






iNeMo: Incremental Neural Mesh Models for Robust Class-Incremental Learning

Tom Fischer¹ , Yaoyao Liu² , Artur Jesslen³, Noor Ahmed¹ ,
Prakhar Kaushik² , Angtian Wang² , Alan Yuille²,
Adam Kortylewski^{3,4}, and Eddy Ilg¹

¹ Saarland University, Saarbrücken, Germany

{fischer, ilg}@cs.uni-saarland.de

² Johns Hopkins University, Baltimore, USA

{yliu538, angtianwang, ayuille1}@jhu.edu

³ University of Freiburg, Freiburg, Germany

{jesslen, kortylew}@cs.uni-freiburg.de

⁴ Max-Planck Institute for Informatics, Saarbrücken, Germany

akortyle@mpi-inf.mpg.de

Abstract. Different from human nature, it is still common practice today for vision tasks to train deep learning models only initially and on fixed datasets. A variety of approaches have recently addressed handling continual data streams. However, extending these methods to manage out-of-distribution (OOD) scenarios has not effectively been investigated. On the other hand, it has recently been shown that non-continual neural mesh models exhibit strong performance in generalizing to such OOD scenarios. To leverage this decisive property in a continual learning setting, we propose incremental neural mesh models that can be extended with new meshes over time. In addition, we present a latent space initialization strategy that enables us to allocate feature space for future unseen classes in advance and a positional regularization term that forces the features of the different classes to consistently stay in respective latent space regions. We demonstrate the effectiveness of our method through extensive experiments on the Pascal3D and ObjectNet3D datasets and show that our approach outperforms the baselines for classification by 2 – 6% in the in-domain and by 6 – 50% in the OOD setting. Our work also presents the first incremental learning approach for pose estimation. Our code and model can be found at github.com/Fischer-Tom/iNeMo.

Keywords: Class-incremental learning · 3D pose estimation

1 Introduction

Humans inherently learn in an incremental manner, acquiring new concepts over time, with little to no forgetting of previous ones. In contrast, trying to mimic the same behavior with machine learning suffers from *catastrophic forgetting* [22, 37, 38], where learning from a continual stream of data can destroy the knowledge that was previously acquired. In this context, the problem was formalized as *class-incremental learning* and a variety of approaches

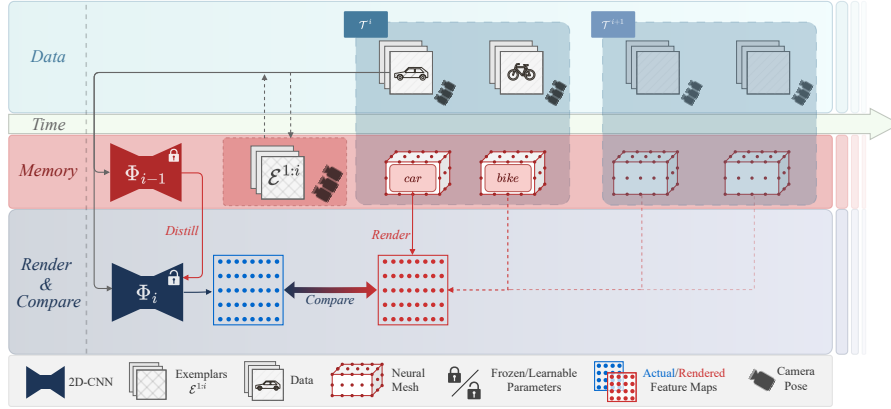


Fig. 1: We present iNeMo that can perform class-incremental learning for pose estimation and classification, and performs well in out-of-distribution scenarios. Our method receives tasks \mathcal{T}^i over time that consist of images with camera poses for new classes. We build up on Neural Mesh Models (NeMo) [53] and abstract objects with simple cuboid 3D meshes, where each vertex carries a neural feature. The neural meshes are optimized together with a 2D feature extractor Φ_i and render-and-compare can then be used to perform pose estimation and classification. We introduce a memory that contains an old feature extractor Φ_{i-1} for distillation, a replay buffer $\mathcal{E}^{1:(i-1)}$ and a growing set of neural meshes \mathcal{M} . Our results show that iNeMo outperforms all baselines for incremental learning and is significantly more robust than previous methods.

have been proposed to address catastrophic forgetting for models that work in-distribution [12, 18, 27, 31, 32, 46]. However, extending these methods to effectively manage out-of-distribution (OOD) scenarios [65] to the best of our knowledge has not been investigated.

Neural mesh models [53] embed 3D object representations explicitly into neural network architectures, and exhibit strong performance in generalizing to such OOD scenarios for classification and 3D pose estimation. However, as they consist of a 2D feature extractor paired with a generative model, their extension to a continual setting with existing techniques is not straight forward. If one would only apply those techniques to the feature extractor, the previously learned neural meshes would become inconsistent and the performance of the model would drop.

In this paper, we therefore present a strategy to learn neural mesh models incrementally and refer to them as incremental Neural Mesh Models (iNeMo). As shown in Figure 1, in addition to the conventional techniques of knowledge distillation and maintaining a replay buffer, our approach introduces a memory that contains a continuously growing set of meshes that represent object categories. To establish the learning of the meshes in an incremental setting, we extend the contrastive learning from [53] by a latent space initialization strategy that enables us to allocate feature space for future unseen classes in advance, and

a positional regularization term that forces the features of the different classes to consistently stay in respective latent space regions. Through extensive evaluations on the Pascal3D [62] and ObjectNet3D [61] datasets, we demonstrate that our method outperforms existing continual learning techniques and furthermore surpasses them by a large margin for out-of-distribution samples. Overall, our work motivates future research on joint 3D object-centric representations. In summary, the contributions of our work are:

1. For the first time, we adapt the conventional continual learning techniques of knowledge distillation and replay to the 3D neural mesh setting.
2. We propose a novel architecture, that can grow by adding new meshes for object categories over time.
3. To effectively train the features of the meshes, we introduce a strategy to partition the latent space and maintain it when new tasks are integrated.
4. We demonstrate that incremental neural mesh models can outperform 2D baselines that use existing 2D continual learning techniques by 2 – 6% in the in-domain and by 6 – 50% in the OOD setting.
5. Finally, we introduce the first incremental approach for pose estimation and show that the neural mesh models outperform 2D baselines.

