# Learning Canonical 3D Object Representation for Fine-Grained Recognition

Sunghun Joung[1], Seungryong Kim[2], Minsu Kim[1], Ig-Jae Kim[3], Kwanghoon Sohn[1,*]

[1]Yonsei University, [2]Korea University, [3]Korea Institute of Science and Technology (KIST)

{sunghunjoung,minsukim320,khsohn}@yonsei.ac.kr, seungryong_kim@korea.ac.kr, drjay@kist.re.kr

## Abstract

*We propose a novel framework for fine-grained object recognition that learns to recover object variation in 3D space from a single image, trained on an image collection without using any ground-truth 3D annotation. We accomplish this by representing an object as a composition of 3D shape and its appearance, while eliminating the effect of camera viewpoint, in a canonical configuration. Unlike conventional methods modeling spatial variation in 2D images only, our method is capable of reconfiguring the appearance feature in a canonical 3D space, thus enabling the subsequent object classifier to be invariant under 3D geometric variation. Our representation also allows us to go beyond existing methods, by incorporating 3D shape variation as an additional cue for object recognition. To learn the model without ground-truth 3D annotation, we deploy a differentiable renderer in an analysis-by-synthesis framework. By incorporating 3D shape and appearance jointly in a deep representation, our method learns the discriminative representation of the object and achieves competitive performance on fine-grained image recognition and vehicle re-identification. We also demonstrate that the performance of 3D shape reconstruction is improved by learning fine-grained shape deformation in a boosting manner.*

## 1. Introduction

Object recognition [34, 23, 17, 76] is one of the most fundamental and essential tasks in computer vision fields, which has achieved steady progress by the advent of deep convolutional neural networks. However, it still remains a challenging problem, especially when an object undergoes severe geometric deformations, *e.g.*, by object scale, pose and part variations, which frequently occur across different instances, or by camera viewpoint changes [16, 25, 8, 27].
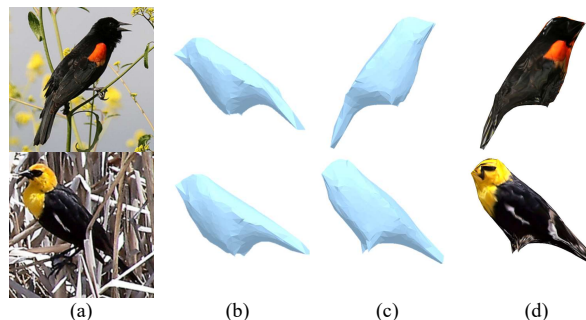
Figure 1. **Intuition of our method:** Given (a) 2D image, we recover object variation in 3D space by estimating (b) 3D shape deformation, (c) camera viewpoint change and (d) appearance variation. It allows for using 3D shape and appearance in a canonical space, while eliminating camera viewpoint variation, enabling us to deal with 3D object variations and facilitating the subsequent object classifier.

To overcome these challenges, recent works [25, 39, 53, 7, 8] seek to handle such geometric variations based on an assumption that the object variation can be decomposed into *appearance* and *2D spatial variation*. They first estimate an appearance flow from an input and then warp the input into a canonical configuration so as to remove the spatial variation, from which the appearance feature is extracted to facilitate the subsequent classifier's task. The appearance flow is generally estimated by modeling 2D transformation [25, 39, 53], *e.g.*, affine transformation, or by learning offset of sampling locations in the convolutional operators [7, 8]. These methods, however, do not account for the fact that the object variation, given an image, is due to the variations in *appearance*, *3D shape* and *camera viewpoint* as in Fig. 1. While the effect of camera viewpoint should be eliminated for achieving geometric invariance, 3D shape variation can be used as an additional cue to extract a shape feature that is able to supplement an appearance feature, but none of the existing methods utilize this.

However, estimating 3D object information from a single image is challenging, since collecting the ground-truth 3D shape is notoriously difficult and time-consuming [69], thus limiting the supervised learning for this task. To overcome this, some methods implicitly consider the 3D object

structure by learning a discriminative feature representation for fine-grained recognition. Formally, they use an extra module to localize the discriminative object parts [33, 20], by using explicit part detectors [1, 32] or implicit attention mechanisms [17, 75]. However, these methods can only localize a few semantic parts without understanding the holistic object structure, and can be limited if the network fails to consistently localize object parts across multiple instances.

In this paper, we present a method that estimates 3D object information in a canonical configuration, including 3D object shape and appearance, with camera viewpoint. It enables the subsequent classifier to directly work on the 3D object information, from which both appearance and shape features are simultaneously extracted. It allows for handling subtle intra-class variations by means of both appearance and 3D shape features, which is not available at the existing approaches. To this end, we deploy a differentiable renderer [29, 43], to infer 3D shape, without ground-truth 3D annotation, in an analysis-by-synthesis framework, as in recent 3D shape reconstruction methods [28, 24, 56, 49]. In particular, we incorporate this framework into an encoder-decoder architecture that disentangles the object variation to 3D shape, appearance, and camera viewpoint. To this end, an image is embedded into a low-dimensional latent code that is fed into separate decoders to estimate the aforementioned factors independently. We also exploit multiple hypothesis camera prediction to avoid local minima during training as in [24, 35, 19].

Unlike conventional methods [25, 39, 53] that model 2D spatial variation only, our method is capable of reconfiguring an appearance feature in a canonical space, where each semantic part of an object is mapped to the same location in a canonical space. Moreover, our method enables dense semantic alignment [73] into a canonical configuration, where positional encoding [42, 2] can further improve recognition performance, while the conventional methods [1, 32, 17, 75] are limited by highlighting only a few salient parts of an object. To improve recognition ability between subtle object variations, we further introduce a shape encoder to utilize 3D shape deformation as an additional cue. By incorporating 3D shape and appearance jointly in a deep representation, our method consistently boosts the discriminative representation learning on fine-grained image recognition and vehicle re-identification tasks. In addition, our joint learning framework enables us to improve 3D shape reconstruction capability by discriminating shape variations between different fine-grained categories.

## 2. Related Works

**Spatial Invariance.** Vanilla CNN [54, 23] provided limited performance under severe geometric variations. STN [25] offered a way to explicitly handle geometric variation by spatially warping the input to a canonical configuration

to facilitate the recognition task. Inspired by STN [25], many variants were proposed by using recurrent formalism [39], deformable convolution [7], polar transformation [13], and attention based warping [52, 76]. These methods typically employ an additional localization network, to predict appearance flow, which is then applied to intermediate features to remove spatial variation. Since they do not share the template shape with different images, modeling shape variability within a category is limited. On the other hand, several methods attempted to address the geometric variation of an object by decomposing the image into shape and appearance with instance-agnostic template shape. Thewlis *et al*. exploited dense coordinate frame [57, 58] in order to recover the deformation of an object with equivariance. Deforming autoencoders [53] proposed a generative model to predict a deformation field by disentangling shape and appearance. While effective, all of these methods model geometric deformation in 2D space, so they lack robustness against geometric variation in 3D space including 3D shape deformation and camera viewpoint variation.

