

# Low-Shot Learning from Imaginary 3D Model

Frederik Pahde<sup>1</sup>, Mihai Puscas<sup>1,2</sup>, Jannik Wolff<sup>1,3</sup>

Tassilo Klein<sup>1</sup>, Nicu Sebe<sup>2</sup>, Moin Nabi<sup>1</sup>

<sup>1</sup>SAP SE., Berlin, <sup>2</sup>University of Trento, <sup>3</sup>TU Berlin

{frederik.pahde, tassilo.klein, m.nabi}@sap.com

{mihaimarian.puscas, nicu.sebe}@unitn.it, jannik.wolff@campus.tu-berlin.de

## Abstract

Since the advent of deep learning, neural networks have demonstrated remarkable results in many visual recognition tasks, constantly pushing the limits. However, the state-of-the-art approaches are largely unsuitable in scarce data regimes. To address this shortcoming, this paper proposes employing a 3D model, which is derived from training images. Such a model can then be used to hallucinate novel viewpoints and poses for the scarce samples of the few-shot learning scenario. A self-paced learning approach allows for the selection of a diverse set of high-quality images, which facilitates the training of a classifier. The performance of the proposed approach is showcased on the fine-grained CUB-200-2011 dataset in a few-shot setting and significantly improves our baseline accuracy.

**Keywords:** Low-Shot Object Recognition, 3D Model, Mesh Reconstruction, 3D Shape Learning, Meta-Learning

## 1. Introduction

Since the successful introduction of deep learning techniques in countless computer vision applications, considerable research has been conducted to reduce the amount of annotated data needed for training such systems. Commonly, this data requirement problem has been approached systematically by developing algorithms which either require less expensive annotations such as semi-supervised or weakly supervised approaches, or more rigorously no annotations at all such as unsupervised systems. Although in theory quite appealing, the usual trade-off in these systems when applicable, is the overall reduced performance.

More importantly, there exist situations where the availability of annotated data is heavily skewed, reflecting the tail distribution found in the wild. In consequence, research in the domain of low-shot learning, i.e. learning and generalizing from only few training samples, has gained more and more interest (e.g. [23, 25, 29]). As such, generative approaches for artificially increasing the training set in low-

shot learning scenarios have been shown to be effective. Specifically, it was shown that with increasing quality and diversity of the generation output the overall performance of the low-shot learning system can be boosted [6, 18, 19, 20].

In this context, we propose to maximize the visual generative capabilities. Specifically, we assume a scenario where the base classes have a large amount of annotated data whereas the data for novel categories are scarce. To alleviate the data shortage we employ a high quality generation stage by learning a 3D structure [10] of the novel class. A curriculum-based discriminative sample selection method further refines the generated data, which promotes learning more explicit visual classifiers.

Learning the 3D structure of the novel class facilitates low-shot learning by allowing us to hallucinate images from different viewpoints of the same object. Simultaneously, learning the novel objects' texture map allows us for a controlled transfer of the novel objects' appearance to new poses seen in the base class samples. Freely hallucinating w.r.t. different poses and viewpoints of a single novel sample then in turn allows us to guarantee novel class data diversity. The framework by Kanazawa et al. [10] has proven to be very effective for learning both 3D models and texture maps without expensive 3D model annotations. While reconstructing a 3D model from single images in a given category has been achieved in the past [28, 11], these methods lack easy applicability to a hallucinatory setup and specifically miss any kind of texture and appearance reconstruction. The intuition behind our idea is visualized in Fig. 1

With a broad range of images generated for varying viewpoints and poses for the novel class, a selection algorithm is applied. To this end, we follow the notion of *self-paced learning* strategy, which is a general concept that has been applied in many other studies [15, 24]. It is related to curriculum learning [1], and is biologically inspired by the common human process of gradual learning, starting with the simplest concepts and increasing complexity. We employ this strategy to select a subset of images generated from the imaginary 3D model, which are associated with