2 Related Work

2.1 Robust Image Classification and Pose Estimation

Image Classification has always been a cornerstone of computer vision. Ground-breaking models such as ResNets [14], Transformers [52], and Swin Transformers [33] have been specifically designed for this task. However, these models predominantly target the in-distribution setting, leading to a significant gap in performance when faced with challenging benchmarks that involve synthetic corruptions [15], occlusions [56], and out-of-distribution (OOD) images [65]. Attempts to close this performance gap have included data augmentation [16] and innovative architectural designs, such as the analysis-by-synthesis approach [23]. Along this line of research, recently neural mesh models emerged as a family of models [36, 53–55] that learn a 3D pose-conditioned model of neural features and predict 3D pose and object class [20] by minimizing the reconstruction error between the actual and rendered feature maps using render-and-compare. Such models have shown to be significantly more robust to occlusions and OOD data. However, they can so far only be trained on fixed datasets. In this work, we present the first approach to learn them in a class-incremental setting.

Object Pose Estimation has been approached primarily as a regression problem [40, 51] or through keypoint detection and reprojection [67] in early methods. More recent research [19, 26] addresses object pose estimation in complex scenarios like partial occlusion. NeMo [53] introduces render-and-compare techniques for category-level object pose estimation, showcasing enhanced robustness in OOD conditions. Later advancements in differentiable rendering [57] and data

augmentation [24] for NeMo have led to further improvements in robust category-level object pose estimation, achieving state-of-the-art performance. However, these approaches are confined to specific object categories and are designed for fixed training datasets only. In contrast, our method for the first time extends them to the class-incremental setting.

2.2 Class-Incremental Learning

Class-incremental learning (also known as continual learning [2, 11, 34] and life-long learning [1, 8, 9]) aims at learning models from sequences of data. The foundational work of [6, 46] replays exemplary data from previously seen classes. The simple strategy has inspired successive works [7, 59]. However, for such methods, sampling strategies and concept drift can impact overall performance. As a mitigation, more recent methods [18, 60] combine replay with other notable regularization schemes like knowledge distillation [27]. In general, class-incremental methods leverage one or more principles from the following three categories: (1) exemplar replay methods build a reservoir of samples from old training rounds [4, 29, 32, 35, 44, 46, 48] and replay them in successive training phases as a way of recalling past knowledge, (2) regularization-based (distillation-based) methods try to preserve the knowledge captured in a previous version of the model by matching logits [27, 46], feature maps [12], or other information [21, 28, 43, 49, 50, 58] in the new model, and (3) network-architecture-based methods [30, 58] design incremental architectures by expanding the network capacity for new class data or freezing partial network parameters to retain the knowledge about old classes.

In our work, we make use of principles from all three of the above by leveraging a replay memory, presenting a novel regularization scheme and adding newly trained neural meshes to the model over time. To the best of our knowledge, our method is the first to combine a 3D inductive bias with these strategies.

3 Prerequisites

3.1 Class Incremental Learning (CIL)

Conventionally, classification models are trained on a single training dataset \mathcal{T} that contains all classes. Multi-class incremental learning departs from this setting by training models on sequentially incoming datasets of new classes that are referred to as tasks $\mathcal{T}^1, \mathcal{T}^2, \dots, \mathcal{T}^{N_{task}}$, where each task may contain more than one new class. After training on a new task \mathcal{T}^i , the model may be evaluated on a test dataset $\mathcal{D}^{1:i}$ that contains classes from all tasks up to i .

When being trained on new tasks through a straightforward fine-tuning, models suffer from *catastrophic forgetting* [22], which leads to bad performance on the previously seen classes. An intuitive approach to mitigate this effect is to use a *replay buffer* [46] that stores a few exemplars $|\mathcal{E}^i| \ll |\mathcal{T}^i|$ from previous tasks and includes them with training data of the new task. Another common

technique is *knowledge distillation* [17, 27] that keeps a copy of the model before training on the new task and ensures that distribution of the feature space from the old and new models are similar when presented the new data.

3.2 Neural Mesh Models

Neural mesh models combine a 2D feature extractor with generative 3D models, as shown by Wang et al. [53] in their Figure 1. The generative models are simple 3D abstractions in the form of cuboids for each class c that are represented as meshes $\mathfrak{N}_c = (\mathcal{V}_c, \mathcal{A}_c, \Theta_c)$, where \mathcal{V}_c denotes the vertices, \mathcal{A}_c denotes the triangles and Θ_c denotes the neural vertex features. The meshes are additionally accompanied by a set of background features \mathcal{B} . Given camera intrinsics and extrinsics, a mesh can then be rendered to a 2D feature map. The 2D feature extractor is usually a 2D CNN $\Phi(I)$ that takes the image as input to extract a feature map and is shared among all classes [20]. Render-and-compare can then be used to check if the features rendered from the mesh align with the features extracted from the image to perform pose estimation [53] or classification [20]. We denote a normalized feature vector at vertex k as θ_c^k , its visibility in the image as o_c^k , its projected integer image coordinates as $\pi_c(k)$, and $f_{\pi_c(k)}$ as the normalized feature vector from the 2D feature extractor that corresponds to the rendered vertex k .