**Fine-grained Object Recognition.** Since modeling severe deformations of 3D object in 2D image is challenging, conventional methods for fine-grained object recognition [15, 33, 1, 32, 40] aim to learn discriminative feature representation of object parts and then classify the object based on the discriminative regions. This, however, requires large human efforts as it needs extra annotation of bounding boxes or parts. To alleviate this, recent methods proposed to automatically localize the discriminative object parts without part annotation using attention mechanisms [17, 75, 55, 76, 31] in an unsupervised manner. However, a deeper understanding of the holistic 3D object shape is essential as the network may not consistently localize the discriminative object parts across multiple instances. Another line of research focuses on end-to-end feature encoding [18, 47, 11, 12, 74] to encourage feature discriminability using deeper representations, high-order feature interactions or metric learning.

**Vehicle Re-identification.** Vehicle re-identification task has gained more attention in recent years, following the release of several benchmarks [44, 41]. The main focus of the task has been on addressing viewpoint variation from 2D images. Given vehicle images under arbitrary camera viewpoints, recent methods aim to transform the input appearance feature into a viewpoint independent representation. Therefore, they utilize either local region based feature learning on pre-defined distinctive regions or keypoints [65, 21, 48, 4], or attention mechanisms [78, 30].

**Single-view 3D Reconstruction.** Single-view 3D reconstruction methods aim to reconstruct 3D object shape from a single image. While conventional methods utilize ground-truth 3D shape for training [63, 66], it requires tremen-
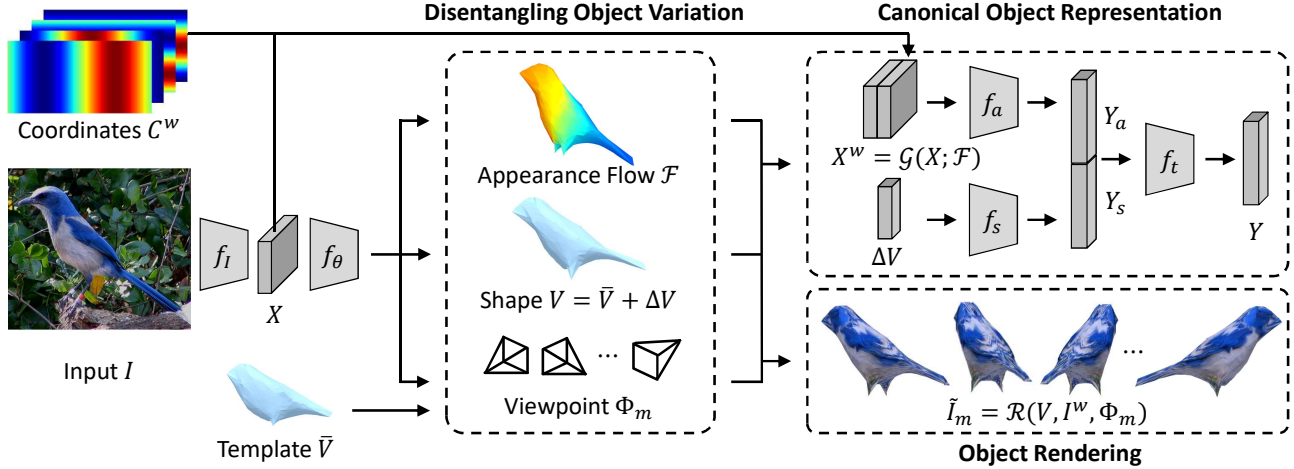
Figure 2. **Overview of our framework:** The input feature $X$, from an image encoder $f_I$, is fed into a module $f_\theta$ for disentangling object variation into appearance flow $\mathcal{F}$, shape deformation $\Delta V$ and camera viewpoint $\Phi$ as well as parameterized template shape $\overline{V}$, respectively. Then appearance feature $Y_a$ and shape feature $Y_s$ are obtained by applying the appearance encoder $f_a$ and shape encoder $f_s$ to get the final object representation $Y$. In our framework, the entire network is trained in a joint and boosting manner through fine-grained object recognition and 3D shape reconstruction tasks.

dous human efforts for annotation [68], or their applicability is restricted to synthetic data [3]. Therefore, several methods have proposed a differentiable renderer [29, 43] to train 3D reconstruction networks using either multi-view images or ground-truth camera viewpoints in an analysis by synthesis framework. To relax such constraints on supervision, Kanazawa *et al*. proposed CMR [28] to explore 3D reconstruction from a collection of images with different instances, by exploiting a learnable 3D template shape. Since CMR [28] requires annotated 2D keypoints for training, several works were proposed to mitigate this by using camera-multiplex [19], semantic consistency [37] or temporal consistency [36], which can estimate 3D mesh with foreground mask for training.

## 3. Method

### 3.1. Preliminaries

Let us denote an intermediate CNN feature representation as $X \in \mathbb{R}^{H \times W \times K}$, with height $H$, width $W$ and $K$ channels. To reduce the spatial variations among different instances within the representation, recent works [25, 39] predict an appearance flow $\mathcal{F}$ to transform the input feature into a canonical configuration, producing a warped feature $X^w = \mathcal{G}(X; \mathcal{F})$ through sampling function $\mathcal{G}$, *e.g*., a bilinear sampler [25]. The appearance flow is estimated via a module $f_\theta$, that takes $X$ as an input and outputs the transformation parameter (*e.g*., affine transformation) to produce appearance flow $\mathcal{F} \in \mathbb{R}^{H^w \times W^w \times 2}$, where each value in $\mathcal{F}(u, v)$ at point $(u, v)$ indicates the coordinates of the input to be sampled, and $H^w$, $W^w$ are the height and width of $X^w$. Conventional methods solely consider 2D geometric deformation and extract the features for appearance only

[25, 8]. However, since the geometric variation of objects occurs in 3D space, the warped feature is often inconsistent across instances under different camera viewpoints.

### 3.2. Motivation and Overview

In this paper, we conjecture that the object variation can be further decomposed into variations of 3D shape and appearance as well as camera viewpoint changes, and by only removing the latter, we can effectively handle complex object variation for fine-grained recognition. In other words, both 3D shape and appearance variation have to be encoded into object representation to discriminate subtle intra-class variation. The main bottleneck is learning to reconstruct the 3D shape. Traditional methods for 3D shape estimation, in training, relied on datasets labeled with 3D shapes [6, 14] or multiple views of the same object [59, 24]. Recent works [28, 19, 37] have relaxed the constraints of supervision by utilizing an analysis-by-synthesis framework that renders an instance-agnostic template shape into 2D images via a differentiable renderer [29, 43].