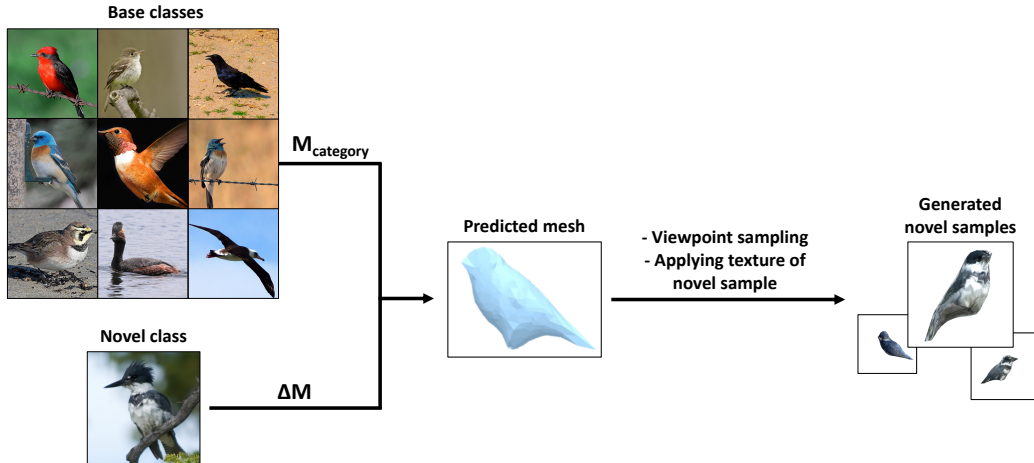


Figure 1: This figure illustrates one of our two generative methods, which is based on [10]: We first learn a generic mesh of the bird category. This mesh is then altered to fit the appearance of the target bird. We rotate the predicted 3D mesh to capture various viewpoints resulting in many 2D images that resemble the target bird. Those meshes are then coated with the novel bird’s texture. To cope with the varying quality, we subsequently apply a self-paced learning mechanism, which is elaborately outlined in figure 2 and in the remainder of the paper. For the second approach to sample generation, we exploit the pose variety of the base birds visible on the top left to enhance diversity. This approach is visualized in Figure 3.

high confidence w.r.t. “class discriminativeness” by the discriminator. Specifically the self-pacedness allows to handle the uncertainty related to the quality of generated samples. Here the notion of “easy” is interpreted as “high quality”. Training is then performed using only the subset consisting of images of sufficient quality. This set is then in turn progressively increased in the subsequent iterations when the model becomes more mature and is able to capture more complexity.

The main contributions of this work are: **First**, we massively expand the diversity of generating data from sparse samples of novel classes through learning 3D structure and texture maps. **Second**, we leverage a self-paced learning strategy facilitating reliable sample selection.

Our approach features robustness and outperforms the baseline in the challenging low-shot scenario.

## 2. Related Work

In this section we briefly review previous work considering: (1) low-shot learning, (2) 3D model learning and inference and (3) self-paced learning.

### 2.1. Low-Shot Learning

For learning deep networks using limited amounts of data, different approaches have been developed. Following Taigman et al. [27], Koch et al. [13] interpreted this task as a verification problem, i.e. given two samples, it has to be verified, whether both samples belong to the same class. Therefore, they employed siamese neural networks [4] to

compute the distance between the two samples and perform nearest neighbor classification in the learned embedding space. Some recent works approach few-shot learning by striving to avoid overfitting by modifications to the loss function or the regularization term. Yoo et al. [32] proposed a clustering of neurons on each layer of the network and calculated a single gradient for all members of a cluster during the training to prevent overfitting. The optimal number of clusters per layer is determined by a reinforcement learning algorithm. A more intuitive strategy is to approach few-shot learning on data-level, meaning that the performance of the model can be improved by collecting additional related data. Douze et al. [5] proposed a semi-supervised approach in which a large unlabeled dataset containing similar images was included in addition to the original training set. This large collection of images was exploited to support label propagation in the few-shot learning scenario. Hariharan et al. [6] combined both strategies (data-level and algorithm-level) by defining the squared gradient magnitude loss, that forces models to generalize well from only a few samples, on the one hand and generating new images by hallucinating features on the other hand. For the latter, they trained a model to find common transformations between existing images that can be applied to new images to generate new training data (see also [31]). Other recent approaches to few-shot learning have leveraged meta-learning strategies. Ravi et al. [23] trained a long short-term memory (LSTM) network as meta-learner that learns the exact optimization algorithm to train a learner neural network that performs the classification in a few-shot learning setting. This method

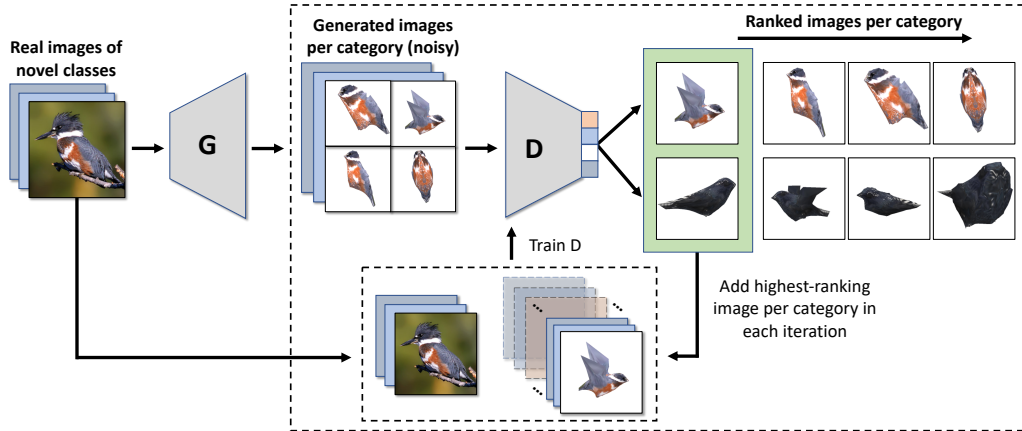


Figure 2: Self-paced fine-tuning on novel classes: For each novel class, noisy samples are generated with different viewpoints and poses by  $G$ . Those images are ranked by  $D$  based on their class-discriminatory power. The highest-ranking images are added to the novel samples and used to update  $D$ , which is trained using a simple cross-entropy loss. This process is repeated multiple times. Initially,  $D$  has been pre-trained on all base class data.

