

Few-Shot Viewpoint Estimation

Hung-Yu Tseng¹

htseng6@ucmerced.edu

Shalini De Mello²

shalinig@nvidia.com

Jonathan Tremblay²

jtremblay@nvidia.com

Sifei Liu²

sifeil@nvidia.com

Stan Birchfield²

sbirchfield@nvidia.com

Ming-Hsuan Yang¹

mhyang@ucmerced.edu

Jan Kautz²

jkautz@nvidia.com

¹ University of California, Merced

² NVIDIA



Figure 1: **Few-shot viewpoint estimation.** Given only a few images of a novel category with annotated viewpoints (left images with rendered CAD models²), we aim to learn to predict the viewpoint of arbitrary objects from the same category.

Abstract

Viewpoint estimation for known categories of objects has been improved significantly thanks to deep networks and large datasets, but generalization to *unknown* categories is still very challenging. With an aim towards improving performance on unknown categories, we introduce the problem of category-level few-shot viewpoint estimation. We design a novel framework to successfully train viewpoint networks for new categories with few examples (10 or less). We formulate the problem as one of learning to estimate category-specific 3D canonical shapes, their associated depth estimates, and semantic 2D keypoints. We apply meta-learning to learn weights for our network that are amenable to category-specific few-shot fine-tuning. Furthermore, we design a flexible meta-Siamese network that maximizes information sharing during meta-learning. Through extensive experimentation on the ObjectNet3D and Pascal3D+ benchmark datasets, we demonstrate that our framework, which we call MetaView, significantly outperforms fine-tuning the state-of-the-art models with few examples, and that the specific architectural innovations of our method are crucial to achieving good performance.

1 Introduction

Estimating the viewpoint (azimuth, elevation, and cyclorotation) of rigid objects, relative to the camera, is a fundamental problem in three-dimensional (3D) computer vision. It is vital to applications such as robotics [34], 3D model retrieval [8], and reconstruction [10]. With convolutional neural networks (CNNs) and the availability of many labeled examples [8, 10, 20], much progress has been made in estimating the viewpoint of *known* categories of objects [6, 10, 20]. However, it remains challenging for even the best methods [13] to generalize well to *unknown* categories that the system has not encountered during training [10, 17, 13]. In such a case, re-training the viewpoint estimation network on an unknown category would require annotating thousands of new examples, which is labor-intensive.

To improve the performance of viewpoint estimation on unknown categories with little annotation effort, we introduce the problem of *few-shot viewpoint estimation*, in which a few (10 or less) labeled training examples are used to train a viewpoint estimation network for each novel category. We are inspired by the facts that (a) humans are able to perform mental rotations of objects [28] and can successfully learn novel views from a few examples [18]; and (b) recently, successful few-shot learning methods for several other vision tasks have been proposed [8, 10, 14].

However, merely fine-tuning a viewpoint estimation network with a few examples of a new category can easily lead to over-fitting. To overcome this problem, we formulate the viewpoint estimation problem as one of learning to estimate category-specific 3D canonical keypoints, their 2D projections, and associated depth values from which viewpoint can be estimated. We use meta-learning [10, 1] to learn weights for our viewpoint network that are optimal for category-specific few-shot learning. Furthermore, we propose meta-Siamese, which is a flexible network design that maximizes information sharing during meta-learning and adapts to an arbitrary number of keypoints. Through extensive evaluation on the ObjectNet3D [10] and Pascal3D+ [10] benchmark datasets, we show that our proposed method helps to significantly improve performance on unknown categories and outperforms fine-tuning the state-of-the-art models with a few examples of new categories.

To summarize, the main scientific contributions of our work are:

- We introduce the problem of category-level few-shot viewpoint estimation, thus bridging viewpoint estimation and few-shot learning.
- We design a novel meta-Siamese architecture and adapt meta-learning to learn weights for it that are optimal for category-level few-shot learning.

2 Related work

Viewpoint estimation. Many viewpoint estimation networks have been proposed for single [10, 10, 16] or multiple [8, 10] categories; or individual instances [10, 10] of objects. They use different network architectures, including those that estimate angular values directly [10, 15, 10, 16, 10]; encode images in latent spaces to match them against a dictionary of ground truth viewpoints [10, 10]; or detect projections of 3D bounding boxes [8, 10, 10, 10] or of semantic keypoints [10, 10], which along with known [10] or estimated [8, 10] 3D object structures are used to compute viewpoint. Zhou et al. propose the state-of-the-art StarMap method that detects multiple visible general keypoints [10] similar to SIFT [10] or SURF [10] via a learned CNN, and estimates category-level canonical 3D shapes. The existing viewpoint estimation methods are designed for known object categories

²We do not use the CAD models in our method, and we show them here for the purpose of illustrating viewpoint.

and hence very few works report performance on unknown ones [40, 31, 43]. Even highly successful techniques such as [43] perform significantly worse on unknown categories versus known ones. To our knowledge, no prior work has explored few-shot learning as a means to improve performance on novel categories and our work is the first to do so.

The existing viewpoint estimation networks also require large training datasets and two of them: Pascal3D+ [44] and ObjectNet3D [45] with 12 and 100 categories, respectively, have helped to move the field forward. At the instance level, the LineMOD [8], T-LESS [9], OPT [33], and YCB-Video [46] datasets that contain images of no more than 30 known 3D objects are widely used. Manual annotation of object viewpoint by aligning 3D CAD models to images (e.g., Figure (1)); or of 2D keypoints is a significant undertaking. To overcome this limitation, viewpoint estimation methods based on unsupervised learning [32]; general keypoints [23]; and synthetic images [20, 30, 31, 35, 40] have been proposed.

Few-shot learning. Successful few-shot learning algorithms for several vision tasks, besides viewpoint estimation, have been proposed recently. These include object recognition [4, 24, 25, 26, 29, 38], segmentation [23, 24], online adaptation of trackers [29], and human motion prediction [8]. Several of these methods use meta-learning [10] to learn a “learner” that is amenable to few-shot learning of a specific task from a set of closely related tasks. The learner may take the form of (a) a training algorithm [4, 10, 24]; (b) a metric-space for representing tasks [29, 33]; or (c) a meta-recurrent network [29, 26]. The MAML [10] meta-learning algorithm that learns a set of network initialization weights that are optimal for few-shot fine-tuning, is shown to be useful for many vision tasks.