Inspired by this, we present a framework that learns the mapping function $f_\theta$ to recover 3D object information from a single image as illustrated in Fig. 2. By disentangling object variation into a composition of independent factors of 3D shape, appearance, and camera viewpoint variation, each component can be trained in a joint and boosting manner through object recognition and 3D shape reconstruction tasks. In order to learn these mappings without ground-truth 3D annotation, we represent the template shape of an object as 3D mesh in a canonical space, where the predicted object variation projects the mesh from this canonical space to the image coordinates. The use of canonical 3D space allows for assigning dense semantic correspondences across dif-
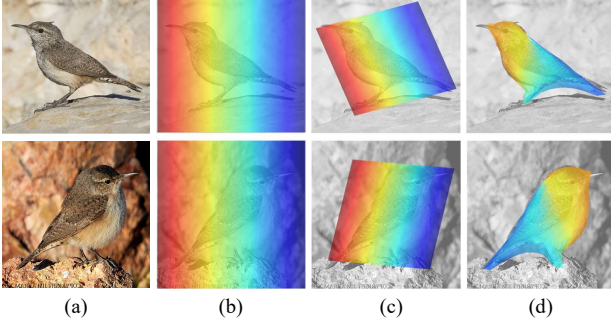
| (a) | (b) | (c) | (d) |

Figure 3. **Comparison of appearance flow:** (a) input images, (b) source coordinates in an image, and appearance flow obtained using (c) STN [25], and (d) ours. Points with the same color in different images are projected to the same point in a canonical space. This shows that our method can make the warped feature to be densely consistent across different camera viewpoints, while the existing methods fail (Best viewed in color).

ferent instances to be consistent. At the same time, we estimate an appearance flow to transform the 2D image into the canonical space, where the warped feature $X^w$ to be spatially invariant under 3D geometric variation. Since shape deformation itself can be an additional cue to discriminate object classes, we further present the shape encoder to maximize the classification performance.

### 3.3. Disentangling Object Variation

We model $f_\theta$ to project the input feature $X$ into a low dimensional latent code that is fed into three decoders to predict an appearance flow $\mathcal{F}$, shape deformation $\Delta V$, and camera viewpoint $\Phi$, such that $\{\mathcal{F}, \Delta V, \Phi\} = f_\theta(X)$.

**Appearance flow $\mathcal{F}$.** We first learn an appearance flow $\mathcal{F}$ to transform the input image defined in 2D to the canonical configuration so as to remove 3D geometric variation and enable to align appearance feature with similar semantics to the same location in a canonical space. As the topology of 3D mesh is fixed, we use $I^w = \mathcal{G}(I; \mathcal{F})$ to model the texture of the template mesh for object rendering as in [28, 37]. For appearance representation in a canonical space, we use $X^w = \mathcal{G}(X; \mathcal{F})$, where we take advantage of a pre-trained network, *i.e.*, ResNet-50 [23].

**Shape deformation $\Delta V$.** We represent the 3D shape in a form of mesh $M = (V, F)$ with vertices $V \in \mathbb{R}^{|V| \times 3}$ and faces $F$. The set of faces defines the connectivity of vertices in the spherical mesh, and we assume it remains fixed. The vertex positions of a deformable object are determined as the summation of an instance-specific deformation $\Delta V$ predicted from an image to a learned instance-agnostic mean shape $\overline{V}$ such that $V = \Delta V + \overline{V}$. Since most object categories exhibit bilateral symmetry [67, 27], we further constrain the predicted shape and deformation to be mirror-symmetric, following [28]. This symmetric constraint can be utilized to reduce the number of parameters for shape

representation as well as infer the invisible surface of an object from the visible features. We leverage $\Delta V$ to extract shape variation features.

**Camera viewpoint $\Phi$.** For camera parameters, we assume a weak-perspective projection, parameterized by scale $\mathbf{s} \in \mathbb{R}$, translation $\mathbf{t} \in \mathbb{R}^2$, and rotation captured by quaternion $\mathbf{r} \in \mathbb{R}^4$. We use $\Phi = (\mathbf{s}, \mathbf{t}, \mathbf{r})$ to denote the projection of 3D points sets on template shape into 2D image coordinates via the weak perspective projection. We optimize the overall network over the multi-hypothesis camera viewpoint [24, 35, 19], which has been well-known to be robust, by maintaining a set of possible camera hypotheses for each training instance, where $\mathcal{C} = \{\Phi_1, ..., \Phi_M\}$ denotes viewpoints with $M$ cameras.

### 3.4. Canonical Object Representation

In this section, we introduce canonical object representation that exploits the feature from appearance and shape deformation, invariant under arbitrary camera viewpoints.

**Embedding appearance.** As argued above, we explicitly learn to map the pixels in the object to their corresponding locations on template shape. It makes each point in $X^w$ densely and semantically aligned to a canonical 3D space as shown in Fig. 3. Note that conventional methods cannot enforce the warped feature to be consistent under different camera viewpoint [25, 39, 13], or only can localize a few semantic parts [1, 32, 17, 75].

In addition, to take advantage of the positional sensitivity, where semantically similar parts are densely mapped to a canonical space, we utilize positional encoding (PE), which has been shown to be effective in natural language processing [60, 10] and object recognition [42, 2]. The straightforward way is to exploit a coordinate map of the same spatial dimensions as $X^w$, normalized to be $u, v \in [-1, 1]$, following CoordConv [42]. This, however, cannot model the continuity of $X^w$ on a 3D mesh surface that is not in 2D planar space. To overcome this, we modify the original 2-dimensional $(u, v)$ pixel coordinates into 4-dimensional coordinates $(\cos(\pi u), \sin(\pi u), \cos(\pi v), \sin(\pi v))$, to provide periodic continuity of the 2D coordinates on 3D the surface. This canonical position map $C^w$ is then concatenated to the input $X^w$ and passed to the appearance encoder $f_a$ with two convolutional layers to output appearance representation $Y_a$ in a vectorized form. With this simple technique, the network can disambiguate different positions of an object instance, and also model the structural composition of object parts.

**Embedding shape deformation.** Unlike existing methods [25, 39, 13] that estimate geometric variation to remove spatial variation, we take advantage of modeling 3D shape. Note that the intra-class shape variation of a particular category has been utilized for 3D shape recognition
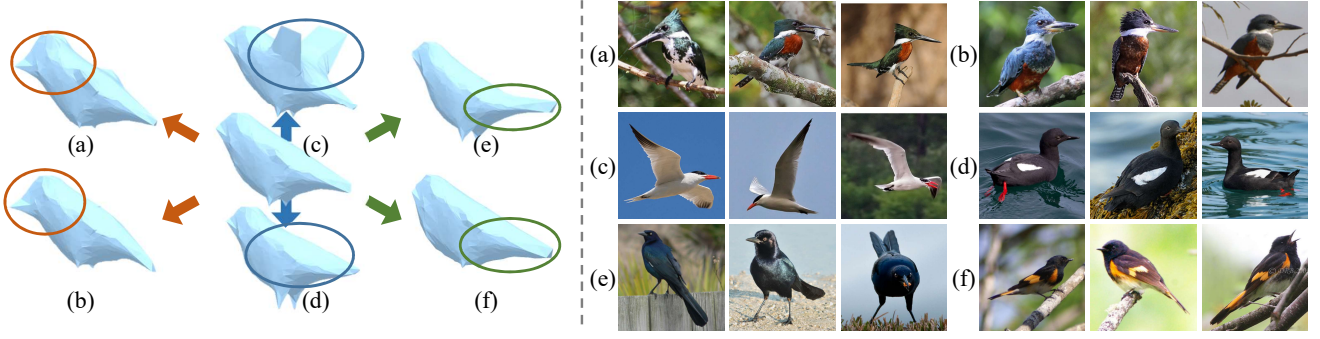
Figure 4. **Visualization of the learned 3D object shape deformations:** We visualize a template shape $\overline{V}$ and averaged shapes of 6 different fine-grained categories. Each shape characterizes a fine-grained category of, for instance, birds on (a,b) head, (c,d) body, and (e,f) tail types. We utilize such shape deformation as an additional cue to discriminate subtle intra-class variation for the object recognition task.