During training, images and object poses are provided, and the vertex features Θ , background features \mathcal{B} , and the 2D feature extractor Φ are trained. We model the probability distribution of a feature f being generated from a vertex v_c^k by defining $P(f|\theta_c^k)$ using a von Mises-Fisher (vMF) distribution to express the likelihood:

$$P(f|\theta_c^k, \kappa) = C(\kappa) e^{\kappa(f^\top \cdot \theta_c^k)}, \quad (1)$$

with mean θ_c^k , concentration parameter κ , and normalization constant $C(\kappa)$ [53]. In the next step, the extracted feature $f_{\pi_c(k)}$ is inserted into $P(f|\theta_c^k, \kappa)$ and maximized using contrastive learning. Simultaneously, the likelihood of all other vertices and background features is minimized:

$$\max P(f_{\pi_c(k)}|\theta_c^k, \kappa), \quad (2)$$

$$\min \sum_{\theta^m \in \bar{\theta}_c^k} P(f_{\pi_c(k)}|\theta^m, \kappa), \quad (3)$$

where the alternative vertices are defined as $\bar{\theta}_c^k = \{\mathcal{B} \cup \Theta_{\bar{c}} \cup (\Theta_c \setminus \mathcal{N}_c^k)\}$ with the neighborhood $\mathcal{N}_c^k = \{\theta_i | \|v_i^k - v_c^k\| < R \wedge v_i^k \in \mathcal{V}_c \setminus v_c^k\}$ around v_c^k determined by some pre-defined distance threshold R . We formulate the Equations 2 and 3 into a single loss by taking the negative log-likelihood:

$$\mathcal{L}_{\text{train}} = - \sum_k o_c^k \cdot \log \left(\frac{e^{\kappa(f_{\pi_c(k)}^\top \cdot \theta_c^k)}}{\sum_{\theta^m \in \bar{\theta}_c^k} e^{\kappa(f_{\pi_c(k)}^\top \cdot \theta^m)}} \right), \quad (4)$$

where considering κ as a global hyperparameter allows cancelling out the normalization constants $C(\kappa)$.

The concentration parameter κ determines the spread of the distribution and can be interpreted as an inverse temperature parameter. In practice, the neural vertex features Θ and the background features \mathcal{B} are unknown and need to be optimized jointly with the feature extractor Φ . This makes the training process initially ambiguous, where a good initialization of Φ and Θ is critical to avoid divergence. After each update of Φ , we therefore follow Bai *et al.* [3] and use the momentum update strategy to train the foreground model Θ_c of a class c , as well as the background model \mathcal{B} :

$$\theta_c^{k,new} \leftarrow o_c^k(1 - \eta) \cdot f_{\pi_c(k)} + (1 - o_c^k + \eta \cdot o_c^k)\theta_c^k, \quad (5)$$

where η is the momentum parameter. The background model \mathcal{B} is updated by sampling $N_{bgupdate}$ feature vectors at pixel positions that are not matched to any vertex of the mesh and replace the $N_{bgupdate}$ oldest features in \mathcal{B} . Both $N_{bgupdate}$ and η are hyperparameters. For a more detailed description of this process, we refer to the supplementary material.

4 Incremental Neural Mesh Models (iNeMo)

Our goal is to learn a model that generalizes robustly in OOD scenarios, while being capable of performing class-incremental learning. To achieve this, we build up on neural mesh models [53] and present a novel formulation for class-incremental learning for classification and object pose estimation that we call iNeMo. An overview is provided in Figure 1.

Challenges in CIL. In the non-incremental setting, the contrastive loss in Equation 4 does not explicitly enforce separating classes, although in practice it is observed that the classes are separated well and accurate classification can be achieved [20]. A naive extension of neural mesh models to class-incremental learning is to simply add a mesh \mathfrak{N}_c for each new class. However, the challenge lies in updating the shared 2D feature extractor Φ . If adding classes naively, achieving a discriminative latent space requires restructuring it as a whole and therefore implies significant changes in both, the CNN backbone Φ and the neural meshes Θ , leading to catastrophic forgetting if no old training samples are available or other measures are taken. Therefore, in the following we present a novel class-incremental learning strategy that maintains a well structured latent space from the beginning.

4.1 Initialization

Latent Space. As the features θ_c^k are normalized, they lie on a unit sphere. We therefore define an initial population $\mathbf{H} = \{\mathbf{h}_j \mid \mathbf{h}_j \in \mathbb{R}^d \wedge \|\mathbf{h}_j\| = 1\}$ of the latent space for all vertices and classes by uniformly sampling the sphere. To partition the latent space, we define a fixed upper bound of classes N . We then generate centroids $\mathbf{E} = [\mathbf{e}_1, \dots, \mathbf{e}_N]$ for all the classes on the unit sphere that are