Relative to the existing work, in this work we train networks for category-level viewpoint estimation. We further assume that we do not have access to 3D CAD models of any object or category. Lastly, we endeavor to train viewpoint networks for new categories with very few examples—a task that has not been attempted previously.

3 Few-shot Viewpoint Estimation

Our proposed MetaView framework for category-level few-shot viewpoint estimation is shown in the top row of Figure 2. It consists of two main components: a category-agnostic feature extraction block designed to extract general features from images that help to improve the accuracy of the downstream viewpoint estimation task; and a category-specific viewpoint estimation block designed to compute the viewpoint of all objects of a specific category. The latter block, in turn, computes viewpoint by detecting a unique set of semantic keypoints (containing 3D, 2D and depth values) via a category-specific feature extraction module ($f_{\theta_{cat}}$) and a category-specific keypoint detection module ($f_{\theta_{key}}$).

Our system operates in the following manner. We first train each of our feature extraction and viewpoint estimation blocks using a training set S^{train} containing a finite set of object categories. We use standard supervised learning to train the feature extraction block and fix its weights for all subsequent training stages. We then use meta-learning to train our viewpoint estimation block. It uses an alternative training procedure designed to make the viewpoint estimation block an effective few-shot “learner”. This means that when our trained viewpoint estimation block is further fine-tuned with a few examples of an unknown category, it can generalize well to other examples of that category.

At inference time, we assume that our system encounters a new category (not present during training) along with a few of its labeled examples from another set S^{test} (e.g., the category “monitor” shown in the lower part of Figure 2). We construct a unique viewpoint estimation network for it, initialize its weights with the optimal weights θ_{cat}^* and θ_{key}^* learned

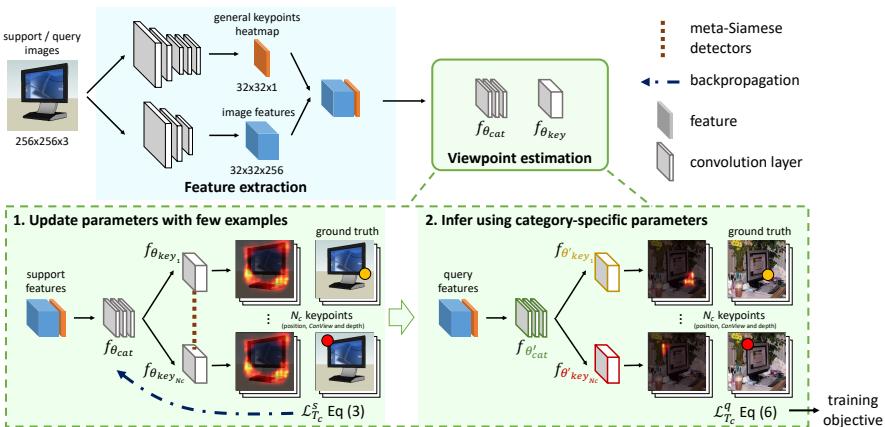


Figure 2: Method overview. Our MetaView framework is composed of a category-agnostic feature extraction (top-left) and category-specific viewpoint estimation (top-right) blocks. The bottom components show the different steps for training our viewpoint estimation block via meta-learning or for adapting it to a new category (bottom left only), which are described in detail in Section 3.

during meta-learning, and fine-tune it with the new category’s few labeled examples (lower left of Figure 2). This results in a category-specific viewpoint network that generalizes well to other examples of this new category (lower right of Figure 2). In the following sections, we describe the architecture and the training procedure of each component in more detail.

3.1 Feature Extraction

The first stage of our pipeline is a feature extraction block (top left of Figure 2), which we train and use to extract features without regard to an object’s category. It consists of two ResNet-18-style [7] networks: one trained as described in [4] to extract a multi-peak heatmap for the locations of many visible general keypoints (see examples in the supplementary material); and another whose first four convolutional blocks compute an identically-sized set of high-level convolutional features and is trained to detect 8 semantic keypoints for all categories by optimizing the loss in Eq. (6) described later in Section 3.2.2. We concatenate the multi-peak heatmap and high-level features and input them to the viewpoint estimation block. We train the feature extraction block via standard supervised SGD learning and once trained, we fix its weights for all subsequent steps.

3.2 Viewpoint Estimation

Our viewpoint estimation block (top right in Figure 2) is specific to each category. It computes a 3D canonical shape for each category, along with its 2D image projection and depth values; and relates these quantities to compute an object’s viewpoint. Furthermore, it is trained via meta-learning to be an optimal few-shot “learner” for any new category. We describe its architecture and training procedure in the following sections.

3.2.1 Architecture

Viewpoint estimation via semantic keypoints. We assume that we have no knowledge of the 3D shape of any object in a category. So, to compute viewpoint, inspired by [43], we train our viewpoint estimation block to estimate a set of 3D points $\{(x_k, y_k, z_k) | k = 1 \dots N_c\}$,

which together represent a canonical shape for the entire category \mathcal{T}_c in an object-centric coordinate system (e.g., for the category ‘‘chairs’’ it may comprise of the corners of a stick-figure representation of a prototypical chair with a back, a seat, and 4 legs). Additionally, for each 3D point k , our network detects its 2D image projection (u_k, v_k) and estimates its associated depth d_k . We refer collectively to the values $(x_k, y_k, z_k), (u_k, v_k), d_k$ of a point k as a ‘‘semantic keypoint’’. Finally, we obtain the viewpoint (rotation) of an object by solving the set of equations that relate each of the k rotated and projected 3D canonical points (x_k, y_k, z_k) to its 2D image location and depth estimate (u_k, v_k, d_k) , via orthogonal Procrustes. Note that our viewpoint estimation block is different from that of Zhou et al.’s [33] as they detect the 2D projections of only the *visible* 3D canonical points, whereas we detect projections of *all* visible and invisible ones, thus providing more data for estimating viewpoint.