was proposed due to the observation that the update function of standard optimization algorithms like SGD is similar to the update of the cell state of a LSTM. Bertinetto et al. [2] trained a meta-learner feed-forward neural network that predicts the parameters of another, discriminative feed-forward neural network in a few-shot learning scenario. Another tool that has been applied successfully to few-shot learning recently is attention. Vinyals et al. [29] introduced matching networks for one-shot learning tasks. This network is able to apply an attention mechanism over embeddings of labeled samples in order to classify unlabeled samples. One further outcome of this work is that it is helpful to mimic the one-shot learning setting already during training by defining mini-batches, called episodes with subsampled classes. Snell et al. [25] generalize this approach by proposing prototypical networks. Prototypical networks search for a non-linear embedding space (the prototype) in which classes can be represented as the mean of all corresponding samples. Classification is then performed by finding the closest prototype in the embedding space. In the one-shot scenario, prototypical networks and matching networks are equivalent.

## 2.2. 3D Shape Learning

Inferring the 3D shape of an object from differing viewpoints has long been a topic of interest in computer vision. Based on the idea that there exists a categorical-specific canonical shape, and that class-specific deformations of it can be learned, systems such as SMPL [17] and "Keep it SMPL" [3] model a human 3D shape space, while Zuffi et al. [35] perform a similar task for quadruped animals. However, even though these methods are able to use synthetic training data, they still rely on a 3D shape ground truth. In

contrast, Kanazawa et al. [10] make use of much cheaper keypoint and segmentation mask annotations, which allows both 3D mesh and texture inference for images.

## 2.3. Self-Paced Learning

Recently, many studies have shown the benefits of organizing the training examples in a meaningful order (e.g., from simple to complex) for model training. Bengio et al. [1] first proposed a general learning strategy: curriculum learning. They show that suitably sorting the training samples, from the easiest to the most difficult, and iteratively training a classifier starting with a subset of easy samples (which is progressively augmented with more and more difficult samples), can be useful to find better local minima. Note that in this and in all the other curriculum-learning-based approaches, the order of the samples is provided by an external supervisory signal, taking into account human domain-specific expertise.

Curriculum learning was extended to self-paced learning by Kumar et al. [15]. They proposed the respective framework, automatically expanding the training pool in an easy-to-hard manner by converting the curriculum mechanism into a concise regularization term. Curriculum learning uses human design to organize the examples, and self-paced learning can automatically choose training examples according to the loss. Supancic et al. [26] adopt a similar framework in a tracking scenario and train a detector using a subset of video frames, showing that this selection is important to avoid drifting. Jiang et al. [9] pre-cluster the training data in order to balance the selection of the easiest samples with a sufficient inter-cluster diversity. Pentina et al. [22] propose a method in which a set of learning tasks is automatically sorted in order to allow a gradual sharing of infor-

mation among tasks. In Zhang et al.’s [33] model saliency is used to progressively select samples in weakly supervised object detection. In context of visual categorization some of these self-paced learning methods use CNN-based features to represent samples [16] or use a CNN as the classifier directly [24].

### 3. Method

#### 3.1. Preliminaries

In this subsection we introduce the necessary notation. Let  $\mathcal{I}$  denote the image space,  $\mathcal{T}$  the texture space,  $\mathcal{M}$  the 3D mesh space and  $\mathcal{C} = \{1, \dots, L\}$  the discrete label space. Further, let  $x_i \in \mathcal{I}$  be the  $i$ -th input data point, and  $y_i \in \mathcal{C}$  its label. In the low-shot setting, we consider two subsets of the label space:  $\mathcal{C}_{\text{base}}$  for labels for which we have access to a large number of samples, and the novel classes  $\mathcal{C}_{\text{novel}}$ , which are underrepresented in the data. Note that both subsets exhaust the label space  $\mathcal{C}$ , i.e.  $\mathcal{C} = \mathcal{C}_{\text{base}} \cup \mathcal{C}_{\text{novel}}$ . We further assume that in general  $|\mathcal{C}_{\text{novel}}| \ll |\mathcal{C}_{\text{base}}|$ . The dataset  $\mathcal{S}$  decomposes as follows:  $\mathcal{S} = \mathcal{S}_{\text{train}} \cup \mathcal{S}_{\text{test}}$ ,  $\mathcal{S}_{\text{train}} \cap \mathcal{S}_{\text{test}} = \emptyset$ . The training data  $\mathcal{S}_{\text{train}}$  consists of 2-tuples  $\{(x_i, y_i)\}_{i=1}^N$  taken from the whole data set containing both image samples and labels. Furthermore, for 3D model prediction we also attach 3-tuples  $\{(l_i, k_i, m_i)\}_{i=1}^N$ , with  $l_i$  being a foreground object segmentation mask and  $k_i$  a 15-point keypoint vector representing the pose of the object. Additionally,  $m_i$  denotes the weak-perspective camera, which is estimated by leveraging structure-from-motion on the training instances’ keypoints  $k_i$ . The test data is drawn from the novel classes and does not contain any 3D information, but solely images and their labels. Next, there is also  $\mathcal{S}_{\text{train}}^{\text{novel}} = \{(x_i, y_i, l_i, k_i, m_i) : (x_i, y_i, l_i, k_i, m_i) \in \mathcal{S}_{\text{train}}, y_i \in \mathcal{C}_{\text{novel}}\}_{i=1}^M \subset \mathcal{S}_{\text{train}}$ , which denotes the training data for the novel categories. For each class in  $\mathcal{C}_{\text{novel}}$ ,  $k$  samples can be used for training (k-shot), resulting in  $|\mathcal{S}_{\text{train}}^{\text{novel}}| \ll |\mathcal{S}_{\text{train}}|$ .