[69, 3] using 3D model as input, but none of the existing methods utilize it for 2D images. Rather than directly encoding shape representation $V$, we exploit instance-specific shape deformation $\Delta V$ since it shares the learned instance-agnostic mean shape $\overline{V}$, which is more suitable to discriminate subtle differences. We thus present an additional shape encoder, where the vectorized $\Delta V$ is fed into the shape encoder $f_s$ with several fully-connected layers to output shape representation $Y_s$.

**Fusing appearance and shape representation.** Given feature representation of appearance $Y_a$ and shape $Y_s$ in a canonical space, we concatenate them together to combine appearance and shape variations and passed to the target network $f_t$ with two fully-connected layers for the final object representation $Y$.

### 3.5. Loss Function

Since we train the networks on an image collection without any ground-truth 3D shape, multi-view, camera viewpoints or keypoint supervision, we follow unsupervised learning of category-specific mesh reconstruction by utilizing analysis-by-synthesis framework [19]. We predict multiple camera hypotheses to overcome a local minima, by making every instance maintain its own $\mathcal{C}$ independently, where we compute the loss against each camera viewpoint $\Phi_m$. By jointly learning our networks on fine-grained object classification and 3D shape reconstruction tasks, appearance flow $\mathcal{F}$ and shape deformation $\Delta V$ can be trained in a way that mutually boosts the two tasks. In the following, we describe loss functions to train our network in detail.

**Loss for disentangling.** We denote the rendered image as $\tilde{I}_m = \mathcal{R}(V, I^w, \Phi_m)$ and rendered silhouette mask as $\tilde{S}_m = \mathcal{R}(V, \Phi_m)$. We then compute the silhouette loss $\mathcal{L}_{\text{mask},m}$ and image reconstruction loss $\mathcal{L}_{\text{pixel},m}$ for $m$-th camera viewpoint as follows:

$$\mathcal{L}_{\text{mask},m} = \|S - \tilde{S}_m\|_2^2 + \text{dt}(S) * \tilde{S}_m, \quad (1)$$

$$\mathcal{L}_{\text{pixel},m} = \text{dist}(\tilde{I}_m \odot S, I \odot S), \quad (2)$$

where $\text{dt}(\cdot)$ denotes distance transform, $*$ is matrix multiplication, $\odot$ is element-wise multiplication, and $\text{dist}(\cdot)$ represents a perceptual distance metric [72]. Note that we only use foreground mask $S$ for training.

**Loss for priors.** In order to recover 3D shape of smooth surface, we apply smoothness loss $\mathcal{L}_{\text{smooth}} = \|LV\|_2$, where $L$ is the discrete Laplace-Beltrami operator to minimize the mean curvature [51]. We construct $L$ once using the initial template mesh following our baseline [28]. Furthermore, we apply deformation regularization loss $\mathcal{L}_{\text{reg}} = \|\Delta V\|_2$ to prevent large deformations.

**Overall training objective.** Since we predict a set of multiple hypothesis for camera projection, we first define the summation of silhouette and image reconstruction loss $\mathcal{L}_m = \mathcal{L}_{\text{mask},m} + \mathcal{L}_{\text{pixel},m}$ over the losses of $M$ camera pose prediction. We then compute $p_m = \frac{e^{-L_m/\sigma}}{\sum_n e^{-L_n/\sigma}}$, the probability of being the optimal camera, to associate with each $L_m$ with its probability $p_m$. The overall training objective is as follows, where we normalize each energy term according to its magnitudes as in [28] with task-specific loss function $\mathcal{L}_{\text{task}}$:

$$\mathcal{L}_{\text{total}} = \sum_m p_m \mathcal{L}_m + \mathcal{L}_{\text{smooth}} + \mathcal{L}_{\text{reg}} + \mathcal{L}_{\text{task}}. \quad (3)$$

### 3.6. Implementation Details

We implement our method using the Pytorch library [50]. In our experiments, we utilize ResNet-50 [23] pre-trained on ImageNet [9] as backbone. We build our module $f_\theta$ on the last convolutional layers of ResNet-50 [23]. For mesh representation, we use 642 vertices and 1280 faces correspond to Icosphere. We exploit 305 symmetric vertex pairs and 32 vertices without symmetry, resulting in $|V| = 337$. For the encoder, we follow [28], by first applying a convolutional layer to downsample the spatial and channel dimensions into 1/4 and 1/8, respectively, which is then vectorized to form a 4,096-D vector. We then apply

| Methods | CUB-Birds [61] | Stanford-Cars [33] |
|---|---|---|
| Base [23] | 74.6 | 70.4 |
| STN [25] | 76.5 | 71.0 |
| DCN [8] | 76.7 | 72.1 |
| SSN [52] | 77.7 | 74.8 |
| ASN [76] | 78.9 | 75.2 |
| VTN [31] | 83.1 | 82.7 |
| Ours w/o PE, $f_s$ | 83.7 | 85.5 |
| Ours w/o $f_s$ | 86.8 | 93.2 |
| Ours | **88.4** | **94.7** |

Table 1. Ablation study for the different components of our method on fine-grained image recognition.

| PE Module | CUB-Birds [61] | Stanford-Cars [33] |
|---|---|---|
| None | 83.7 | 85.5 |
| PE-2 | 85.2 | 89.8 |
| PE-4 | **86.8** | **93.2** |

Table 2. Ablation study for the different positional encoding (PE) modules on fine-grained image recognition.

| # of layers | CUB-Birds [61] | Stanford-Cars [33] |
|---|---|---|
| 0 | 86.8 | 93.2 |
| 1 | 87.0 | 93.6 |
| 2 | 87.6 | 94.3 |
| 3 | **88.4** | **94.7** |
| 4 | 87.9 | 94.4 |

Table 3. Ablation study for the different number of fully connected layers in shape encoder $f_s$ on fine-grained image recognition.

two fully-connected layers to get the shared latent code of size 200. This latent code is then fed into independent decoder networks with linear layers to predict shape deformation $\mathbb{R}^{|V|\times 3}$ and camera projection parameters $\mathbb{R}^7$. For the appearance flow, the latent code is fed into 5 upconvolution layers followed by $\tanh$ function to normalize the output into the $[-1, 1]$ coordinate space.