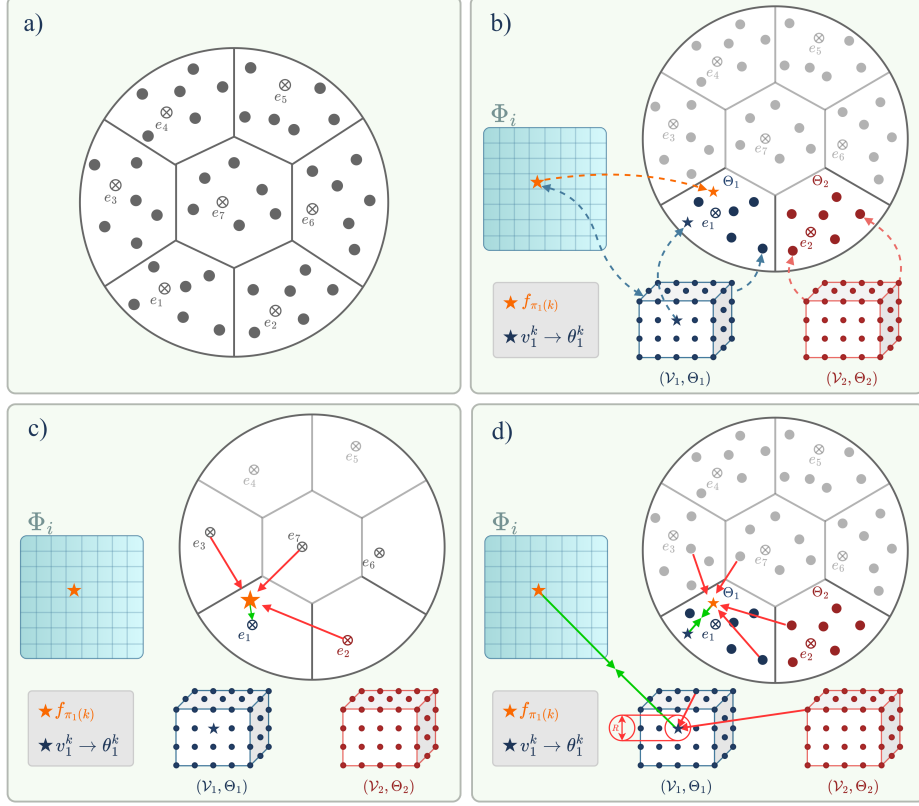


Fig. 2: Overview of Regularization: **a)** The features are constrained to lie on a unit sphere and the latent space is initially uniformly populated. Centroids e_i are then computed to lie maximally far apart, and the feature population is partitioned for a maximum number of classes. **b)** When starting a new task, the vertex features for each new cube from this task are randomly initialized from some class partition. By projecting the locations of the vertices to images, corresponding image features are determined as illustrated by the orange star. **c)** To avoid entanglement, we regularize the latent space by constraining the image feature to stay within the class partition using \mathcal{L}_{etf} . **d)** We then employ the contrastive loss $\mathcal{L}_{\text{cont}}$ that pulls the vertex and image features together and separates the image feature from other features of its own, and the other meshes.

pairwise maximally far apart by solving the equation for a simplex Equiangular Tight Frame (ETF) [41]:

$$\mathbf{E} = \sqrt{\frac{N}{N-1}} \mathbf{U} \left(\mathbf{I}_N - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top \right), \quad (6)$$

where \mathbf{I}_N denotes the n -dimensional identity matrix, $\mathbf{1}_n$ is an all-ones vector, and $\mathbf{U} \in \mathbb{R}^{d \times n}$ is any matrix that allows rotation. The column vectors are of

equal Euclidean norm and any pair has an inner product of $\mathbf{e}_i^\top \cdot \mathbf{e}_j = -\frac{1}{N-1}$ for $i \neq j$, which together ensures pairwise maximum distances. Finally, we assign the features \mathbf{h}_j to classes by determining the respective closest centroid from \mathbf{E} , which leads to a partitioning $\mathbf{H}_1, \dots, \mathbf{H}_N$ of \mathbf{H} . An illustration of this strategy is provided in Figure 2 a).

Task \mathcal{T}^i . At the start of each task, we need to introduce new neural meshes. Following Wang *et al.* [53], for each new class c we initialize \mathfrak{N}_c as a cuboid where its dimensions are determined from ground-truth meshes and vertices are sampled on a regular grid on the surface. As illustrated in Figure 2 b), we then pick the partition \mathbf{H}_c of initial features and randomly assign them to the vertices of the new mesh Θ_c . We initialize the feature extractor Φ_0 with unsupervised pre-training using DINO-v1 [5]. As shown in Figure 1, to train for a new task, we make a copy $\Phi_i = \Phi_{i-1}$ and then leverage Φ_{i-1} for knowledge distillation. If available, we discard any old network Φ_{i-2} .

4.2 Optimization

Positional Regularization. To ensure that our latent space maintains the initial partitioning over time, we introduce a penalty of the distance of the neural features Θ_c to their corresponding class centroid \mathbf{e}_c :

$$\mathcal{L}_{\text{etf}} = - \sum_k o_c^k \cdot \log \left(\frac{e^{\kappa_2 (f_{\pi_c(k)}^\top \cdot \mathbf{e}_c)}}{\sum_{\mathbf{e}_m \in \mathbf{E}} e^{\kappa_2 (f_{\pi_c(k)}^\top \cdot \mathbf{e}_m)}} \right). \quad (7)$$

This is illustrated in Figure 2 c).

Continual Training Loss. We denote any unused partitions in \mathbf{H} with $\bar{\mathbf{H}}$ and limit the spread of the neural meshes in the current task to refrain from $\bar{\mathbf{H}}$ by posing the following additional contrastive loss:

$$\mathcal{L}_{\text{cont}} = - \sum_k o_c^k \cdot \log \left(\frac{e^{\kappa_1 (f_{\pi_c(k)}^\top \cdot \theta_c^k)}}{\sum_{\theta^m \in \bar{\theta}_k} e^{\kappa_1 (f_{\pi_c(k)}^\top \cdot \theta^m)} + \sum_{h_j \in \bar{\mathbf{H}}} e^{\kappa_1 (f_{\pi_c(k)}^\top \cdot h_j)}} \right). \quad (8)$$

The denominator is split into two parts, where the first one minimizes Equation 3 and the second part corresponds to the additional constraint imposed by the features in the unused partitions $\bar{\mathbf{H}}$. This is illustrated in Figure 2 d).

Knowledge Distillation. To mitigate forgetting, as indicated in Figure 1, we additionally use a distillation loss after the initial task. The new inputs are also fed through the frozen backbone Φ_{i-1} of the previous task to obtain its feature map. Specifically, let $\hat{f}_{\pi_c(k)}$ denote the old feature for the vertex k . To

distill classes from previous tasks into Φ_i , we formulate the distillation using the Kullback-Leibler divergence:

$$\mathcal{L}_{\text{kd}} = - \sum_k \sum_m p_m(\hat{f}_{\pi_c(k)}) \log \left(\frac{p_m(\hat{f}_{\pi_c(k)})}{p_m(f_{\pi_c(k)})} \right), \quad (9)$$

where:

$$p_m(f_{\pi_c(k)}) = \frac{e^{\kappa_3(f_{\pi_c(k)}^\top \cdot \theta_m)}}{\sum_{\theta_m \in \Theta^{i-1}} e^{\kappa_3(f_{\pi_c(k)}^\top \cdot \theta_m)}}. \quad (10)$$

Note that, unless we are considering an exemplar of a previous task, the real corresponding feature θ_c^k is not even considered in this formulation. However, the aim here is not to optimize Φ_i for the current task, but to extract the dark knowledge [17] from Φ_{i-1} about classes from previous tasks. Consequently, the concentration $\kappa_3 < 1$ has to be small to get usable gradients from all likelihoods.

Continual Training. During training of Φ_i , we optimize the combined training objective:

$$\mathcal{L} = \mathcal{L}_{\text{cont}} + \lambda_{\text{etf}} \mathcal{L}_{\text{etf}} + \lambda_{\text{kd}} \mathcal{L}_{\text{kd}}, \quad (11)$$

where λ_{etf} and λ_{kd} are weighting parameters.