#### 3.2. 3D Model Based Data Generation

The underlying observation on which our method is based on is that increased diversity of generated images directly translates into higher classification performance for novel categories. The proposed work aims at emulating processes in human cognition that allow for reconstructing different viewpoints and poses through conceptualizing a 3D model of an object of interest. Specifically, we aim to learn such a 3D representation for novel samples appearing during training and leverage it to predict different viewpoints and poses of that object.

We use the architecture proposed by Kanazawa et al. [10] to predict a 3D mesh  $M_i$  and texture  $T_i$  from an image sample  $x_i$ . With the assumption that all  $x_i \in \mathcal{I}$  represent objects of the same category, the shape of each instance is predicted

---

**Algorithm 1** Self-paced learning,  $\text{RANK}()$  is a function that ranks generated images based on their score of  $D'$  and  $\text{TOP}()$  returns the highest ranked images

---

```

1: Input: Pre-trained network  $D, \mathcal{S}_{\text{gen}}^{\text{novel}}, r$ 
2: Output: Fine-tuned classifier  $D'$ 
3: for  $i = 1, \dots, n$  do
4:    $\mathcal{S}_{\text{all}}^{\text{novel}} = \emptyset$ 
5:   for  $c \in \mathcal{C}_{\text{novel}}$  do
6:     candidates =  $\emptyset$ 
7:     for  $x_i^{\text{gen}} \in \mathcal{S}_{\text{gen}}^{\text{novel}}$  do
8:       candidates = candidates  $\cup x_i^{\text{gen}}$ 
9:     candidatesranked =  $\text{RANK}(\text{candidates}, D')$ 
10:    sample =  $\text{TOP}(\text{candidates}_{\text{ranked}}, r)$ 
11:     $\mathcal{S}_{\text{gen}}^{\text{novel}} = \mathcal{S}_{\text{gen}}^{\text{novel}} \cup \text{sample}$ 
12:   $\mathcal{S}_{\text{all}}^{\text{novel}} = \mathcal{S}_{\text{train}}^{\text{novel}} \cup \mathcal{S}_{\text{gen}}^{\text{novel}}$ 
13:  update  $D'$  with  $\mathcal{S}_{\text{all}}^{\text{novel}}$ 

```

---

by deforming a learned category-specific mesh  $M_{\text{cat}}$ . Note that *category* refers to the entire fine-grained bird dataset, as opposed to *class*. All recovered shapes will share a common underlying 3D mesh structure,  $M_i = M_{\text{cat}} + \Delta M_i$ , with  $\Delta M_i$  being the predicted mesh deformation for instance  $x_i$ . Because the mesh  $M$  has the same vertex connectivity as the average categorical mesh  $M_{\text{cat}}$ , and further as  $M_{\text{sphere}}$  representing a sphere, a predicted texture map  $T_i$  can be easily applied over any generated mesh.

An advantage of [10] over related methods is that learning the 3D representation does not require expensive 3D model or multi-view annotations.

Given  $(M_i, T_i, \Theta_i)$  and  $\Theta = (\alpha, \beta, \gamma)$ , where the three camera rotation angles  $\alpha, \beta, \gamma$  are sampled uniformly from  $[0, \pi/6]$ , we can project the reconstructed object using  $f_{\text{gen}}(M_i, T_i, \Theta_i)$  such that  $X_i^{\text{view}} = \{x_i^0, \dots, x_i^L\}$  contains samples of the object seen from different viewpoints.

As  $X_i^{\text{view}}$  only contains different viewpoints of the novel object, it will not contain any novel poses. This is a concern for non-rigid object categories, where it cannot be guaranteed that the unseen samples in a novel class will have similar poses to the known samples in the novel class. To mitigate this, the diversity of the generated data must be expanded to include new object poses.

All meshes predicted from  $x_j \in \mathcal{S}_{\text{base}}$  obtain the spherical texture map  $T_i$  corresponding to  $x_i \in \mathcal{S}_{\text{train}}^{\text{novel}}$  using  $f_{\text{gen}}(M_j, T_i, \Theta_j)$ . This transfers the shape from base class objects to novel class instances resulting in  $X_i^{\text{pose}} = \{x_i^J, \dots, x_i^S\}$ .

Using poses from images of different labels is an inherently noisy approach through inter-class mesh variance. However, a subsequent sample selection strategy allows the algorithm to make use of the most representative poses. Indeed, as seen in Figure 3, meshes  $M_j \in \mathcal{S}_{\text{base}}$  exist for

which the predicted images  $x_i^j$  are visually similar to samples of the unseen classes.

