. Our consistency loss however, would want the observation image to perfectly align with the shape via the predicted pose.

**Interpretation.** The main takeaway from these results is that it is indeed possible to learn shape and pose prediction without direct supervision for either. We empirically and qualitatively observe competitive performances for both these tasks when compared to approaches that leverage stronger forms of supervision. We see that we always learn meaningful shape and pose prediction systems across observation types (mask/depth) and that performance degrades gracefully when using less supervision (known/unknown translation).

## 4.2. Learning from Online Product Images

**Dataset.** We examined the 'chair' object category from the Stanford Online Products Dataset [29] which comprises of automatically downloaded images from eBay.com [1]. Since multiple images (views) of the same product are available, we can leverage our approach to learn from this data. As we also require associated foreground masks for these images, we use an out-of-the-box semantic segmentation system [7] to obtain these. However, the obtained segmentation masks are often incorrect. Additionally, many of the product images were not suited for our setting as they only comprised of a zoom-in of a small portion of the instance (*e.g.* chair wheel). We therefore manually selected images of unoccluded/untruncated instances with a reasonably accurate (though still noisy) predicted segmentation. We then used the object instances with atleast 2 valid views for training. This results in a filtered dataset of $N = 282$ instances with $N_i = 3.65$ views on average per instance.

**Results.** We can apply our approach to learn from this dataset comprising of multiple views with associated (approximate) foreground masks. Since the camera intrinsics are unknown, we assume a default intrinsic matrix (see appendix). We then learn to predict the (unknown) translation and rotation via $f_p$ and the (unknown) shape via $f_s$ using the available multi-view supervision. Note that the learned CNNs are trained from scratch, and that we use the same architecture/hyperparameters as in the ShapeNet experiments.

Some results (on images of novel instances) using our learned CNN are visualized in Figure 7. We see that we can learn to predict meaningful 3D structure and infer the appropriate shape and pose corresponding to the input image. Since only foreground mask supervision is leveraged, we cannot learn to infer the concavities in shapes. We also observe confusion across poses which result in similar foreground masks. However, we feel that this result using training data derived from a challenging real world setting, concretely demonstrates our method's ability to learn despite the lack of direct shape or pose supervision. To the best of our knowledge, this is the first such result and it represents an encouraging step forward.

## 5. Discussion

We presented a framework that allows learning single-view prediction of 3D structure without direct supervision for shape or pose. While this is an encouraging result that indicates the feasibility of using natural forms of supervision for this task, a number of challenges remain to be addressed. As our supervisory signal, we rely on consistency with validation images of unoccluded objects and it would be useful to deal with unknown occlusions. It would also be interesting to apply similar ideas for learning the 3D structure of general scenes though this might additionally require leveraging alternate 3D representations and allowing for object motion to handle dynamic scenes.

8

# References

[1] https://www.ebay.com/. 8

[2] V. Blanz and T. Vetter. A morphable model for the synthesis of 3d faces. In *SIGGRAPH*, 1999. 2

[3] A. Broadhurst, T. W. Drummond, and R. Cipolla. A probabilistic framework for space carving. In *ICCV*, 2001. 2

[4] M. Brown and D. G. Lowe. Unsupervised 3d object recognition and reconstruction in unordered datasets. In *3DIM*, 2005. 2

[5] T. J. Cashman and A. W. Fitzgibbon. What shape are dolphins? building 3d morphable models from 2d images. *TPAMI*, 2013. 2

[6] A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song, H. Su, J. Xiao, L. Yi, and F. Yu. ShapeNet: An Information-Rich 3D Model Repository. Technical Report arXiv:1512.03012 [cs.GR], 2015. 5

[7] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *TPAMI*, 2017. 8, ii

[8] C. B. Choy, D. Xu, J. Gwak, K. Chen, and S. Savarese. 3d-r2n2: A unified approach for single and multi-view 3d object reconstruction. In *ECCV*, 2016. 1, 3, 5

[9] B. Curless and M. Levoy. A volumetric method for building complex models from range images. In *SIGGRAPH*, 1996. 2

[10] J. De Bonet and P. Viola. Roxels: Responsibility weighted 3d volume reconstruction. In *ICCV*, 1999. 2