4.3 Exemplar Selection

At each training stage, we randomly remove exemplars from old classes to equally divide our replay buffer for the current number of classes. Xiang *et al.* [62] showed that certain classes are heavily biased towards certain viewing angles. Therefore, to increase the robustness and accuracy for rarely appearing view directions, we propose an exemplar selection strategy that takes viewing angles into account. Assuming we want to integrate a new class and the available slots for it are m , we build a b -bin histogram across the azimuth angles and randomly select $\lfloor m/b \rfloor$ exemplars for each bin. When insufficient exemplars are available for a bin we merge it together with a neighboring one. In case the process yields less than m exemplars in total, we fill up remaining slots with random samples. When reducing the exemplar sets, we evenly remove samples from each bin to maintain the balance across the azimuth angle distribution.

4.4 Inference

Classification. Following Jesslen *et al.* [20], we perform classification via a vertex matching approach. For each feature f_i in the produced feature map of Φ , we compute its similarities to the foreground (Θ) and background (\mathcal{B}) models. We define the background score s_β^i and the class scores s^i for each class c as

$$s_c^i = \max_{\theta_c^l \in \Theta_c} f_i^\top \cdot \theta_c^l, \quad (12)$$

$$s_\beta^i = \max_{\beta^l \in \mathcal{B}} f_i^\top \cdot \beta^l, \quad (13)$$

where we identify a feature as being in the foreground \mathcal{F} , if there is at least one $s_c^i > s_\beta^i$ and classify based on the foreground pixels only.

In contrast to Jesslen *et al.* [20], we additionally include an uncertainty term to reduce the influence of features that can not be identified with high confidence. In the following, we denote the n -th largest class score for feature f_i as $\max_{s_i}^{(n)}$. The final score of class y is then given as

$$s_y = \max_{i \in \mathcal{F}} \left[s_y^i - (1 - (\max_{s_i}^{(1)} - \max_{s_i}^{(2)})) \right], \quad (14)$$

where the subtracted term indicates a measure of confusion estimated based on the difference of the two highest class scores for foreground feature f_i . The predicted category is then simply the class c that maximizes this score.

Pose Estimation. For pose estimation we use the same render-and-compare approach as Wang *et al.* [53] together with the template matching proposed by Jesslen *et al.* [20] for speedup. For more information about the pose estimation, we refer the reader to the supplemental material.

5 Experiments

In the following, we explain the experimental setup and then discuss the results of our incremental neural mesh models for image classification and 3D pose estimation on both, in-domain and OOD datasets. For a comprehensive ablation study of all components of our model, we refer to the supplemental material.

5.1 Datasets and Implementation Details

In-Domain-Datasets. **PASCAL3D+** [63] (P3D) has high-quality camera pose annotations with mostly unoccluded objects, making it ideal for our setting. However, with only 12 classes it is small compared to other datasets used in continual learning [25, 47]. **ObjectNet3D** [61] (O3D) contains 100 classes and presents a significantly more difficult setting. Camera pose annotations are less reliable and the displayed objects can be heavily occluded or truncated, making both the vertex mapping and the update process noisy.

OOD-Datasets. The **Occluded-PASCAL3D+** [56] (O-P3D) and **corrupted-PASCAL3D+** (C-P3D) datasets are variations of original P3D and consist of a test dataset only. In the O-P3D dataset, parts of the original test datasets have been artificially occluded by superimposing occluders on the images with three different levels: L1 (20% – 40%), L2 (40% – 60%) and L3 (60% – 80%). The C-P3D dataset, on the other hand, follows [15] and tests robustness against image corruptions. We evaluate 19 different corruptions with a severity of 4 out of 5, using the **imagecorruptions** [39] library. Finally, we consider the **OOD-CV** [65] dataset, which provides a multitude of severe domain shifts.

Table 1: Average classification accuracies on Pascal3D (P3D) and ObjectNet3D (O3D). Training data has been split into a base task (denoted Bn for size n) and evenly sized increments (denoted $+n$ for size n). As visible, our method consistently outperforms the baselines by a significant margin.

Metric	Method	Repr.	<i>P3D</i>		<i>O3D</i>		
			$B0 + 6$	$B0 + 3$	$B0 + 20$	$B0 + 10$	$B50 + 10$
Classification $acc(1 : i)$ in % \uparrow	LWF	R50.	93.83	89.34	67.78	48.82	46.73
	FeTrIL	R50.	95.64	96.82	67.18	70.34	70.43
	FeCAM	R50.	84.85	64.36	67.96	69.59	72.21
	iCaRL	R50	97.1	93.80	72.55	57.46	64.02
	DER	R50	96.69	94.18	78.55	76.33	75.17
	Podnet	R50	95.13	91.71	71.96	65.21	72.98
	Ours	NeMo	98.82	98.21	89.25	88.85	84.20

Implementational Details. We choose a ResNet50 architecture for our feature extractor Φ with two upsampling layers and skip connections, resulting in a final feature map at $\frac{1}{8}$ of the input resolution. Each neural mesh \mathfrak{N}_y contains approximately 1,100 uniformly distributed vertices with a neural texture of dimension $d = 128$. We train for 50 epochs per task, with a learning rate of $1e - 5$ that is halved after 10 epochs. The replay buffer can store up to 240 and 2,000 samples for P3D and O3D respectively. Our feature extractor is optimized using Adam with default parameters and the neural textures Θ are updated with momentum of $\eta = 0.9$. During pose estimation, we initialize the camera pose using template matching as proposed by [20] and optimize it with PyTorch3D’s differentiable rasterizer [45]. The initial camera pose is refined by minimizing the reconstruction loss between the feature map produced by Φ and the rendered mesh. We use Adam with a learning rate of 0.05 for 30 total epochs and a distance threshold $R = 48$ to measure the neighborhood \mathcal{N}_c^k in Equation 8. Each term in the combined loss in Equation 11 is assigned a weighting and concentration parameter. The weighting parameters are $\lambda_{\text{etf}} = 0.2$ and $\lambda_{\text{kd}} = 2.0$ and as concentration parameters we choose $\kappa_1 = 1/0.07 \approx 14.3$, $\kappa_2 = 1$, and $\kappa_3 = 0.5$. We provide further details on the training settings of the baselines in the supplemental material.