Thus, for each sample  $x_i \in S_{train}^{novel}$ , a set of images  $S_{gen}^{novel} = X_i^{view} \cup X_i^{pose}$  is generated. This generated data captures both different viewpoints of the novel class and the appearance of the novel class applied to differing poses from the base classes.

### 3.3. Pre-Training of Classifier

In the low-shot learning framework proposed by Hariharan and Girshick [7], a representation of the base categorical data must be learned beforehand. This is achieved by learning a classifier on the samples available in the base classes, i.e.  $x_i \in S_{train}^{base}$ . For this task we make use of an architecture identical to the StackGAN discriminator [34], modified to serve as a classifier. This discriminator  $D$  is learned on  $S_{train}^{base}$  by minimizing  $L_{class}$  defined as a cross-entropy loss.

However, to accommodate for the different amount of classes in base and novel,  $D$  has to be adapted. Specifically, the class-aware layer with  $|C_{base}|$  output neurons is replaced and reduced to  $|C_{novel}|$  output neurons, which are randomly initialized. We refer to this adapted classifier as  $D'$ . Subsequently, the network can be fine-tuned using the available novel class data.

### 3.4. Self-Paced Learning

As seen in section 3.2, for a given novel sample  $x_i \in S_{train}^{novel}$  we can generate  $S_{gen}^{novel} = X_i^{view} \cup X_i^{pose}$ , containing new viewpoints and poses of the given object.

For the self-paced learning stage, we fine-tune with the novel samples, as well as the samples generated through projecting the predicted 3D mesh and texture maps. i.e. with the data given by  $S_{train}^{novel} \cup S_{gen}^{novel}$ .

Unfortunately, the samples contained in  $S_{gen}^{novel}$  can be noisy for a variety of reasons: failure in predicting the 3D mesh deformation due to a too large difference between the categorical mesh and the object mesh, or even viewpoints that are not representative to the novel class.

To mitigate this we propose a self-paced learning strategy ensuring that only the best generated samples within  $S_{gen}^{novel}$  are used.

Again taking into account the setting of low-shot learning, we restrict the number of samples per class available to  $k$ . Due to the limited amount of samples, the initialized  $D'$  will be weak on the classification task, but sufficiently powerful for performing an initial ranking of the generated images. For this task we employ the softmax activation for class-specific confidence scoring. As  $D'$  learns to generalize better, more difficult samples will be selected.

This entails iteratively choosing generated images that have highest probability in  $D'$  for  $C_{novel}$ , yielding a curated set of generated samples  $S_{gen}^{novel}$ . An issue in selecting the highest scoring sample in each iteration is the possibility of

not making full use of the available data w.r.t. its diversity - the highest scoring images being of a very similar pose and viewpoint to the original sample.

We address this shortcoming by using a clustering-and-discard strategy: For the novel class training sample  $x_i$ , we generate  $X_i^{gen} = \{x_i^0, \dots, x_i^{L+S}\}$  new images, representing new viewpoints and poses of the object.  $X_i^{gen}$  is then further associated with  $K_i^{gen} = \{k_i^0, \dots, k_i^Q\}$ , representing all the predicted keypoints of the associated generated samples.  $K_i^{gen}$  is clustered using a simple k-means implementation [21]. On every self-paced iteration, the pose cluster associated to the selected top-ranked sample is discarded to increase data diversity.

Finally, we aggregate original samples and generated images  $S_{train}^{novel} \cup S_{gen}^{novel}$  for training, during which we update  $D'$ . Doing so yields both a more accurate ranking as well as higher class prediction accuracy as the number of samples increases. Ultimately, the approach learns a reliable classifier that performs well in low-shot learning scenarios. It is summarized in algorithm 1.

## 4. Experiments

### 4.1. Datasets

We test the applicability of our method on CUB-200-2011 [30], a fine-grained classification datasets containing 11,788 images of 200 different bird species of size  $\mathcal{I} \subset \mathbb{R}^{256 \times 256}$ . The data is split equally into training and test data. As a consequence, samples are roughly equally distributed, with training and test each containing  $\approx 30$  images per class. Additionally, foreground masks, semantic keypoints and angle predictions are provided by [10]. Note that nearly 300 images are removed where the number of visible keypoints is less or equal than 6.

Following Zhang et al. [34], we split the data such that  $|C_{base}| = 150$  and  $|C_{novel}| = 50$ . To simulate low-shot learning,  $k \in \{1, 2, 5, 10, 20\}$  images of  $C_{novel}$  are used for training, as proposed by [6].

### 4.2. Algorithmic Details

During representation learning, we train an initial classifier on the base classes for 600 epochs and use Adam [12] for optimization. We set the learning rate  $\tau$  to  $10^{-3}$  and the batch size for  $D$  to 32. In the initialization phase for self-paced learning, we construct  $D'$  by replacing the last layer of  $D$  by a linear softmax layer of size  $|C_{novel}|$ . The resulting network is then optimized using the cross-entropy loss function and an Adam optimizer with the same parameters. Batch size is set to 32 and training proceeds for 20 epochs. Self-paced learning of  $D'$  continues to use the same settings, i.e. the Adam optimizer minimizing a cross-entropy loss. In every iteration we choose exactly one generated image per class and perform training for 10 epochs.