[11] D. Eigen and R. Fergus. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In *ICCV*, 2015. 3

[12] S. A. Eslami, N. Heess, T. Weber, Y. Tassa, D. Szepesvari, G. E. Hinton, et al. Attend, infer, repeat: Fast scene understanding with generative models. In *NIPS*, 2016. 2

[13] M. Gadelha, S. Maji, and R. Wang. Unsupervised 3d shape induction from 2d views of multiple objects. In *3DV*, 2017. 2, 5

[14] R. Garg and I. Reid. Unsupervised cnn for single view depth estimation: Geometry to the rescue. In *ECCV*, 2016. 3

[15] J. J. Gibson. The ecological approach to visual perception. 1979. 1

[16] R. Girdhar, D. Fouhey, M. Rodriguez, and A. Gupta. Learning a predictable and generative vector representation for objects. In *ECCV*, 2016. 1, 3, 5

[17] C. Godard, O. Mac Aodha, and G. J. Brostow. Unsupervised monocular depth estimation with left-right consistency. In *CVPR*, 2017. 3

[18] J. Gwak, C. B. Choy, A. Garg, M. Chandraker, and S. Savarese. Weakly supervised 3d reconstruction with adversarial constraint. In *3DV*, 2017. 3

[19] A. Kar, S. Tulsiani, J. Carreira, and J. Malik. Category-specific object reconstruction from a single image. In *CVPR*, 2015. 2

[20] D. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980.* 4

[21] A. Kundu, Y. Li, F. Dellaert, F. Li, and J. M. Rehg. Joint semantic segmentation and 3d reconstruction from monocular video. In *ECCV*, 2014. 2

[22] A. Laurentini. The visual hull concept for silhouette-based image understanding. *TPAMI*, 1994. 2

[23] S. Liu and D. B. Cooper. Ray markov random fields for image-based 3d modeling: model and efficient inference. In *CVPR*, 2010. 2

[24] W. Matusik, C. Buehler, R. Raskar, S. J. Gortler, and L. McMillan. Image-based visual hulls. In *SIGGRAPH*, 2000. 2

[25] D. J. Rezende, S. A. Eslami, S. Mohamed, P. Battaglia, M. Jaderberg, and N. Heess. Unsupervised learning of 3d structure from images. In *NIPS*, 2016. 1, 3

[26] N. Savinov, C. Hane, L. Ladicky, and M. Pollefeys. Semantic 3d reconstruction with continuous regularization and ray potentials using a visibility consistency constraint. In *CVPR*, 2016. 2

[27] N. Savinov, C. Häne, M. Pollefeys, et al. Discrete optimization of ray potentials for semantic 3d reconstruction. In *CVPR*, 2015. 2

[28] N. Snavely, S. M. Seitz, and R. Szeliski. Photo tourism: exploring photo collections in 3d. In *ACM transactions on graphics (TOG)*, 2006. 2

[29] H. O. Song, Y. Xiang, S. Jegelka, and S. Savarese. Deep metric learning via lifted structured feature embedding. In *CVPR*, 2016. 8, ii

[30] S. Tulsiani, T. Zhou, A. A. Efros, and J. Malik. Multi-view supervision for single-view reconstruction via differentiable ray consistency. In *CVPR*, 2017. 1, 3, 4, 5, i, ii

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

[32] A. O. Ulusoy, A. Geiger, and M. J. Black. Towards probabilistic volumetric reconstruction using ray potentials. In *3DV*, 2015. 2

[33] R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 1992. 5

[34] J. Wu, Y. Wang, T. Xue, X. Sun, W. T. Freeman, and J. B. Tenenbaum. MarrNet: 3D Shape Reconstruction via 2.5D Sketches. In *NIPS*, 2017. 3

[35] X. Yan, J. Yang, E. Yumer, Y. Guo, and H. Lee. Perspective transformer nets: Learning single-view 3d object reconstruction without 3d supervision. In *NIPS*, 2016. 1, 3, 4, 5, ii

[36] T. Zhou, M. Brown, N. Snavely, and D. Lowe. Unsupervised learning of depth and ego-motion from video. In *CVPR*, 2017. 3