In addition, the input images are resized to a fixed resolution of $512 \times 512$, and $H^w, W^w$ are set as $256, 512$. Since the spatial resolution is different between $I$ and $X$, we set $H^w, W^w$ to learn texture mapping $I^w$, and downsample it into $1/32$ for warping feature $X$ into $X^w$. For task-specific loss function $\mathcal{L}_{\text{task}}$, we use the cross-entropy loss for fine-grained image recognition, and use cross-entropy and triplet loss for vehicle re-identification as in [48]. We use 3D template meshes in [35] as initial meshes of birds and cars, since it speeds up to convergence of the networks [28, 19]. For the mask label, we use ground-truth mask on CUB-Birds [61], and obtain fore-ground masks using off-the-shelf segmentation [22] for Stanford-Cars [33] and Veri-776 [44]. For all experiments, we use SoftRasterizer [43] as our 3D mesh rendering module $\mathcal{R}$ and follow the experimental protocols of [19] for training.

## 4. Experiments

### 4.1. Experimental Setup

In this section, we comprehensively analyze and evaluate our method on fine-grained image recognition and vehicle re-identification. First, we analyze the influence of the different components of our method on fine-grained image recognition. We then evaluate our method compared to the state-of-the-art methods. Finally, we evaluate 3D shape reconstruction performance compared to existing methods.

### 4.2. Ablation Study

We first analyze our method with the ablation studies, with respect to the different components, different PE modules and the different number of layers in the shape encoder,

followed by visual analysis of shape deformation.

**Analysis of the different components.** To validate the geometric invariance of our method, we compare with previous spatial deformation modeling methods, such as STN [25], DCN [8], SSN [52], ASN [76] and VTN [31] on fine-grained image recognition benchmarks including CUB-Birds [61] and Stanford-Cars [33]. For a fair comparison, we apply these methods at the same layers as ours, *i.e.*, the last convolutional layers of ResNet-50 [23]. As an ablation study, we evaluate our method with different components, only with spatial deformation modeling without positional encoding and shape encoder, denoted by Ours w/o PE, $f_s$ and without shape encoder, denoted by Ours w/o $f_s$. The results are provided in Tab. 1, where our method consistently outperforms the conventional methods. It is noticeable that the positional encoding module can only be used in our method since conventional methods [25, 52, 76] cannot map pixels on the same semantic part to a canonical space across different camera viewpoints. Furthermore, our method can use the holistic 3D structure of an instance thanks to its invariant nature under 3D geometric variation including 3D object shape and camera viewpoint variation.

**The effects on different PE module.** Tab. 2 shows experiments to validate the effect of various PE modules. We denote the method using 2-dimensional $u, v$ pixel coordinates [42] as PE-2, and the proposed methods using 4-dimensional coordinates as PE-4. Both PE-2 and PE-4 have shown higher accuracy by favoring a position sensitivity in a canonical appearance space. These results share the same properties as previous studies [42, 2], indicating that positional encoding benefits to clarify the spatial representation. In addition, PE-4 has shown better performance by utilizing sinusoidal functions to provide periodic continuity on 3D surface, rather than 2D space as in PE-2.
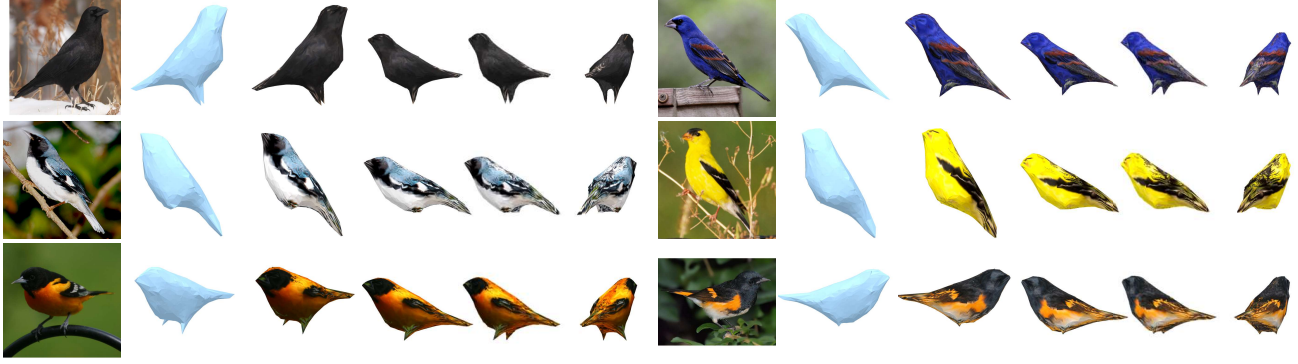
Figure 5. **Qualitative examples on CUB-Birds [61]:** (from left) input image $I$, its 3D shape $V$, and its rendered results $\tilde{I}_m$ from predicted camera viewpoint and from three other camera viewpoints.

**The effects on different shape encoder.** Since there is no reference model to utilize shape deformation from 2D images, we evaluate the performance with respect to the different number of layers in the model. For simplicity, we fixed the channel dimension of each fully connected layer to be 512, followed by ReLU, with a different number of layers. As the result with 3 fully connected layers has shown the best performance in Tab. 3, we set the number of layers as 3 for the remaining experiments.

**Visual analysis of shape deformation.** To analyze the discriminate capability of shape feature, we visualize the learned shape deformations in the validation set of CUB-Birds [61] as exemplified in Fig. 4. We averaged the estimated $\Delta V$ for each fine-grained category, which is associated with learned mean shape $\overline{V}$ such that $V = \overline{V} + \Delta V$ for visualization. We can see that each shape corresponds to the natural factors of fine-grained categories, such as long beak or round tail, and the statistics of captured scenes, such as a bird flying in the sky or floating on the water.

### 4.3. Comparison to Other Methods

**Fine-grained image recognition.** In the following, we evaluate our method with state-of-the-art methods on fine-grained image recognition benchmarks including CUB-Birds [61] and Stanford-Cars [33]. We compare with the methods using object part annotations, such as PN-CNN [1] and SPDA-CNN [71], using bounding box annotations, such as PA-CNN [32], MG-CNN [62], B-CNN [40] and FCAN [45] and using only images, such as RA-CNN [17], MA-CNN [75], DT-RAM [38], DFL-CNN [64], MAMC [55], NTSN [70], DCL [5], TASN [5], ACNet [26] and LIO [77]. The results are provided in Tab. 4 with a description of the backbone network, where our method achieves competitive performance. It is noticeable that mask annotations or segmentation results from off-the-shelf algorithms [22] allow us to model with 3D geometric variations while others fail even with part annotations.