Model	k				
	1	2	5	10	20
Baseline	27.55	30.75	54.25	58.51	71.62
Views + poses	33.40	43.72	54.81	65.27	74.06
SPL w/ views	33.54	41.49	54.88	65.48	<b>74.97</b>
SPL w/ poses	33.82	42.47	54.95	64.85	73.64
SPL w/ poses + clustering	33.40	45.05	57.74	65.69	74.62
SPL w/ poses + views	35.29	41.98	55.37	66.04	71.48
SPL w/ poses + views (balanced)	35.77	44.56	54.60	64.30	74.83
SPL w/ all	<b>36.96</b>	<b>45.40</b>	<b>58.09</b>	<b>66.53</b>	74.83

Table 1: Ablation study of our model in a top-5, 50-way scenario on the CUB-200-2011 dataset in different k-shot settings, best results are in bolt. We observe that each of the proposed extensions increases the accuracy in at least one setting which justifies their usage. This regards to both, methods for generating additional data and the approach to only select generated samples of sufficient quality for training the classifier.

### 4.3. Models

In order to assess the performance of individual components, we perform an ablation study.

The simplest transfer learning approach is making use of a pre-trained representation and then fine-tuning that model on the novel data. A first baseline (**Baseline**) uses this strategy: we pre-train a classifier  $D$  on the base classes, following by fine-tuning with  $k$  novel class instances  $x_i \in \mathcal{S}_{\text{train}}^{\text{novel}}$ . This strategy makes use of the fine-grained character of the dataset, learning initial representations on  $\mathcal{C}_{\text{base}}$  and performing classification on  $\mathcal{C}_{\text{novel}}$ .

A second model **views + poses** studies the validity of the generated viewpoint and pose data. For  $r$  sampling iterations, a single uniformly sampled  $x_i \in \mathcal{S}_{\text{gen}}^{\text{novel}}$  is attached to a novel sample set.

We then introduce sample selection to our method. Note that viewpoint generation is achieved through 3D Mesh  $M_i$  and texture  $T_i$  of the same sample  $x_i$ , while the different poses are generated through applying the novel class instance texture  $T_i$  to base class meshes  $M_j$ . The **SPL w/ views** and **SPL w/ poses** sample the generated data from the generated viewpoints  $X^{\text{view}}$  and  $X^{\text{pose}}$  respectively.

**SPL w/ poses + views** makes use of the entirety of  $\mathcal{S}_{\text{gen}}^{\text{novel}}$ , while **SPL w/ poses + views (balanced)** tackles the data imbalance between different viewpoint samples and different pose samples by ranking the two branches separately, and selecting one sample from each such that for one novel sample,  $x_i^{\text{max,pose}}$  and  $x_i^{\text{max,view}}$  are used in fine-tuning.

The clustering-and-dismissal mechanism detailed in 3.4 is evaluated in the **SPL w/ poses + clustering** model, while **SPL w/ all** makes use of the method in its entirety.

### 4.4. Results of Ablation Study

The results of the ablation study outlined in the previous section are shown in Table 1, presenting 50-way, top-5 accuracies for k-shot learning with  $k \in \{1, 2, 5, 10, 20\}$ .

We first evaluate the baseline model, which is trained on the base classes and fine-tuned on the novel classes. Due to using a relatively shallow classification network, and the sparsity of the novel samples, the network rapidly overfits.

Introducing more data diversity to the fine-tuning stage through 3D model inference provides a significant boost in performance in all  $k \in \{1, 2, 5, 10, 20\}$ . With the generated samples selected randomly, the network does not easily overfit, but this selection method provides no protection against noisy generated samples.

Subsequent models evaluate different selection strategies across the two defined generated data splits for new viewpoints and poses, i.e.  $X^{\text{view}}$  and  $X^{\text{pose}}$ . The contribution of the self paced learning strategy can be evaluated directly comparing the top-5 accuracies of the **view + poses** model and the **SPL w/ views + poses** model. The increase of performance when k is small shows that the selection strategy can achieve better performance, but inconsistently across different  $k$  values.

One cause of this problem is how the generated data is split, and whether the classifier has access to the most valuable generated samples. In **SPL w/ poses** and **SPL w/ views**, we only select samples from  $X^{\text{pose}}$  and  $X^{\text{view}}$  respectively. The experimental results of both models are similar and inferior to **SPL w/ views + poses**, where both sets are used. Even with higher performance, the aggregate model selects from  $X^{\text{view}}$  almost exclusively, hinting on a type of mode collapse.

To further diversify the possible data picks, we "balance" the two sets: For each sample,  $x_i^{\text{max,pose}}$  and  $x_i^{\text{max,view}}$  are selected as the highest scoring samples in their respective



sets. This disentangling of pose and viewpoint data offers an across-the-board improvement, as seen in **SPL w/ views + poses (balanced)**.

While normally each sample that was selected in self-paced iteration  $r$  is discarded, this will likely leave a number of samples that are similar in pose, such that the classifier may rank them as maximum. This does not add significant new information to the learning process, and as such the clustering-by-pose method guiding the sample dismissal is introduced. Indeed, as observed in **SPL w/ all**, both the sample-discard strategy, and the balancing strategy are similar useful for selections in self-paced learning. With all discussed techniques introduced, the model achieves a significant performance boost compared to the baseline.