[37] R. Zhu, H. Kiani, C. Wang, and S. Lucey. Rethinking reprojection: Closing the loop for pose-aware shape reconstruction from a single image. In *ICCV*, 2017. 3

# Appendix : Multi-view Consistency as Supervisory Signal
# for Learning Shape and Pose Prediction

Shubham Tulsiani, Alexei A. Efros, Jitendra Malik

University of California, Berkeley

{shubhtuls,efros,malik}@eecs.berkeley.edu

## A1. Loss Formulation

We briefly described, in the main text, the formulation of a view consistency loss $L(\bar{x}, C; V)$ that measures the inconsistency between a shape $\bar{x}$ viewed according to camera $C$ and a depth/mask image $V$. Crucially, this loss was differentiable w.r.t both, pose and shape. As indicated in the main text, our formulation builds upon previously proposed differentiable ray consistency formulation [30] with some innovations to make it differentiable w.r.t pose. For presentation clarity, we first present our full formulation, and later discuss its relation to the previous techniques (a similar discussion can also be found in the main text).

**Notation.** The (predicted) shape representation $\bar{x}$ is parametrized as occupancy probabilities of cells in a 3D grid. We use the convention that a particular value in the tensor $x$ corresponds to the probability of the corresponding voxel being *empty*. The verification image $V$ that we consider can be a depth or foreground mask image. Finally, the camera $C$ is parametrized via the intrinsic matrix $K$, and extrinsic matrix defined using a translation $t$ and rotation $R$.

**Per-pixel Error as Ray Consistency Cost.** We consider the verification image $V$ one pixel at a time and define the per-pixel error using a (differentiable) ray consistency cost. Each pixel $p \equiv (u, v)$ has an associated value $v_p$ *e.g.* in the case of a depth image, $v_p$ is the recorded depth at the pixel $p$. Additionally, each pixel corresponds to a ray originating from the camera centre and crossing the image plane at $(u, v)$. Given the camera parameters $C$ and shape $\bar{x}$, we can examine the ray corresponding to this pixel and check whether it is consistent with the observation $o_p$. We define a ray consistency cost function $L_p(\bar{x}, C; v_p)$ to capture the error associated with the pixel $p$. The view consistency loss can then be defined as the sum of per-pixel errors $L(\bar{x}, C; V) \equiv \sum_p L_p(\bar{x}, C; v_p)$.

**Sampling Occupancies along a Ray.** To define the consistency cost function $L_p(\bar{x}, C; v_p)$, we need to consider the ray as it is passing through the probabilistically occupied voxel grid $\bar{x}$. We do so by looking at discrete points sampled along the ray. Concretely, we sample points at a pre-defined set of $N = 80$ depth values $\{d_i | 1 \le i \le N\}$

along each ray. We denote by $x_i^p$ the occupancy value at the $i^{th}$ sample along this ray. To determine $x_i^p$, we look at the 3D coordinate of the corresponding point. Note that this can be determined using the camera parameters. Given the camera intrinsic parameters $(f_u, f_v, u_0, v_0)$, the ray corresponding to the image pixel $(u, v)$ travels along the direction $(\frac{u-u_0}{f_u}, \frac{v-v_0}{f_v}, 1)$ in the camera frame. Therefore, the $i^{th}$ point along the ray, in the camera coordinate frame, is located at $l_i \equiv (\frac{u-u_0}{f_u} d_i, \frac{v-v_0}{f_v} d_i, d_i)$. Then, given the camera extrinsics $(R, t)$, we can compute the location of his point in the coordinate frame of the predicted shape $\bar{x}$. Finally, we can use trilinear sampling to determine the occupancy at this point by sampling the value at this using the occupancies $\bar{x}$. Denoting by $T(G, pt)$ a function that samples a volumetric grid $G$ at a location $pt$, we can compute the occupancy sampled at the $i^{th}$ as below.

$$x_i^p = \mathcal{T}(\bar{x}, R \times (l_i + t)); \tag{4}$$

$$l_i \equiv (\frac{u - u_0}{f_u} d_i, \frac{v - v_0}{f_v} d_i, d_i) \tag{5}$$

Note that since the trilinear sampling function $T$ is differentiable w.r.t its arguments, the sampled occupancy $x_i^p$ is differentiable w.r.t the shape $\bar{x}$ and the camera $C$.