Evaluation. We evaluate our method and its baselines on both, class-incremental classification and class-incremental pose estimation. For the methods trained on P3D, we evaluate on the P3D test dataset, the O-P3D dataset, and the C-P3D dataset. When training on O3D or OOD-CV, we evaluate on their corresponding test dataset only. For classification, we follow previous work [27, 32, 46] and consider the mean accuracy over all tasks $acc(1 : i)$ of Φ_i , after training on \mathcal{T}^i on test dataset \mathcal{D} for classes $1..i$. The 3D pose of an object can be represented with azimuth, elevation, and roll angle. We measure the deviation of predicted and ground-truth pose in terms of these angles according to the error of the predicted

Table 2: Average classification and pose estimation accuracies on Pascal3D (P3D) and its variants. As visible, iNeMo outperforms all 2D baselines consistently for classification and by an especially large margin for the OOD and strong occlusion cases. We also present the first approach for incremental pose estimation and outperform other methods in most cases for $\pi/6$, while we consistently outperform them for the tighter error bound $\pi/18$. Note that for all evaluations except OOD-CV, we use the model trained on 4 tasks that is also displayed in Figure 3. As OOD-CV provides a separate training set of 10 classes, we consider 2 tasks with 5 classes.

Metric	Method	Repr.	<i>P3D</i>	<i>Occluded P3D</i>			<i>C-P3D</i>	<i>OOD-CV</i>
				<i>L1</i>	<i>L2</i>	<i>L3</i>		
Classification $acc(1:i)$ in % \uparrow	LwF	R50	89.34	30.58	21.64	14.65	66.17	57.61
	FeTrIL	R50	96.82	88.34	76.28	55.89	39.02	63.74
	FeCAM	R50	84.85	52.53	42.75	34.28	42.38	56.05
	iCaRL	R50	93.80	34.95	26.00	16.93	76.24	61.80
	DER	R50	94.18	49.70	36.86	22.76	69.56	56.35
	PODNet	R50	91.91	42.40	32.99	22.43	68.46	57.10
	Ours	NeMo	98.21	94.19	87.20	71.55	83.09	80.82
Pose $\pi/6$ $acc(1:i)$ in % \uparrow	LwF	R50	53.47	44.58	39.77	36.61	53.55	30.65
	iCaRL	R50	57.74	44.03	38.15	33.52	54.57	28.71
	Ours	NeMo	79.28	64.71	52.26	34.01	47.30	33.75
Pose $\pi/18$ $acc(1:i)$ in % \uparrow	LwF	R50	20.33	12.03	8.38	5.52	17.29	8.04
	iCaRL	R50	22.76	11.04	7.33	4.56	17.81	8.04
	Ours	NeMo	51.73	35.53	26.88	10.672	23.02	12.8

and the ground-truth rotation matrix $\Delta(R_{\text{pred}}, R_{\text{gt}}) = \|\log m(R_{\text{pred}}^{\top} R_{\text{gt}})\|_F / \sqrt{2}$ as proposed by [67]. Following previous work [53, 67], we report the accuracy according to the thresholds $\pi/6$ and $\pi/18$.

Baselines. For the task of class-incremental learning, we compare against a collection of replay-based and replay-free methods. For the replay-based methods, we choose the seminal work iCaRL [46], and the more recent PODNet [12] and DER [64]. For replay-free methods, we choose the seminal work LwF [27] and the two state-of-the-art methods FeTrIL [42] and FeCAM [13]. All approaches are implemented using the PyCIL library [66] and trained with the original hyperparameters as in [66]. For a fair comparison of all methods, we use the ResNet-50 backbone initialized with DINO-v1 [5] pre-trained weights.

To the best of our knowledge, incremental pose estimation with a class-agnostic backbone has not been explored before. We define incremental pose estimation baselines by discretizing the polar camera coordinates and formulate pose estimation as a classification problem following [67]. More specifically, we define 42 bins for each azimuth, elevation and roll angle, making it a $42 \cdot 3 = 126$ class classification problem [67]. This allows a straightforward extension of conventional class-incremental learning techniques to the setting of pose estimation.

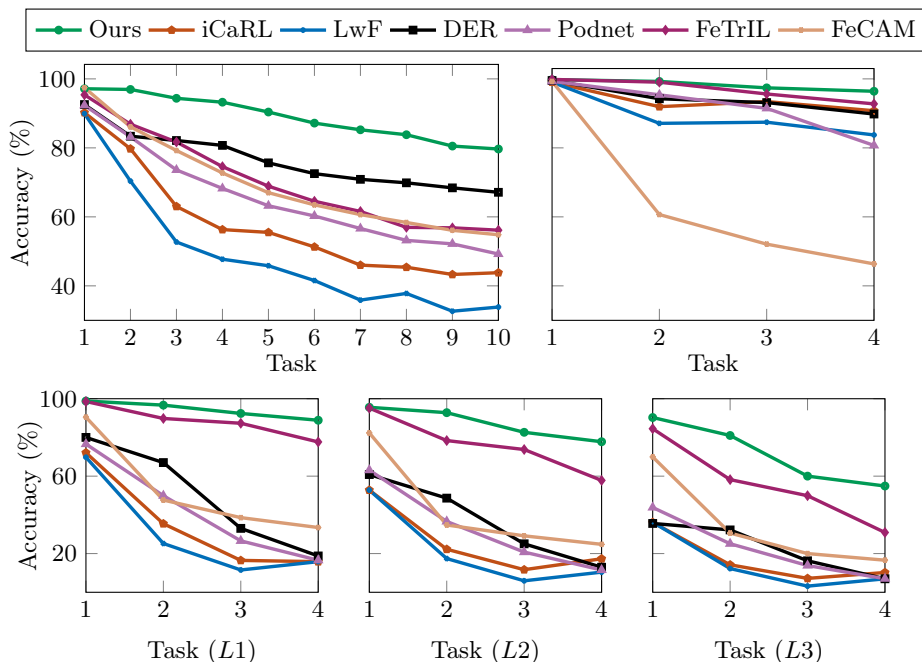


Fig. 3: Comparison of classification performance decay over tasks for our method and the baselines. **Top-Left:** Results for O3D (100 classes) split into 10 even tasks. **Top-Right:** Results for P3D (12 classes) split into 4 even tasks. **Bottom:** Results for O-P3D with occlusion levels L1, L2 and L3 after each task. One can observe that our method outperforms all other methods. Especially in the occluded cases, our method outperforms them by a very large margin up to 70%, even still showing strong performance for the largest occlusion level L3 with 60 – 80% occlusions.