Semantic keypoint estimation. To locate the 2D image projection (u_k, v_k) of each 3D keypoint k , the output of our network is a 2D heatmap $h_k(u, v)$, which predicts the probability of the point being located at (u, v) . It is produced by a spatial softmax layer. We obtain the final image coordinates (u_k, v_k) via a weighted sum of the row (u) and column (v) values as:

$$(u_k, v_k) = \sum_{u,v} (u, v) \cdot h_k(u, v). \quad (1)$$

Our network similarly computes a 2D map of depth values $c_k(u, v)$ that is of the same size as $h_k(u, v)$, along with three more maps $m_k^{i=\{x,y,z\}}(u, v)$ for each dimension of its 3D canonical keypoint. The final depth estimate d_k and 3D keypoint (x_k, y_k, z_k) is computed as:

$$d_k = \sum_{u,v} c_k(u, v), \quad (x_k, y_k, z_k) = \sum_{u,v} m_k^{i=\{x,y,z\}}(u, v) \cdot h_k(u, v). \quad (2)$$

Category-specific keypoints estimation. Given a category \mathcal{T}_c , our viewpoint estimation block must detect its unique N_c semantic keypoints via a category-specific feature extractor $f_{\theta_{cat}}$ followed by a set of category-specific semantic keypoint detectors $\{f_{\theta_{key_k}} | k = 1 \dots N_c\}$ (lower left of Figure 2). Each keypoint detector $f_{\theta_{key_k}}$ detects one unique category-specific semantic keypoint k , while the feature extractor $f_{\theta_{cat}}$ computes the common features required by all of them. Since our viewpoint estimation block must adapt to multiple different categories with different numbers of semantic keypoints, it cannot have a fixed number of pre-defined keypoint detectors. To flexibly change the number of keypoint detectors for each novel category, we propose a meta-Siamese architecture (lower left of Figure 2), which we operate as follows. For each new category \mathcal{T}_c , we replicate a generic pre-trained keypoint detector $(f_{\theta_{key}})$ N_c times and train each copy to detect *one* unique keypoint k of the new category, thus creating a specialized keypoint-detector with a unique and different number of semantic keypoints $\{f_{\theta_{key_k}} | k = 1 \dots N_c\}$ for each new category.

3.2.2 Training

Our goal is to train the viewpoint estimation block to be an effective *few-shot learner*. In other words, its learned feature extractor $f_{\theta_{cat}^*}$ and semantic keypoint detector $f_{\theta_{key}^*}$, after being fine-tuned with a few examples of a new category (lower left in Figure 2), should effectively extract features for the new category and detect each of its unique keypoints, respectively. To learn the optimal weights $\theta^* = \{\theta_{cat}^*, \theta_{key}^*\}$ that make our viewpoint estimation block amenable to few-shot fine-tuning without catastrophically over-fitting for a new category, we adopt the MAML meta-learning algorithm [9].

MAML optimizes a special meta-objective using a standard optimization algorithm, e.g., SGD. In standard supervised learning the objective is to minimize only the *training* loss for a task during each iteration of optimization. However, the meta-objective in MAML is to explicitly minimize, during each training iteration, the *generalization* loss for a task *after* a network has been trained with a few of its labeled examples. Furthermore, it samples a random task from a set of many such related tasks available for training during each iteration. We describe our specific meta-training algorithm to learn the optimal weights $\theta^* = \{\theta_{cat}^*, \theta_{key}^*\}$ for our viewpoint estimation block as follows.

During each iteration of meta-training, we sample a random task from S^{train} . A task comprises of a support set D_c^s and a query set D_c^q , each containing 10 and 3 labeled examples, respectively, of a category \mathcal{T}_c . The term “shot” refers to the number of examples in the support set D_c^s . For this category, containing N_c semantic keypoints, we replicate our generic keypoint detector ($f_{\theta_{key}}$) N_c times to construct its unique meta-Siamese keypoints detector with the parameters $\tilde{\theta} \leftarrow [\theta_{cat}, \theta_{key_1}, \theta_{key_2}, \dots, \theta_{key_{N_c}}]$ (lower left in Figure 2) and initialize each θ_{key_k} with θ_{key} . We use the category-specific keypoint detector to estimate its support set’s semantic keypoints and given their ground truth values, we compute the following loss:

$$\mathcal{L}_{\mathcal{T}_c}^s = \lambda_{2D}\mathcal{L}_{2D} + \lambda_{3D}\mathcal{L}_{3D} + \lambda_d\mathcal{L}_d, \quad (3)$$

where \mathcal{L}_{2D} , \mathcal{L}_{3D} , and \mathcal{L}_d are the average L_2 regression losses for correctly estimating the semantic keypoints’ 2D and 3D positions, and depth estimates, respectively. The λ parameters control the relative importance of each loss term. We compute the gradient of this loss $\mathcal{L}_{\mathcal{T}_c}^s$ w.r.t. to the network’s parameters $\tilde{\theta}$ and use a single step of SGD to update $\tilde{\theta}$ to $\tilde{\theta}'$ with a learning rate of α :

$$\tilde{\theta}' \leftarrow \tilde{\theta} - \alpha \nabla_{\tilde{\theta}} \mathcal{L}_{\mathcal{T}_c}^s. \quad (4)$$

Next, with the updated model parameters $\tilde{\theta}'$, we compute the loss $\mathcal{L}_{\mathcal{T}_c}^q$ for the query set D_c^q of this category (lower right in Figure 2). To compute the query loss, in addition to the loss terms described in (3), we also use a weighted concentration loss term:

$$\mathcal{L}_{con} = \frac{1}{N_c} \sum_{k=1}^{N_c} \sum_{u,v} h_k(u,v) \| [u_k, v_k]^\top - [u, v]^\top \|_2, \quad (5)$$

which forces the distribution of a 2D keypoint’s heatmap $h_k(u, v)$ to be peaky around the predicted position (u_k, v_k) . We find that this concentration loss term helps to improve the accuracy of 2D keypoint detection. Our final query loss is:

$$\mathcal{L}_{\mathcal{T}_c}^q = \lambda_{2D}\mathcal{L}_{2D} + \lambda_{3D}\mathcal{L}_{3D} + \lambda_d\mathcal{L}_d + \lambda_{con}\mathcal{L}_{con}. \quad (6)$$