**Probabilistic Ray Tracing.** We have so far considered the ray associated with a pixel $p$ and computed samples with corresponding occupancy probabilities along it. We now trace this ray as it travels forward and use the samples along the ray as checkpoints. In particular, we assume that when the ray reaches the point corresponding to the $i^{th}$ sample, it either travels forward or terminates at that point. Conditioned on the ray reaching this sample, it travels forward with probability $x_i^p$ and terminates with likelihood $(1 - x_i^p)$. We denote by $\mathbf{z}^p \in \{1, \cdots, N+1\}$ a random variable corresponding to the sample index where the ray (probabilistically) terminates, where $z^p = N + 1$ implies that the ray escapes. We call these probabilistic ray terminations as *ray termination events*

and can compute the probability distribution $q(z_p)$ for these.

$$q(z^p = i) = (1 - x_i^p) \prod_{j=1}^{i-1} x_j^p \quad \forall(i \leq N); \qquad (6)$$

$$q(z^p = N + 1) = \prod_{j=1}^{N} x_j^p; \qquad (7)$$

**Event Costs.** Each event corresponds to the ray terminating at a particular point. It is possible to assign a cost to each event based on how inconsistent it is to w.r.t the pixel value $v_p$. If we have a depth observation $v_p \equiv d_p$, we can penalize the event $z^p = i$ by measuring the difference between $d_p$ and $d_i$. Alternatively, if we have a foreground image observation *i.e.* $v_p \equiv s_p \in \{0, 1\}$ where $s_p = 1$ implies a foreground pixel, we can penalize all events which correspond to a different observation. We can therefore define a cost function $\psi_p(i)$ which computes the cost associated with event $z_p = i$.

$$\psi_p^{depth}(i) = |d_p - d_i|; \qquad (8)$$

$$\psi_p^{mask}(i) = |s_p - \mathbb{1}(i \leq N)|; \qquad (9)$$

**Ray Consistency Cost.** We formulated the concept of ray termination events, and associated a probability and a cost to these. The ray consistency cost is then defined as the expected event cost.