We provide such class-incremental pose estimation results using the training procedure of iCaRL [46] and LwF [27]. Both methods are trained for 100 epochs per task using SGD with a learning rate of 0.01 as proposed by [67].

5.2 Robust Class-Incremental Classification

In Table 1, we provide the in-distribution classification results for P3D and O3D. Our method outperforms the baselines in all cases, including the harder O3D setting with 100 classes. Furthermore, Table 2 shows the comparison of class-incremental results on all P3D variants for both, classification and 3D pose estimation. As visible, our method outperforms the other methods with a large margin under domain shifts: for the L3 occluded case, it is larger than 48%, for the corrupted C-P3D it is larger than 6%, and for the OOD-CV dataset it is larger than 19%. Figure 3 shows task-wise accuracy on the O3D/P3D dataset for 10 and 4 even tasks respectively, as well as the task-wise accuracy on the O-P3D

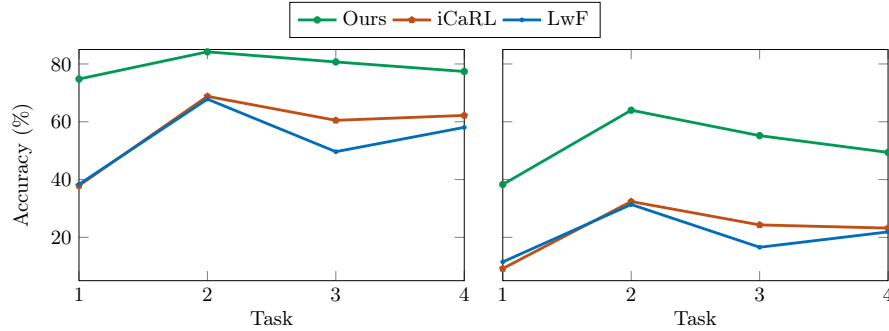


Fig. 4: Comparison of the task-wise pose estimation accuracy on P3D for 4 even tasks, where we show the thresholds **left:** $\pi/6$ and **right:** $\pi/18$. One can observe that our method outperforms all other methods and retains high pose estimation accuracy throughout the incremental training process. One can also observe that for pose estimation, there is a stronger dependence on the difficulty of the considered classes instead of the method’s ability to retain knowledge.

dataset for all occlusion levels. The same observation as before can be made, where our method exhibits significantly less performance decay over new tasks. This overall demonstrates that our incremental neural mesh models outperform their baselines decisively in robustness.

5.3 Class Incremental Pose Estimation

Table 2 also shows that our method significantly outperforms both ResNet-50 based methods for the task of incremental pose estimation. As visible, the feature representation learned by the 2D pose estimation networks is much less affected by both, catastrophic forgetting and domain shifts. Figure 4 shows that the performance decrease is much less severe across all tasks, where the difference in performance is much more dependent on the difficulty of the considered classes instead of the method’s ability to retain knowledge.

6 Conclusions

In this work, we introduce incremental neural mesh models, which enable robust class-incremental learning for both, image classification and 3D pose estimation. For the first time, we present a model that can learn new prototypical 3D representations of object categories over time. The extensive evaluation on Pascal3D and ObjectNet3D shows that our approach outperforms all baselines even in the in-domain setting and surpasses them by a large margin in the OOD case. We also introduced the first approach for class-incremental learning of pose estimation. The results overall demonstrate the fundamental advantage of 3D object-centric representations, and we hope that this will spur a new line of research in the community.

Acknowledgements

We gratefully acknowledge the stimulating research environment of the GRK 2853/1 “Neuroexplicit Models of Language, Vision, and Action”, funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under project number 471607914.

References

1. Aljundi, R., Chakravarty, P., Tuytelaars, T.: Expert gate: Lifelong learning with a network of experts. In: CVPR. pp. 3366–3375 (2017)
2. Aljundi, R., Kelchtermans, K., Tuytelaars, T.: Task-free continual learning. In: CVPR. pp. 11254–11263 (2019)
3. Bai, Y., Wang, A., Kortylewski, A., Yuille, A.: Coke: Localized contrastive learning for robust keypoint detection. Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (2023)
4. Bang, J., Kim, H., Yoo, Y., Ha, J.W., Choi, J.: Rainbow memory: Continual learning with a memory of diverse samples. In: CVPR. pp. 8218–8227 (2021)
5. Caron, M., Touvron, H., Misra, I., Jégou, H., Mairal, J., Bojanowski, P., Joulin, A.: Emerging properties in self-supervised vision transformers. In: Proceedings of the International Conference on Computer Vision (ICCV) (2021)
6. Castro, F.M., Marín-Jiménez, M.J., Guil, N., Schmid, C., Alahari, K.: End-to-end incremental learning. In: ECCV. pp. 233–248 (2018)
7. Chaudhry, A., Dokania, P.K., Ajanthan, T., Torr, P.H.: Riemannian walk for incremental learning: Understanding forgetting and intransigence. In: ECCV. pp. 532–547 (2018)
8. Chaudhry, A., Ranzato, M., Rohrbach, M., Elhoseiny, M.: Efficient lifelong learning with A-GEM. In: ICLR (2019)
9. Chen, Z., Liu, B.: Lifelong machine learning. Synthesis Lectures on Artificial Intelligence and Machine Learning **12**(3), 1–207 (2018)
10. Cimpoi, M., Maji, S., Kokkinos, I., Mohamed, S., Vedaldi, A.: Describing textures in the wild. In: Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR) (2014)
11. De Lange, M., Aljundi, R., Masana, M., Parisot, S., Jia, X., Leonardis, A., Slabaugh, G., Tuytelaars, T.: A continual learning survey: Defying forgetting in classification tasks. TPAMI **44**(7), 3366–3385 (2021)
12. Douillard, A., Cord, M., Ollion, C., Robert, T., Valle, E.: Podnet: Pooled outputs distillation for small-tasks incremental learning. In: ECCV. pp. 86–102 (2020)
13. Goswami, D., Liu, Y., Twardowski, B., van de Weijer, J.: Fecam: Exploiting the heterogeneity of class distributions in exemplar-free continual learning. In: Advances in Neural Information Processing Systems. vol. 36 (2024)
14. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR. pp. 770–778 (2016)
15. Hendrycks, D., Dietterich, T.: Benchmarking neural network robustness to common corruptions and perturbations. In: ICLR (2019)
16. Hendrycks, D., Mu, N., Cubuk, E.D., Zoph, B., Gilmer, J., Lakshminarayanan, B.: Augmix: A simple data processing method to improve robustness and uncertainty. arXiv preprint arXiv:1912.02781 (2019)