The generalization loss of our network $\mathcal{L}_{\mathcal{T}_c}^q$, after it has been trained with just a few examples of a specific category, serves as the final meta-objective that is minimized in each iteration of meta-training and we optimize the network’s initial parameters θ w.r.t. its query loss $\mathcal{L}_{\mathcal{T}_c}^q$ using:

$$\theta_{cat} \leftarrow \theta_{cat} - \beta \nabla_{\theta_{cat}} \mathcal{L}_{\mathcal{T}_c}^q(f_{\tilde{\theta}}), \quad (7)$$

$$\theta_{key} \leftarrow \theta_{key} - \beta \frac{1}{N_c} \sum_{k=1..N_c} \left[\nabla_{\theta_{key_k}} \mathcal{L}_{\mathcal{T}_c}^q(f_{\tilde{\theta}}) \right]. \quad (8)$$

We repeat the meta-training iterations until our viewpoint estimation block converges to f_{θ^*} , as presented in Algorithm 1. Notice that in Eq. (8) we compute the optimal weights for

the generic keypoint detector θ_{key} by averaging the gradients of all the duplicated keypoint detectors θ_{key_k} . We find that this novel design feature of our network along with its shared category-level feature extractor with parameters θ_{cat} help to improve accuracy. They enable efficient use of *all* the available keypoints to learn the optimal values for θ_{cat} and θ_{key} during meta-training, which is especially important when training data is scarce.

Algorithm 1 MetaView Meta-training

```

1: Require: a set of tasks  $S^{train}$ 
2: randomly initialize  $\theta_{key}$  and  $\theta_{cat}$ 
3: while training do
4:   sample one task  $\mathcal{T}_c \sim S^{train}$ 
5:    $\triangleright$  meta-Siamese keypoint detectors
6:    $\theta_{key_1}, \theta_{key_2}, \dots, \theta_{key_{N_c}} \leftarrow \theta_{key}$ 
7:    $\triangleright$  viewpoint estimator
8:    $\tilde{\theta} \leftarrow \left[ \theta_{cat}, \theta_{key_1}, \theta_{key_2}, \dots, \theta_{key_{N_c}} \right]$ 
9:    $\triangleright$  update viewpoint estimator using support set
10:   $\tilde{\theta}' \leftarrow \tilde{\theta} - \alpha \nabla_{\tilde{\theta}} \mathcal{L}_{\mathcal{T}_c^s}(f_{\tilde{\theta}})$ 
11:   $\triangleright$  meta learning optimization using query set
12:   $\theta_{cat} \leftarrow \theta_{cat} - \beta \nabla_{\theta_{cat}} \mathcal{L}_{\mathcal{T}_c^q}^q(f_{\tilde{\theta}'})$ 
13:   $\theta_{key} \leftarrow \theta_{key} - \beta \frac{1}{N_c} \sum_{k=1..N_c} \left[ \nabla_{\theta_{key_k}} \mathcal{L}_{\mathcal{T}_c^q}^q(f_{\tilde{\theta}'}) \right]$ 
14: end while

```

3.2.3 Inference

We evaluate the performance of how well our viewpoint estimation block f_{θ^*} , which is learned via meta-learning performs at the task of adapting to unseen categories. Similar to meta-training, we sample a category from S^{test} with the same shot size as used for training. We construct its unique viewpoint estimation network $f_{\tilde{\theta}^*}$ and fine-tune it with a few of its examples by minimizing the loss in Eq. (3). This results in a optimal few-shot trained network $f_{\tilde{\theta}^{*\prime}}$ for this category. We then evaluate the generalization performance of $f_{\tilde{\theta}^{*\prime}}$ on all testing images of that category. We repeat this procedure for all categories in S^{test} and for multiple randomly selected few-shot training samples per category, and average across all of them.

4 Results

Implementation details. We provide detailed descriptions of our CNN architectures, and their training procedures in the supplementary material, to limit the number of pages.

Experiments. We evaluate our method for two different experimental settings. First, we follow the *intra*-dataset experiment of [43] and split the categories in ObjectNet3D [42] into 76 and 17 for training and testing, respectively. Secondly, we conduct an *inter*-dataset experiment. From ObjectNet3D, we exclude the 12 categories that are also present in Pascal3D+ [41]. We then use the remaining 88 categories in ObjectNet3D for training and test on Pascal3D+. Complying with [56], we discard the images with occluded or truncated objects from the test set in both experiments. We use two metrics for evaluation: 1) *Acc30*, which is the percentage of views with a rotational error less than 30° and 2) *MedErr*, which is the median rotational error across a dataset, measured in degrees. We compute the rota-

Table 1: Intra-dataset experiment. We report $Acc30(\uparrow)$ / $MedErr(\downarrow)$. All models are trained and evaluated on 76 and 17 categories from ObjectNet3D, respectively. The “zero” methods don’t use images of unknown categories for training and all others involve few-shot learning.

Method	bed	bookshelf	calculator	cellphone	computer	f_cabinet	guitar	iron	knife	microwave
StarMap (zero)	0.37 / 45.1	0.69 / 18.5	0.19 / 61.8	0.51 / 29.8	0.74 / 15.6	0.78 / 14.1	0.64 / 20.4	0.02 / 142	0.08 / 136	0.89 / 12.2
StarMap* (zero)	0.31 / 45.0	0.63 / 22.2	0.27 / 52.2	0.51 / 29.8	0.64 / 24.2	0.78 / 15.8	0.52 / 28.0	0.00 / 134	0.06 / 124	0.82 / 16.9
Baseline (zero)	0.26 / 49.1	0.57 / 25.0	0.78 / 53.3	0.38 / 45.5	0.66 / 20.3	0.73 / 18.7	0.39 / 44.6	0.06 / 135	0.08 / 127	0.82 / 16.8
StarMap* + fine-tune	0.32 / 47.2	0.61 / 21.0	0.26 / 50.6	0.56 / 26.8	0.59 / 24.4	0.76 / 17.1	0.54 / 27.9	0.00 / 128	0.05 / 120	0.82 / 19.0
Baseline + fine-tune	0.28 / 43.7	0.67 / 22.0	0.77 / 18.4	0.45 / 34.6	0.67 / 22.7	0.67 / 21.5	0.27 / 52.1	0.02 / 127	0.06 / 108	0.85 / 16.6
StarMap* + MAML	0.32 / 42.2	0.76 / 15.7	0.58 / 26.8	0.59 / 22.2	0.69 / 19.2	0.76 / 15.5	0.59 / 21.5	0.00 / 136	0.08 / 117	0.82 / 17.3
Ours	0.36 / 37.5	0.76 / 17.2	0.92 / 12.3	0.58 / 25.1	0.70 / 22.2	0.66 / 22.9	0.63 / 24.0	0.20 / 76.9	0.05 / 97.9	0.77 / 17.9