$$L_p(\bar{x}, C; v_p) \;=\; \mathop{\mathbb{E}}_{z_p} \psi_p(z_p) \;=\; \sum_{i=1}^{N} q(z_p = i)\psi_p(i) \qquad (10)$$
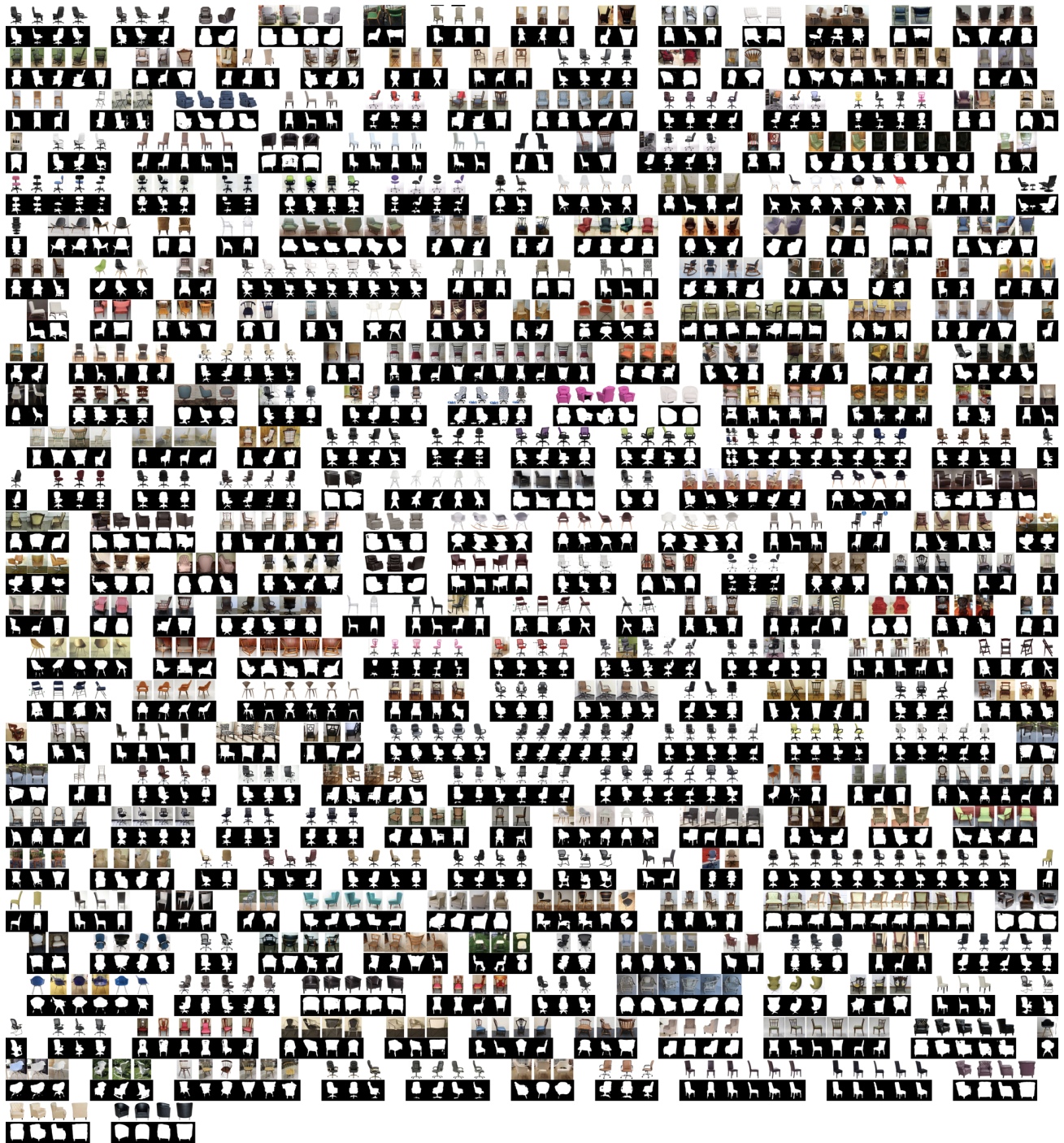
Note that the probabilities $q(z_p = i)$ are a differentiable function of $x_p$ which, in turn, is a differentiable function of shape $\bar{x}$ and camera $C$. The view consistency loss, which is simply a sum of multiple ray consistency terms, is therefore also differentiable w.r.t the shape and pose.

**Relation to Previous Work.** The formulation presented draws upon previous work on differentiable ray consistency [30] and leverages the notions of probabilistic ray termination events and event costs to define the ray consistency loss. A crucial difference however, is that we, using trilinear sampling, compute occupancies for point samples along the ray instead of directly using the occupancies of the voxels in the ray's path. Unlike their formulation, this allows our loss to also be differentiable w.r.t pose which is a crucial requirement for our scenario. Yan *et al.* [35] also use a similar sampling trick but their formulation is restricted to specifically using mask verification images and is additionally not leveraged for learning about pose. Tulsiani *et al.* [30] also discuss how their formulation can be adapted to use more general verification images *e.g.* color, semantics *etc.* using additional per-voxel predictions. While our experiments presented in the main text focus on leveraging mask or depth verification images, a similar generalization is possible for our formulation.

## A2. Online Product Images Dataset

We used the 'chair' object category from the Stanford Online Products Dataset [29]. To obtain associated foreground masks for these images, the semantic segmentation system from Chen *et al.* [7], where for each image, the mask was indicated by the pixels with most likely class label as 'chair'. As the obtained segmentation masks were often incorrect, or objects in the images truncated/occluded, we manually selected images of unoccluded/untruncated instances with a reasonably accurate (though still noisy) predicted segmentation. For our training, we only used the object instances with atleast 2 valid views. This resulting dataset is visualized in Figure 8. The result visualizations shown in the main text are using images from the original online products dataset [29], but correspond to objects instances that were not used for our training (due to lack of a sufficient number of valid views).

**Figure 8:** Training instances for the online products dataset. We visualize all the training images used along with their (approximate) segmentation masks, with images from the same object grouped together.