#### 4.5. Analysis of Self-Paced Fine-Tuning

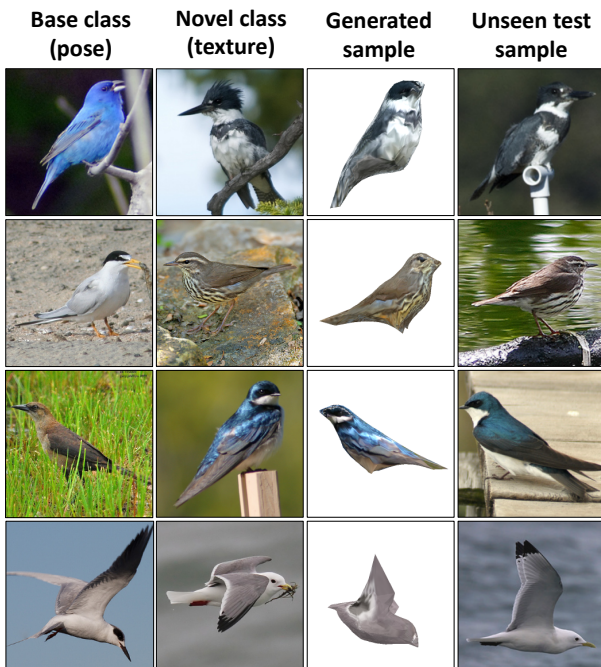


Figure 3: Texture from novel class birds is transferred onto poses from base class birds. The generated samples have been previously selected by the discriminator w.r.t. to their class-discriminatory power in the self-paced learning setting. Those hallucinations are visually similar to unseen test samples, indicating their value for training a classifier.

We run several additional experiments to further analyze the behavior of our method. For the those experiments we use the CUB-200-2011 bird dataset, and compare to the method by Hariharan and Girshick [6] in Table 2.

We first report the baseline model in the top-1, 1-shot scenario. Due to the relative shallowness of the classification network and without any sample selection or hallucina-

Baseline	NN	Our (shallow)	Our (ResNet)	[6]
9.1	9.7	14.4	18.5	19.1

Table 2: Top-1, 50-way, 1-shot accuracies on the CUB-200-2011 dataset. We see that our shallow CNN (trained with self-paced learning) exceeds both baselines. The ResNet (not trained with self-paced learning) is within reach of Hariharan and Girshick’s model with SGM loss [6], for which we have reproduced respective results.

tion, the performance is quite low.

Methods using simple nearest neighbour classifiers can perform well on few-shot learning tasks [14]. We implement a simple nearest-neighbour classifier using the representations learned in our baseline on the base class samples,  $x_i \in S_{train}^{base}$ , specifically making use of the last hidden layer of the network. This model marginally outperforms the baseline.

Improving the novel class data diversity by using self-paced sample selection and k-means clustering-and-dismissal, the performance rises by 5.3 points to 14.4, which equals more than 50% relative improvement.

So far, we have used a classifier with simple architecture and loss function in order to present the most general possible framework and to allow for a fair comparison with baseline methods. However, we expect a significant boost in accuracy using larger classifiers. To test this hypothesis, we fine-tune a modified ResNet-18 [8]. We first reduce the output dimensionality of the last pooling layer from 512 to 256 by lowering the amount of filters. After having trained this model on the base classes, we replace the last, fully-connected layer of size  $|C_{base}|$  with a smaller one of size  $|C_{novel}|$  to account for the different amount of classes. Afterwards, we freeze all layers except the final one, and train with  $S_{train}^{novel} \cup S_{gen}^{novel}$  after having ranked the existing samples with the best shallow network. We observe comparable results to Hariharan and Girshick [6] despite of neither having used the ResNet-18 as a ranking function for self-paced learning, nor performing iterative sampling. Note that our method provides a general framework to augment the training set with class-discriminative generated samples that can potentially be used in conjunction with more sophisticated methods as the SGM loss [6] to obtain better results.

## 5. Conclusion and Future Work

In this paper, we proposed to extend few-shot learning by incorporating image hallucination from 3D models in conjunction with a self-paced learning strategy. Experiments on the CUB dataset demonstrate that learning generative methods employing 3D models reaches performance that significantly outperforms our baseline and is competitive to

popular methods in the field. Thus the proposed approach allows for an efficient compensation of the lack of data in novel categories.

For future work we plan to optimize the pipeline in an end-to-end fashion, discarding the self-paced learning sample selection and replacing it with learnable viewpoint angle parameters.