| Methods | Backbone | [61] | [33] |
|---|---|---|---|
| PN-CNN [1] | AlexNet | 85.4 | - |
| SPDA-CNN [71] | VGG-19 | 85.1 | - |
| PA-CNN [32] | VGG-19 | 82.8 | 92.8 |
| MG-CNN [62] | VGG-19 | 83.0 | - |
| B-CNN [40] | 2×VGG-19 | 84.8 | 90.6 |
| FCAN [45] | ResNet-50 | 84.3 | 91.3 |
| RA-CNN [17] | 3×VGG-19 | 85.3 | 92.5 |
| MA-CNN [75] | 3×VGG-19 | 86.5 | 92.8 |
| DT-RAM [38] | ResNet-50 | 86.0 | 93.1 |
| DFL-CNN [64] | ResNet-50 | 87.4 | 93.1 |
| MAMC [55] | ResNet-50 | 86.5 | 93.0 |
| NTSN [70] | 3×ResNet-50 | 87.5 | 91.4 |
| DCL [5] | VGG-16 | 86.9 | 94.1 |
| | ResNet-50 | 87.8 | 94.5 |
| TASN [76] | VGG-19 | 86.1 | 93.2 |
| | ResNet-50 | 87.9 | 93.8 |
| ACNet [26] | VGG-16 | 87.8 | 94.3 |
| | ResNet-50 | 88.1 | 94.6 |
| LIO [77] | ResNet-50 | 88.0 | 94.5 |
| Ours | ResNet-50 | **88.4** | **94.7** |

Table 4. Comparison with the state-of-the-art fine-grained recognition methods on CUB-Birds [61] and Stanford-Cars [33].

**Vehicle re-identification.** We also evaluate our method on the task of vehicle re-identification using Veri-776 benchmark [44], where addressing subtle object variations from 2D images under different camera viewpoints is the main challenge. Following standard practice, we use the mean average precision (mAP), Cumulative Match Curve (CMC) for top 1 (CMC@1) and top 5 (CMC@5) matches for quantitative evaluation. We only use visual information without using either spatio-temporal information or license plate. We compare with the state-of-the-art methods, such as OIFE [65], VAMI [78], RAM [46], PRN [21], AAVER [30], PVEN [48] and SPAN [4]. As in Tab. 5, our method outperforms the conventional methods, thanks to its invariant nature under 3D geometric variation.

Figure 6. **Qualitative examples on Stanford-Cars [33]:** (from left) input image $I$, its 3D shape $V$, and its rendered result $\tilde{I}_m$ from predicted camera viewpoint.

| Methods | mAP | CMC@1 | CMC@5 |
|---|---|---|---|
| OIFE [65] | 0.480 | 0.659 | 0.877 |
| VAMI [78] | 0.501 | 0.770 | 0.908 |
| RAM [46] | 0.615 | 0.886 | 0.940 |
| PRN [21] | 0.743 | 0.943 | 0.989 |
| AAVER [30] | 0.612 | 0.890 | 0.947 |
| PVEN [48] | 0.795 | 0.956 | 0.984 |
| SPAN [4] | 0.689 | 0.940 | 0.976 |
| Ours | **0.801** | **0.959** | **0.991** |

Table 5. Comparison with the state-of-the-art vehicle re-identification methods on Veri-776 [44].

| Methods | Metric | |
|---|---|---|
| | Mask IoU | PCK |
| CMR [28] | 0.706 | 47.3 |
| U-CMR [19] | 0.712 | 49.4 |
| UMR [37] | 0.734 | 51.2 |
| Ours w/o PE, $f_s$ | 0.729 | 51.5 |
| Ours w/o $f_s$ | 0.732 | 51.9 |
| Ours | **0.737** | **52.1** |

Table 6. Comparison of 3D mesh reconstruction on CUB-Birds dataset [61]. Mask IoU and keypoint transfer (KT) are evaluated.

## 4.4. Discussion

Although promising results have been achieved in various tasks, several limitations still remain. In particular, our approach is not applicable to the categories where 3D shape across instances differ significantly or undergo large articulation, *e.g.*, human. It is also challenging to learn 3D shape from images with heavy occlusion, which limits its applicability to generic object recognition tasks in the wild. Nevertheless, this is an encouraging step towards understanding the underlying object variation in 3D space for image recognition tasks, and we hope it encourages future work by overcoming the aforementioned limitations.

## 5. Conclusion

We have introduced a novel framework to learn the discriminative representation of an object under 3D geometric variations, which accounts for the fact that object variation can be decomposed into the variation of 3D object shape and appearance in a canonical space, as well as camera viewpoint variation. We have developed a framework for disentangling of 3D shape, appearance, and camera viewpoint variation, trained without ground-truth 3D annotation. It allows for reconfiguring the appearance feature into a canonical space, enabling us to utilize a positional encoding for better representation learning. To deal with subtle object variation and improve recognition ability, we have further introduced a shape encoder to utilize 3D shape deformation. Our experiments have shown that our method effectively learns the discriminative object representation on fine-grained recognition and 3D shape reconstruction tasks.

**3D shape reconstruction.** Finally, we compare our method with 3D reconstruction methods including CMR [28], U-CMR [19] and UMR [37] on CUB-Birds [61]. Due to the lack of ground truth 3D shape, we evaluate mask reprojection accuracy using intersection over union (IoU) as in [28]. We also measure the keypoint reprojection accuracy using a percentage of correct keypoints (PCK) with a distance threshold $\alpha = 0.1$, by mapping a set of keypoints from the source image to the learned template and then to the target image via estimated shape deformation and camera viewpoint. Note that Stanford-Cars [33] and Veri-776 [44] only contain bounding box as annotations, which limit quantitative evaluation of 3D shape reconstruction.

We evaluate our method with ablation, without positional encoding and shape encoder, denoted by Ours w/o PE, $f_s$ and without shape encoder, denoted by Ours w/o $f_s$, compared to the state-of-the-art methods [28, 19, 37] in Tab. 6. We account that semantic keypoint [28] or co-part segmentation [37] can be seen as a coarse-level spatial representation of the object's part (*e.g.*, head or tail), while our method enables a dense spatial representation of the object's surface by utilizing a positional encoding module. In addition, our shape encoder allows us to learn accurate shape deformation by discriminating shape variations in different fine-grained categories. Since we aimed to use appearance and shape features for recognition while eliminating camera viewpoint variation, there was no significant performance gain or drop for viewpoint estimation compared to the baseline [19]. Qualitative results on CUB-Birds [61] and Stanford-Cars [33] are in Fig. 5 and Fig. 6.

# References

[1] Steve Branson, Grant Van Horn, Serge Belongie, and Pietro Perona. Bird species categorization using pose normalized deep convolutional nets. In *BMVC*, 2014.

[2] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, pages 213–229, 2020.

[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] Tsai-Shien Chen, Chih-Ting Liu, Chih-Wei Wu, and Shao-Yi Chien. Orientation-aware vehicle re-identification with semantics-guided part attention network. In *ECCV*, pages 330–346, 2020.

[5] Yue Chen, Yalong Bai, Wei Zhang, and Tao Mei. Destruction and construction learning for fine-grained image recognition. In *CVPR*, pages 5157–5166, 2019.

[6] Christopher B Choy, Danfei Xu, JunYoung Gwak, Kevin Chen, and Silvio Savarese. 3d-r2n2: A unified approach for single and multi-view 3d object reconstruction. In *ECCV*, pages 628–644, 2016.