Method	pot	rifle	slipper	stove	toilet	tub	wheelchair	TOTAL		
StarMap (zero)	0.50 / 30.0	0.00 / 104	0.11 / 146	0.82 / 12.0	0.43 / 35.8	0.49 / 31.8	0.14 / 93.8	0.44 / 39.3		
StarMap* (zero)	0.51 / 29.2	0.02 / 97.4	0.10 / 130	0.81 / 13.9	0.44 / 34.4	0.37 / 37.0	0.17 / 74.4	0.43 / 39.4		
Baseline (zero)	0.46 / 38.8	0.00 / 98.6	0.09 / 123	0.82 / 14.8	0.32 / 39.5	0.29 / 50.4	0.14 / 71.6	0.38 / 44.6		
StarMap* + fine-tune	0.51 / 29.9	0.02 / 100	0.08 / 128	0.80 / 16.1	0.38 / 36.8	0.35 / 39.8	0.18 / 80.4	$0.41 \pm 0.00 / 41.0 \pm 0.22$		
Baseline + fine-tune	0.38 / 39.1	0.01 / 107	0.03 / 123	0.72 / 21.6	0.31 / 39.9	0.28 / 48.5	0.15 / 70.8	$0.40 \pm 0.02 / 39.1 \pm 1.79$		
StarMap* + MAML	0.51 / 28.2	0.01 / 100	0.15 / 128	0.83 / 15.6	0.39 / 35.5	0.41 / 38.5	0.24 / 71.5	$0.46 \pm 0.01 / 33.9 \pm 0.16$		
Ours	0.49 / 31.6	0.21 / 80.9	0.07 / 115	0.74 / 21.7	0.50 / 32.0	0.29 / 46.5	0.27 / 55.8	0.48 ± 0.01 / 31.5 ± 0.72		

Table 2: Inter-dataset experiment. We report $Acc30(\uparrow)$ / $MedErr(\downarrow)$. All models are trained on ObjectNet3D and evaluated on Pascal3D+. The “zero” methods don’t use images of unknown categories for training and all others involve few-shot learning.

Method	aero	bike	boat	bottle	bus	car	chair
StarMap (zero)	0.04 / 97.7	0.10 / 90.42	0.14 / 78.42	0.81 / 16.7	0.54 / 29.4	0.25 / 67.8	0.19 / 97.3
StarMap* (zero)	0.02 / 112	0.02 / 102	0.06 / 110	0.44 / 34.3	0.48 / 32.7	0.18 / 87.0	0.29 / 70.0
Baseline (zero)	0.03 / 114	0.06 / 101	0.10 / 95	0.41 / 36.6	0.36 / 42.0	0.14 / 93.7	0.26 / 71.5
StarMap* + fine-tune	0.03 / 102	0.05 / 98.8	0.07 / 98.9	0.48 / 31.9	0.46 / 33.0	0.18 / 80.8	0.22 / 74.6
Baseline + fine-tune	0.02 / 113	0.04 / 112	0.11 / 93.4	0.39 / 37.1	0.35 / 39.9	0.11 / 99.0	0.21 / 75.0
StarMap* + MAML	0.03 / 99.2	0.08 / 88.4	0.11 / 92.2	0.55 / 28.0	0.49 / 31.0	0.21 / 81.4	0.21 / 80.2
Ours	0.12 / 104	0.08 / 91.3	0.09 / 108	0.71 / 24.0	0.64 / 22.8	0.22 / 73.3	0.20 / 89.1

Method	table	mbike	sofa	train	tv	TOTAL	
StarMap (zero)	0.62 / 23.3	0.15 / 70.0	0.23 / 49.0	0.63 / 25.7	0.46 / 31.3	0.32 / 50.1	
StarMap* (zero)	0.43 / 31.7	0.09 / 86.7	0.26 / 42.5	0.30 / 46.8	0.59 / 24.7	0.25 / 71.2	
Baseline (zero)	0.38 / 39.0	0.11 / 82.3	0.39 / 57.5	0.29 / 50.0	0.63 / 24.3	0.24 / 70.0	
StarMap* + fine-tune	0.46 / 31.4	0.09 / 91.6	0.32 / 44.7	0.36 / 41.7	0.52 / 29.1	0.25 ± 0.01 / 64.7 ± 1.07	
Baseline + fine-tune	0.41 / 35.1	0.09 / 79.1	0.32 / 58.1	0.29 / 51.3	0.59 / 29.9	0.22 ± 0.02 / 69.2 ± 1.48	
StarMap* + MAML	0.29 / 36.8	0.11 / 83.5	0.44 / 42.9	0.42 / 33.9	0.64 / 25.3	0.28 ± 0.00 / 60.5 ± 0.10	
Ours	0.39 / 36.0	0.14 / 74.7	0.29 / 46.2	0.61 / 23.8	0.58 / 26.3	0.33 ± 0.02 / 51.3 ± 4.28	

tional error as $E_R = \frac{\|\log(R_{gt}^\top R)\|_F}{\sqrt{2}}$, where $\|\cdot\|_F$ is the Frobenius norm, and R_{gt} and R are the ground truth and predicted rotation matrices, respectively.