## References

- [1] Y. Bengio, J. Louradour, R. Collobert, and J. Weston. Curriculum learning. In *ICML*, pages 41–48, 2009.
- [2] L. Bertinetto, J. F. Henriques, J. Valmadre, P. Torr, and A. Vedaldi. Learning feed-forward one-shot learners. In *Advances in Neural Information Processing Systems*, pages 523–531, 2016.
- [3] F. Bogo, A. Kanazawa, C. Lassner, P. Gehler, J. Romero, and M. J. Black. Keep it smpl: Automatic estimation of 3d human pose and shape from a single image. In *European Conference on Computer Vision*, pages 561–578. Springer, 2016.
- [4] J. Bromley, I. Guyon, Y. LeCun, E. Säckinger, and R. Shah. Signature verification using a” siamese” time delay neural network. In *Advances in Neural Information Processing Systems*, pages 737–744, 1994.
- [5] M. Douze, A. Szlam, B. Hariharan, and H. Jégou. Low-shot learning with large-scale diffusion. *CoRR*, 2017.
- [6] B. Hariharan and R. Girshick. Low-shot Visual Recognition by Shrinking and Hallucinating Features. In *ICCV*, 2017.
- [7] B. Hariharan and R. B. Girshick. Low-shot visual object recognition. *CoRR*, abs/1606.02819, 2016.
- [8] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [9] L. Jiang, D. Meng, S.-I. Yu, Z. Lan, S. Shan, and A. Hauptmann. Self-paced learning with diversity. In *Advances in Neural Information Processing Systems*, pages 2078–2086, 2014.
- [10] A. Kanazawa, S. Tulsiani, A. A. Efros, and J. Malik. Learning category-specific mesh reconstruction from image collections. In *ECCV*, 2018.
- [11] A. Kar, S. Tulsiani, J. Carreira, and J. Malik. Category-specific object reconstruction from a single image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1966–1974, 2015.
- [12] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [13] G. Koch, R. Zemel, and R. Salakhutdinov. Siamese neural networks for one-shot image recognition. In *ICML Deep Learning Workshop*, volume 2, 2015.
- [14] R. G. Krishnan, A. Khanelwal, R. Ranganath, and D. Songtao. Max-margin learning with the bayes factor. In *Proceedings of the Conference on Uncertainty in Artificial Intelligence (UAI)*, 2018.
- [15] M. P. Kumar, B. Packer, and D. Koller. Self-paced learning for latent variable models. In *NIPS*, pages 1189–1197, 2010.
- [16] X. Liang, S. Liu, Y. Wei, L. Liu, L. Lin, and S. Yan. Towards computational baby learning: A weakly-supervised approach for object detection. In *ICCV*, pages 999–1007, 2015.
- [17] M. Loper, N. Mahmood, J. Romero, G. Pons-Moll, and M. J. Black. Smpl: A skinned multi-person linear model. *ACM Transactions on Graphics (TOG)*, 34(6):248, 2015.
- [18] G. H. Navaneeth Bodla and R. Chellappa. Semi-supervised fusedgan for conditional image generation. *arXiv preprint*.
- [19] A. Odena, C. Olah, and J. Shlens. Conditional image synthesis with auxiliary classifier gans. *arXiv preprint arXiv:1610.09585*, 2016.
- [20] F. Pahde, M. Nabi, T. Klein, and P. Jahnichen. Discriminative hallucination for multi-modal few-shot learning. In *2018 25th IEEE International Conference on Image Processing (ICIP)*, pages 156–160. IEEE, 2018.
- [21] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [22] A. Pentina, V. Sharmanska, and C. H. Lampert. Curriculum learning of multiple tasks. In *CVPR*, pages 5492–5500, 2015.
- [23] S. Ravi and H. Larochelle. Optimization as a model for few-shot learning. In *International Conference on Learning Representations*, 2017.
- [24] E. Sangineto, M. Nabi, D. Culibrk, and N. Sebe. Self paced deep learning for weakly supervised object detection. *arXiv preprint arXiv:1605.07651*, 2016.
- [25] J. Snell, K. Swersky, and R. Zemel. Prototypical networks for few-shot learning. In *NIPS*, pages 4080–4090, 2017.
- [26] J. S. Supancic III and D. Ramanan. Self-paced learning for long-term tracking. In *CVPR*, pages 2379–2386, 2013.
- [27] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf. Deepface: Closing the gap to human-level performance in face verification. In *CVPR*, pages 1701–1708, 2014.
- [28] S. Vicente, J. Carreira, L. Agapito, and J. Batista. Reconstructing pascal voc. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 41–48, 2014.
- [29] O. Vinyals, C. Blundell, T. Lillicrap, D. Wierstra, et al. Matching networks for one shot learning. In *NIPS*, pages 3630–3638, 2016.
- [30] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The Caltech-UCSD Birds-200-2011 Dataset. Technical report, 2011.
- [31] Y.-X. Wang, R. Girshick, M. Hebert, and B. Hariharan. Low-Shot Learning from Imaginary Data. In *CVPR*, 2018.
- [32] D. Yoo, H. Fan, V. N. Boddeti, and K. M. Kitani. Efficient K-Shot Learning with Regularized Deep Networks. In *AAAI*, 2018.
- [33] D. Zhang, D. Meng, L. Zhao, and J. Han. Bridging saliency detection to weakly supervised object detection based on self-paced curriculum learning. *arXiv preprint arXiv:1703.01290*, 2017.
- [34] H. Zhang, T. Xu, H. Li, S. Zhang, X. Wang, X. Huang, and D. Metaxas. Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. In *ICCV*, 2017.
- [35] S. Zuffi, A. Kanazawa, D. W. Jacobs, and M. J. Black. 3d menagerie: Modeling the 3d shape and pose of animals.