[7] Jifeng Dai, Yi Li, Kaiming He, and Jian Sun. R-fcn: Object detection via region-based fully convolutional networks. In *NeurIPS*, pages 379–387, 2016.

[8] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *ICCV*, pages 764–773, 2017.

[9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, pages 248–255, 2009.

[10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT (1)*, pages 4171–4186, 2019.

[11] Yao Ding, Yanzhao Zhou, Yi Zhu, Qixiang Ye, and Jianbin Jiao. Selective sparse sampling for fine-grained image recognition. In *ICCV*, pages 6599–6608, 2019.

[12] Ruoyi Du, Dongliang Chang, Ayan Kumar Bhunia, Jiyang Xie, Zhanyu Ma, Yi-Zhe Song, and Jun Guo. Fine-grained visual classification via progressive multi-granularity training of jigsaw patches. In *ECCV*, pages 153–168, 2020.

[13] Carlos Esteves, Christine Allen-Blanchette, Xiaowei Zhou, and Kostas Daniilidis. Polar transformer networks. In *ICLR*, 2018.

[14] Haoqiang Fan, Hao Su, and Leonidas J Guibas. A point set generation network for 3d object reconstruction from a single image. In *CVPR*, pages 605–613, 2017.

[15] Ryan Farrell, Om Oza, Ning Zhang, Vlad I Morariu, Trevor Darrell, and Larry S Davis. Birdlets: Subordinate categorization using volumetric primitives and pose-normalized appearance. In *ICCV*, pages 161–168, 2011.

[16] Pedro F Felzenszwalb, Ross B Girshick, David McAllester, and Deva Ramanan. Object detection with discriminatively trained part-based models. *IEEE TPAMI*, 32(9):1627–1645, 2009.

[17] Jianlong Fu, Heliang Zheng, and Tao Mei. Look closer to see better: Recurrent attention convolutional neural network for fine-grained image recognition. In *CVPR*, pages 4438–4446, 2017.

[18] Weifeng Ge, Xiangru Lin, and Yizhou Yu. Weakly supervised complementary parts models for fine-grained image classification from the bottom up. In *CVPR*, pages 3034–3043, 2019.

[19] Shubham Goel, Angjoo Kanazawa, and Jitendra Malik. Shape and viewpoint without keypoints. In *ECCV*, pages 88–104, 2020.

[20] Pei Guo and Ryan Farrell. Aligned to the object, not to the image: A unified pose-aligned representation for fine-grained recognition. In *WACV*, pages 1876–1885, 2019.

[21] Bing He, Jia Li, Yifan Zhao, and Yonghong Tian. Part-regularized near-duplicate vehicle re-identification. In *CVPR*, pages 3997–4005, 2019.

[22] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *ICCV*, pages 2961–2969, 2017.

[23] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016.

[24] Eldar Insafutdinov and Alexey Dosovitskiy. Unsupervised learning of shape and pose with differentiable point clouds. In *NeurIPS*, pages 2802–2812, 2018.

[25] Max Jaderberg, Karen Simonyan, Andrew Zisserman, and Koray Kavukcuoglu. Spatial transformer networks. In *NeurIPS*, pages 2017–2025, 2015.

[26] Ruyi Ji, Longyin Wen, Libo Zhang, Dawei Du, Yanjun Wu, Chen Zhao, Xianglong Liu, and Feiyue Huang. Attention convolutional binary neural tree for fine-grained visual categorization. In *CVPR*, pages 10468–10477, 2020.

[27] Sunghun Joung, Seungryong Kim, Hanjae Kim, Minsu Kim, Ig-Jae Kim, Junghyun Cho, and Kwanghoon Sohn. Cylindrical convolutional networks for joint object detection and viewpoint estimation. In *CVPR*, pages 14163–14172, 2020.

[28] Angjoo Kanazawa, Shubham Tulsiani, Alexei A Efros, and Jitendra Malik. Learning category-specific mesh reconstruction from image collections. In *ECCV*, pages 371–386, 2018.

[29] Hiroharu Kato, Yoshitaka Ushiku, and Tatsuya Harada. Neural 3d mesh renderer. In *CVPR*, pages 3907–3916, 2018.

[30] Pirazh Khorramshahi, Amit Kumar, Neehar Peri, Sai Saketh Rambhatla, Jun-Cheng Chen, and Rama Chellappa. A dual-path model with adaptive attention for vehicle re-identification. In *ICCV*, pages 6132–6141, 2019.

[31] Seungryong Kim, Sabine Süsstrunk, and Mathieu Salzmann. Volumetric transformer networks. In *ECCV*, pages 561–578, 2020.

[32] Jonathan Krause, Hailin Jin, Jianchao Yang, and Li Fei-Fei. Fine-grained recognition without part annotations. In *CVPR*, pages 5546–5555, 2015.

[33] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *ICCVW*, pages 554–561, 2013.

[34] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *NeurIPS*, pages 1097–1105, 2012.

[35] Nilesh Kulkarni, Abhinav Gupta, and Shubham Tulsiani. Canonical surface mapping via geometric cycle consistency. In *ICCV*, pages 2202–2211, 2019.

[36] Xueting Li, Sifei Liu, Shalini De Mello, Kihwan Kim, Xiaolong Wang, Ming-Hsuan Yang, and Jan Kautz. Online adaptation for consistent mesh reconstruction in the wild. In *NeurIPS*, pages 15009–15019, 2020.

[37] Xueting Li, Sifei Liu, Kihwan Kim, Shalini De Mello, Varun Jampani, Ming-Hsuan Yang, and Jan Kautz. Self-supervised single-view 3d reconstruction via semantic consistency. In *ECCV*, pages 677–693, 2020.

[38] Zhichao Li, Yi Yang, Xiao Liu, Feng Zhou, Shilei Wen, and Wei Xu. Dynamic computational time for visual attention. In *ICCVW*, pages 1199–1209, 2017.

[39] Chen-Hsuan Lin and Simon Lucey. Inverse compositional spatial transformer networks. In *CVPR*, pages 2568–2576, 2017.

[40] Tsung-Yu Lin, Aruni RoyChowdhury, and Subhransu Maji. Bilinear cnn models for fine-grained visual recognition. In *ICCV*, pages 1449–1457, 2015.

[41] Hongye Liu, Yonghong Tian, Yaowei Yang, Lu Pang, and Tiejun Huang. Deep relative distance learning: Tell the difference between similar vehicles. In *CVPR*, pages 2167–2175, 2016.

[42] Rosanne Liu, Joel Lehman, Piero Molino, Felipe Petroski Such, Eric Frank, Alex Sergeev, and Jason Yosinski. An intriguing failing of convolutional neural networks and the coordconv solution. In *NeurIPS*, pages 9605–9616, 2018.

[43] Shichen Liu, Tianye Li, Weikai Chen, and Hao Li. Soft rasterizer: A differentiable renderer for image-based 3d reasoning. In *ICCV*, pages 7708–7717, 2019.