Comparisons. We compare several viewpoint estimation networks to ours. These include:

- **StarMap:** The original StarMap method [43]. It contains two stages of an Hourglass network [16] as the backbone and computes a multi-peak heatmap of general visible keypoints, and their depth and canonical 3D points.
- **StarMap*:** Our re-implementation of StarMap [43] with one stage of ResNet-18 [9] as the backbone for a fair comparison to ours.
- **StarMap* + MAML:** The StarMap* network trained with MAML for few-shot viewpoint estimation.
- **Baseline:** The ResNet-18 network trained to detect a fixed number (8) of semantic keypoints for all categories via standard supervised learning.

For methods that involve few-shot fine-tuning on unknown categories (*i.e.*, StarMap* or Baseline with fine-tuning, StarMap + MAML, and Ours), we use a shot size of 10. We repeat each experiment ten times with random initial seeds and report their average performance. Note that we also attempted to train viewpoint estimation networks that estimate angular values directly (e.g., [40]); or those that detect projections of 3D bounding boxes (e.g., [8]).

Table 3: Ablation study. The table shows the individual contributions of our meta-Siamese design (MS), the concentration loss (L_{con}), and general keypoints heatmap (KP) on the performance of MetaView in the intra-dataset experiment. We report $Acc30(\uparrow)/MedErr(\downarrow)$.

Method	bed	bookshelf	calculator	cellphone	computer	f_cabinet	guitar	iron	knife	microwave
Ours	0.28 / 42.3	0.68 / 23.1	0.87 / 15.3	0.47 / 32.1	0.63 / 24.9	0.71 / 22.1	0.03 / 100	0.15 / 76.0	0.01 / 121	0.69 / 23.2
Ours (MS)	0.27 / 42.4	0.77 / 22.2	0.74 / 24.0	0.54 / 28.3	0.64 / 24.9	0.63 / 25.3	0.13 / 76.9	0.05 / 103	0.65 / 26.2	
Ours (MS, L_{con})	0.31 / 41.3	0.79 / 19.0	0.84 / 17.4	0.53 / 28.0	0.62 / 25.9	0.66 / 23.6	0.35 / 35.8	0.16 / 86.5	0.05 / 101	0.81 / 17.7
Ours (MS, L_{con} , KP)	0.36 / 37.5	0.76 / 17.2	0.92 / 12.3	0.58 / 25.1	0.70 / 22.2	0.66 / 22.9	0.63 / 24.0	0.20 / 76.9	0.05 / 97.9	0.77 / 17.9

Method	pot	rifle	slipper	stove	toilet	tub	wheelchair	TOTAL
Ours	0.46 / 32.1	0.04 / 119	0.02 / 125	0.81 / 19.5	0.15 / 51.2	0.26 / 45.9	0.02 / 109	$0.35 \pm 0.01 / 42.5 \pm 1.15$
Ours (MS)	0.34 / 37.4	0.18 / 78.8	0.05 / 111	0.71 / 21.5	0.37 / 34.8	0.24 / 44.8	0.10 / 76.1	$0.41 \pm 0.01 / 36.0 \pm 0.78$
Ours (MS, L_{con})	0.49 / 31.2	0.16 / 90.5	0.05 / 111	0.75 / 21.7	0.41 / 34.4	0.31 / 42.4	0.22 / 60.8	$0.45 \pm 0.01 / 33.6 \pm 0.94$
Ours (MS, L_{con} , KP)	0.49 / 31.6	0.21 / 80.9	0.07 / 115	0.74 / 21.7	0.50 / 32.0	0.29 / 46.5	0.27 / 55.8	0.48 \pm 0.01 / 31.5 \pm 0.72

Table 4: Shot size. We report $Acc30(\uparrow)/MedErr(\downarrow)$. The table shows the effect of varying the number of support images (“shot size”) during meta-training and testing in the intra-dataset experiments with ObjectNet3D.

Method	bed	bookshelf	calculator	cellphone	computer	f_cabinet	guitar	iron	knife	microwave
Ours (1 shot)	0.24 / 45.7	0.16 / 70.8	0.26 / 56.7	0.19 / 57.3	0.41 / 32.6	0.48 / 31.4	0.06 / 76.8	0.02 / 125	0.01 / 120	0.18 / 48.4
Ours (5 shots)	0.31 / 39.9	0.50 / 29.6	0.67 / 25.1	0.34 / 48.6	0.67 / 23.7	0.66 / 24.0	0.34 / 40.4	0.09 / 91.7	0.04 / 110	0.81 / 16.7
Ours (10 shots)	0.36 / 37.5	0.76 / 17.2	0.92 / 12.3	0.58 / 25.1	0.70 / 22.2	0.66 / 22.9	0.63 / 24.0	0.20 / 76.9	0.05 / 97.9	0.77 / 17.9

Method	pot	rifle	slipper	stove	toilet	tub	wheelchair	TOTAL
Ours (1 shot)	0.39 / 36.8	0.00 / 102	0.05 / 121	0.36 / 35.8	0.33 / 39.1	0.11 / 75.1	0.12 / 81.5	$0.21 \pm 0.05 / 55.2 \pm 6.82$
Ours (5 shots)	0.51 / 29.4	0.05 / 107	0.04 / 110	0.74 / 21.1	0.38 / 35.7	0.27 / 46.7	0.23 / 61.3	$0.41 \pm 0.03 / 36.2 \pm 1.58$
Ours (10 shots)	0.49 / 31.6	0.21 / 80.9	0.07 / 115	0.74 / 21.7	0.50 / 32.0	0.29 / 46.5	0.27 / 55.8	0.48 \pm 0.01 / 31.5 \pm 0.72

with MAML, but they either failed to converge or performed very poorly. So, we do not report results for them. The results of the intra-dataset and inter-dataset experiments are presented in Table 1 and Table 2, respectively.

Zero-shot performance. For both experiments, methods trained using standard supervised learning solely on the training categories (*i.e.*, StarMap, StarMap* and Baseline denoted by “zero”) are limited in their ability to generalize to unknown categories. For the original StarMap method [43] in the intra-dataset experiment (Table 1), the overall $Acc30$ and $MedErr$ worsen from 63% and 17°, respectively, when the test categories are known to the system to 44% and 39.3°, respectively, when they are unknown. This indicates that the existing state-of-the-art viewpoint estimation networks require information that is unique to each category to infer its viewpoint. Since the original StarMap [43] uses a larger backbone network than ResNet-18 [7] it performs better than our implementation (StarMap*) of it.

Few-shot performance. Among the methods that involve few-shot fine-tuning for unknown categories, methods that are trained via meta-learning (StarMap + MAML and our MetaView) perform significantly better than the methods that are not (StarMap* or Baseline with fine-tuning) in both the intra- and inter-dataset experiments. These results are the first demonstration of the effectiveness of meta-learning at the task of category-level few-shot viewpoint learning. Furthermore, in both experiments, our MetaView framework results in the best overall performance of all the zero- and few-shot learning methods. It outperforms StarMap* + MAML, which shows the effectiveness of our novel design components that differentiate it from merely training StarMap* with MAML. They include our network’s ability to (a) detect the 2D locations and depth values of all 3D canonical points and not just the visible ones; (b) share information during meta-learning via the meta-Siamese design; and (c) flexibly construct networks with a different number of keypoints for each category. Lastly, observe that even with a smaller backbone network, our method performs better than the current best performance for the task of viewpoint estimation of unknown categories, *i.e.* of StarMap [43] “zero” and thus helps to improve performance on unknown categories with

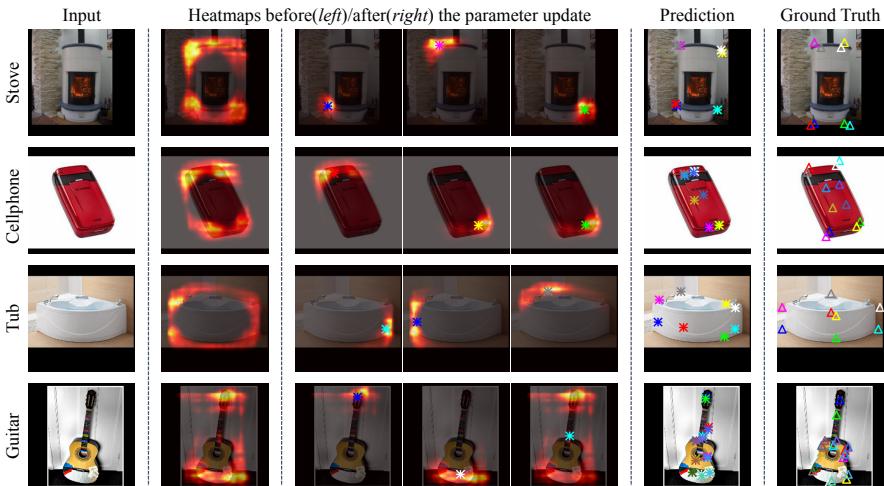


Figure 3: Qualitative results of the intra-dataset experiment. We show the keypoint detection results of MetaView on unknown categories, before and after few-shot fine-tuning. The images from left to right are: the input image, the 2D keypoint heatmap before fine-tuning with Eq. (3), three example heatmaps for specific keypoints after fine-tuning, all the predicted keypoints, and their ground truth values.

very little additional labeling effort.

The effectiveness of MetaView is also evident from Figure 3, which shows examples of the 2D keypoint heatmaps $h_k(u, v)$ (described in Section 3.1) produced by it before and after few-shot fine-tuning with examples of new categories. The keypoint detector, prior to few-shot fine-tuning, is not specific to any keypoint and generates heatmaps that tend to have high responses on corners, edges or regions of the foreground object. After fine-tuning, however, it successfully learns to detect keypoints of various new categories and produces heatmaps with more concentrated peaks.

Ablation study. To validate the effectiveness of our various novel design components including our meta-Siamese design, concentration loss term (L_{con}), and of using the general keypoints’ multi-peak map as input, we show the results of an ablation study for the inter-dataset experiment in Table 3. While each component individually contributes to the overall performance, the concentration loss and the meta-Siamese design contribute the most.

Shot size. We vary the number of support images (*i.e.*, shot size to 1, 5 and 10) for each new category during meta-training and -testing. The results of this experiment for the intra-dataset setting are presented in Table 4. We observe that as more training images per category are available for training, the accuracy of our MetaView approach scales up correspondingly.

5 Conclusion

To improve performance on unknown categories, we introduce the problem of category-level few-shot viewpoint estimation. We propose the novel MetaView framework that successfully adapts to unknown categories with few labeled examples and helps to improve performance on them with little additional annotation effort. Our meta-Siamese keypoint detector is general and can be explored in the future for other few-shot tasks requiring keypoints detection.

References

- [1] Marcin Andrychowicz, Misha Denil, Sergio Gomez, Matthew W Hoffman, David Pfau, Tom Schaul, Brendan Shillingford, and Nando De Freitas. Learning to learn by gradient descent by gradient descent. In *NIPS*, 2016.
- [2] Herbert Bay, Andreas Ess, Tinne Tuytelaars, and Luc Van Gool. Speeded-up robust features (surf). *CVIU*, 110(3):346–359, 2008.
- [3] 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.
- [4] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *ICML*, 2017.
- [5] Alexander Grabner, Peter M Roth, and Vincent Lepetit. 3D pose estimation and 3D model retrieval for objects in the wild. In *CVPR*, 2018.
- [6] Liang-Yan Gui, Yu-Xiong Wang, Deva Ramanan, and José MF Moura. Few-shot human motion prediction via meta-learning. In *ECCV*, 2018.
- [7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [8] S. Hinterstoisser, V. Lepetit, S. Ilic, S. Holzer, G. Bradski, K. Konolige, and N. Navab. Model based training, detection and pose estimation of texture-less 3D objects in heavily cluttered scenes. In *ACCV*, 2012.
- [9] Tomáš Hodaň, Pavel Haluza, Štěpán Obdržálek, Jiří Matas, Manolis Lourakis, and Xenophon Zabulis. T-LESS: An RGB-D dataset for 6D pose estimation of texture-less objects. In *WACV*, 2017.
- [10] Wadim Kehl, Fabian Manhardt, Federico Tombari, Slobodan Ilic, and Nassir Navab. SSD-6D: Making RGB-based 3D detection and 6D pose estimation great again. In *ICCV*, 2017.
- [11] Abhijit Kundu, Yin Li, and James M Rehg. 3D-RCNN: Instance-level 3D object reconstruction via render-and-compare. In *CVPR*, 2018.
- [12] Alina Kuznetsova, Sung Ju Hwang, Bodo Rosenhahn, and Leonid Sigal. Exploiting view-specific appearance similarities across classes for zero-shot pose prediction: A metric learning approach. In *AAAI*, 2016.
- [13] David G Lowe. Distinctive image features from scale-invariant key-points. *IJCV*, 60(2):91–110, 2004.
- [14] Francisco Massa, Bryan C Russell, and Mathieu Aubry. Deep exemplar 2D-3D detection by adapting from real to rendered views. In *CVPR*, 2016.
- [15] Arsalan Mousavian, Dragomir Anguelov, John Flynn, and Jana Košecká. 3D bounding box estimation using deep learning and geometry. In *CVPR*, 2017.