[44] Xinchen Liu, Wu Liu, Tao Mei, and Huadong Ma. A deep learning-based approach to progressive vehicle re-identification for urban surveillance. In *ECCV*, pages 869–884, 2016.

[45] Xiao Liu, Tian Xia, Jiang Wang, Yi Yang, Feng Zhou, and Yuanqing Lin. Fully convolutional attention networks for fine-grained recognition. *arXiv preprint arXiv:1603.06765*, 2016.

[46] Xiaobin Liu, Shiliang Zhang, Qingming Huang, and Wen Gao. Ram: a region-aware deep model for vehicle re-identification. In *ICME*, pages 1–6, 2018.

[47] Wei Luo, Xitong Yang, Xianjie Mo, Yuheng Lu, Larry S Davis, Jun Li, Jian Yang, and Ser-Nam Lim. Cross-x learning for fine-grained visual categorization. In *ICCV*, pages 8242–8251, 2019.

[48] Dechao Meng, Liang Li, Xuejing Liu, Yadong Li, Shijie Yang, Zheng-Jun Zha, Xingyu Gao, Shuhui Wang, and Qingming Huang. Parsing-based view-aware embedding network for vehicle re-identification. In *CVPR*, pages 7103–7112, 2020.

[49] KL Navaneet, Ansu Mathew, Shashank Kashyap, Wei-Chih Hung, Varun Jampani, and R Venkatesh Babu. From image collections to point clouds with self-supervised shape and pose networks. In *CVPR*, pages 1132–1140, 2020.

[50] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. In *NeurIPS Autodiff Workshop*, 2017.

[51] Ulrich Pinkall and Konrad Polthier. Computing discrete minimal surfaces and their conjugates. *Experimental mathematics*, 2(1):15–36, 1993.

[52] Adria Recasens, Petr Kellnhofer, Simon Stent, Wojciech Matusik, and Antonio Torralba. Learning to zoom: a saliency-based sampling layer for neural networks. In *ECCV*, pages 51–66, 2018.

[53] Zhixin Shu, Mihir Sahasrabudhe, Riza Alp Guler, Dimitris Samaras, Nikos Paragios, and Iasonas Kokkinos. Deforming autoencoders: Unsupervised disentangling of shape and appearance. In *ECCV*, pages 650–665, 2018.

[54] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.

[55] Ming Sun, Yuchen Yuan, Feng Zhou, and Errui Ding. Multi-attention multi-class constraint for fine-grained image recognition. In *ECCV*, pages 805–821, 2018.

[56] Maxim Tatarchenko, Stephan R Richter, René Ranftl, Zhuwen Li, Vladlen Koltun, and Thomas Brox. What do single-view 3d reconstruction networks learn? In *CVPR*, pages 3405–3414, 2019.

[57] James Thewlis, Hakan Bilen, and Andrea Vedaldi. Unsupervised learning of object frames by dense equivariant image labelling. In *NeurIPS*, pages 844–855, 2017.

[58] James Thewlis, Hakan Bilen, and Andrea Vedaldi. Modelling and unsupervised learning of symmetric deformable object categories. In *NeurIPS*, pages 8178–8189, 2018.

[59] Shubham Tulsiani, Tinghui Zhou, Alexei A Efros, and Jitendra Malik. Multi-view supervision for single-view reconstruction via differentiable ray consistency. In *CVPR*, pages 2626–2634, 2017.

[60] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, pages 5998–6008, 2017.

[61] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. In *California Institute of Technology*, 2011.

[62] Dequan Wang, Zhiqiang Shen, Jie Shao, Wei Zhang, Xiangyang Xue, and Zheng Zhang. Multiple granularity descriptors for fine-grained categorization. In *ICCV*, pages 2399–2406, 2015.

[63] Nanyang Wang, Yinda Zhang, Zhuwen Li, Yanwei Fu, Wei Liu, and Yu-Gang Jiang. Pixel2mesh: Generating 3d mesh models from single rgb images. In *ECCV*, pages 52–67, 2018.

[64] Yaming Wang, Vlad I Morariu, and Larry S Davis. Learning a discriminative filter bank within a cnn for fine-grained recognition. In *CVPR*, pages 4148–4157, 2018.

[65] Zhongdao Wang, Luming Tang, Xihui Liu, Zhuliang Yao, Shuai Yi, Jing Shao, Junjie Yan, Shengjin Wang, Hongsheng

Li, and Xiaogang Wang. Orientation invariant feature embedding and spatial temporal regularization for vehicle re-identification. In *ICCV*, pages 379–387, 2017.

[66] Chao Wen, Yinda Zhang, Zhuwen Li, and Yanwei Fu. Pixel2mesh++: Multi-view 3d mesh generation via deformation. In *ICCV*, pages 1042–1051, 2019.

[67] Shangzhe Wu, Christian Rupprecht, and Andrea Vedaldi. Unsupervised learning of probably symmetric deformable 3d objects from images in the wild. In *CVPR*, pages 1–10, 2020.

[68] 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*, pages 160–176, 2016.

[69] Yu Xiang, Roozbeh Mottaghi, and Silvio Savarese. Beyond pascal: A benchmark for 3d object detection in the wild. In *WACV*, pages 75–82, 2014.

[70] Ze Yang, Tiange Luo, Dong Wang, Zhiqiang Hu, Jun Gao, and Liwei Wang. Learning to navigate for fine-grained classification. In *ECCV*, pages 420–435, 2018.

[71] Han Zhang, Tao Xu, Mohamed Elhoseiny, Xiaolei Huang, Shaoting Zhang, Ahmed Elgammal, and Dimitris Metaxas. Spda-cnn: Unifying semantic part detection and abstraction for fine-grained recognition. In *CVPR*, pages 1143–1152, 2016.

[72] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, pages 586–595, 2018.

[73] Zhizheng Zhang, Cuiling Lan, Wenjun Zeng, and Zhibo Chen. Densely semantically aligned person re-identification. In *CVPR*, pages 667–676, 2019.

[74] Yifan Zhao, Ke Yan, Feiyue Huang, and Jia Li. Graph-based high-order relation discovery for fine-grained recognition. In *CVPR*, pages 15079–15088, 2021.

[75] Heliang Zheng, Jianlong Fu, Tao Mei, and Jiebo Luo. Learning multi-attention convolutional neural network for fine-grained image recognition. In *ICCV*, pages 5209–5217, 2017.

[76] Heliang Zheng, Jianlong Fu, Zheng-Jun Zha, and Jiebo Luo. Looking for the devil in the details: Learning trilinear attention sampling network for fine-grained image recognition. In *CVPR*, pages 5012–5021, 2019.

[77] Mohan Zhou, Yalong Bai, Wei Zhang, Tiejun Zhao, and Tao Mei. Look-into-object: Self-supervised structure modeling for object recognition. In *CVPR*, pages 11774–11783, 2020.

[78] Yi Zhou and Ling Shao. Viewpoint-aware attentive multi-view inference for vehicle re-identification. In *CVPR*, pages 6489–6498, 2018.