- [16] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation. In *ECCV*, 2016.
- [17] Alex Nichol, Joshua Achiam, and John Schulman. On first-order meta-learning algorithms. In *arXiv:1803.02999*, 2018.
- [18] Stephen E Palmer. *Vision science: Photons to phenomenology*. MIT press, 1999.
- [19] Eunbyung Park and Alexander C Berg. Meta-tracker: Fast and robust online adaptation for visual object trackers. In *ECCV*, 2018.
- [20] Georgios Pavlakos, Xiaowei Zhou, Aaron Chan, Konstantinos G Derpanis, and Kostas Daniilidis. 6-dof object pose from semantic key-points. In *ICRA*, 2017.
- [21] Mahdi Rad and Vincent Lepetit. BB8: A scalable, accurate, robust to partial occlusion method for predicting the 3D poses of challenging objects without using depth. In *ICCV*, 2017.
- [22] Mahdi Rad, Markus Oberweger, and Vincent Lepetit. Feature mapping for learning fast and accurate 3D pose inference from synthetic images. In *CVPR*, 2018.
- [23] Kate Rakelly, Evan Shelhamer, Trevor Darrell, Alexei A Efros, and Sergey Levine. Few-shot segmentation propagation with guided networks. *arXiv preprint arXiv:1806.07373*, 2018.
- [24] Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. In *ICLR*, 2017.
- [25] Danilo Jimenez Rezende, Shakir Mohamed, Ivo Danihelka, Karol Gregor, and Daan Wierstra. One-shot generalization in deep generative models. *JMLR*, 48, 2016.
- [26] Adam Santoro, Sergey Bartunov, Matthew Botvinick, Daan Wierstra, and Timothy Lillicrap. Meta-learning with memory-augmented neural networks. In *ICML*, 2016.
- [27] Amirreza Shaban, Shravy Bansal, Zhen Liu, Irfan Essa, and Byron Boots. One-shot learning for semantic segmentation. *arXiv preprint arXiv:1709.03410*, 2017.
- [28] Roger N Shepard and Jacqueline Metzler. Mental rotation of three-dimensional objects. *Science*, 171(3972):701–703, 1971.
- [29] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In *NIPS*, 2017.
- [30] Hao Su, Charles R. Qi, Yangyan Li, and Leonidas J. Guibas. Render for CNN: Viewpoint estimation in images using CNNs trained with rendered 3D model views. In *ICCV*, 2015.
- [31] Martin Sundermeyer, Zoltan-Csaba Marton, Maximilian Durner, Manuel Brucker, and Rudolph Triebel. Implicit 3D orientation learning for 6D object detection from RGB images. In *ECCV*, 2018.
- [32] Supasorn Suwajanakorn, Noah Snavely, Jonathan Tompson, and Mohammad Norouzi. Discovery of latent 3D key-points via end-to-end geometric reasoning. In *NIPS*, 2018.

- [33] Bugra Tekin, Sudipta N Sinha, and Pascal Fua. Real-time seamless single shot 6D object pose prediction. In *CVPR*, 2017.
- [34] Jonathan Tremblay, Thang To, Artem Molchanov, Stephen Tyree, Jan Kautz, and Stan Birchfield. Synthetically trained neural networks for learning human-readable plans from real-world demonstrations. In *ICRA*, 2018.
- [35] Jonathan Tremblay, Thang To, Balakumar Sundaralingam, Yu Xiang, Dieter Fox, and Stan Birchfield. Deep object pose estimation for semantic robotic grasping of household objects. In *CoRL*, 2018.
- [36] Shubham Tulsiani and Jitendra Malik. Viewpoints and key-points. In *CVPR*, 2015.
- [37] Shubham Tulsiani, Joao Carreira, and Jitendra Malik. Pose induction for novel object categories. In *ICCV*, 2015.
- [38] Oriol Vinyals, Charles Blundell, Tim Lillicrap, Daan Wierstra, et al. Matching networks for one shot learning. In *NIPS*, 2016.
- [39] Po-Chen Wu, Yueh-Ying Lee, Hung-Yu Tseng, Hsuan-I Ho, Ming-Hsuan Yang, and Shao-Yi Chien. A benchmark dataset for 6dof object pose tracking. In *ISMAR Adjunct*, 2017.
- [40] Y. Xiang, T. Schmidt, V. Narayanan, and D. Fox. PoseCNN: A convolutional neural network for 6D object pose estimation in cluttered scenes. In *RSS*, 2018.
- [41] Yu Xiang, Roozbeh Mottaghi, and Silvio Savarese. Beyond PASCAL: A benchmark for 3D object detection in the wild. In *WACV*, 2014.
- [42] Yu Xiang, Wonhui Kim, Wei Chen, Jingwei Ji, Christopher Choy, Hao Su, Roozbeh Mottaghi, Leonidas Guibas, and Silvio Savarese. ObjectNet3D: A large scale database for 3D object recognition. In *ECCV*, 2016.
- [43] Xingyi Zhou, Arjun Karpur, Linjie Luo, and Qixing Huang. Starmap for category-agnostic key-point and viewpoint estimation. In *ECCV*, 2018.