17. Hinton, G., Vinyals, O., Dean, J., et al.: Distilling the knowledge in a neural network. In: NIPS Workshops (2014)
18. Hou, S., Pan, X., Loy, C.C., Wang, Z., Lin, D.: Learning a unified classifier incrementally via rebalancing. In: CVPR. pp. 831–839 (2019)
19. Iwase, S., Liu, X., Khirodkar, R., Yokota, R., Kitani, K.M.: Repose: Fast 6d object pose refinement via deep texture rendering. In: ICCV. pp. 3303–3312 (2021)
20. Jesslen, A., Zhang, G., Wang, A., Yuille, A., Kortylewski, A.: Robust 3d-aware object classification via discriminative render-and-compare. arXiv preprint arXiv:2305.14668 (2023)
21. Joseph, K.J., Khan, S., Khan, F.S., Anwer, R.M., Balasubramanian, V.N.: Energy-based latent aligner for incremental learning. In: CVPR. pp. 7452–7461 (2022)
22. Kirkpatrick, J., Pascanu, R., Rabinowitz, N., Veness, J., Desjardins, G., Rusu, A.A., Milan, K., Quan, J., Ramalho, T., Grabska-Barwinska, A., et al.: Overcoming catastrophic forgetting in neural networks. PNAS pp. 3521–3526 (2017)
23. Kortylewski, A., Liu, Q., Wang, A., Sun, Y., Yuille, A.: Compositional convolutional neural networks: A robust and interpretable model for object recognition under occlusion. IJCV pp. 1–25 (2020)
24. Kouros, G., Shrivastava, S., Picron, C., Nagesh, S., Chakravarty, P., Tuytelaars, T.: Category-level pose retrieval with contrastive features learnt with occlusion augmentation. arXiv preprint arXiv:2208.06195 (2022)
25. Krizhevsky, A., Hinton, G., et al.: Learning multiple layers of features from tiny images. Technical Report TR-2009 (2009)
26. Li, Y., Wang, G., Ji, X., Xiang, Y., Fox, D.: Deepim: Deep iterative matching for 6d pose estimation. In: ECCV. pp. 683–698 (2018)
27. Li, Z., Hoiem, D.: Learning without forgetting. TPAMI **40**(12), 2935–2947 (2018)
28. Liu, Y., Li, Y., Schiele, B., Sun, Q.: Online hyperparameter optimization for class-incremental learning. In: AAAI (2023)
29. Liu, Y., Li, Y., Schiele, B., Sun, Q.: Wakening past concepts without past data: Class-incremental learning from online placebos. In: WACV. pp. 2226–2235 (January 2024)
30. Liu, Y., Schiele, B., Sun, Q.: Adaptive aggregation networks for class-incremental learning. In: CVPR. pp. 2544–2553 (2021)
31. Liu, Y., Schiele, B., Sun, Q.: RMM: reinforced memory management for class-incremental learning. In: NeurIPS. pp. 3478–3490 (2021)
32. Liu, Y., Su, Y., Liu, A., Schiele, B., Sun, Q.: Mnemonics training: Multi-class incremental learning without forgetting. In: CVPR. pp. 12245–12254 (2020)
33. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. In: ICCV. pp. 10012–10022 (2021)
34. Lopez-Paz, D., Ranzato, M.: Gradient episodic memory for continual learning. In: NIPS. pp. 6467–6476 (2017)
35. Luo, Z., Liu, Y., Schiele, B., Sun, Q.: Class-incremental exemplar compression for class-incremental learning. In: CVPR. pp. 11371–11380. IEEE (2023)
36. Ma, W., Wang, A., Yuille, A.L., Kortylewski, A.: Robust category-level 6d pose estimation with coarse-to-fine rendering of neural features. In: ECCV. pp. 492–508 (2022)
37. McCloskey, M., Cohen, N.J.: Catastrophic interference in connectionist networks: The sequential learning problem. In: Psychology of Learning and Motivation, vol. 24, pp. 109–165. Elsevier (1989)
38. McRae, K., Hetherington, P.: Catastrophic interference is eliminated in pre-trained networks. In: CogSci (1993)

39. Michaelis, C., Mitzkus, B., Geirhos, R., Rusak, E., Bringmann, O., Ecker, A.S., Bethge, M., Brendel, W.: Benchmarking robustness in object detection: Autonomous driving when winter is coming. arXiv preprint arXiv:1907.07484 (2019)
40. Mousavian, A., Anguelov, D., Flynn, J., Kosecka, J.: 3d bounding box estimation using deep learning and geometry. In: CVPR. pp. 7074–7082 (2017)
41. Pappayan, V., Han, X., Donoho, D.L.: Prevalence of neural collapse during the terminal phase of deep learning training. *Proceedings of the National Academy of Sciences* **117**(40), 24652–24663 (2020)
42. Petit, G., Popescu, A., Schindler, H., Picard, D., Delezoide, B.: Fetril: Feature translation for exemplar-free class-incremental learning. In: CVPR (2023)
43. PourKeshavarzi, M., Zhao, G., Sabokrou, M.: Looking back on learned experiences for class/task incremental learning. In: ICLR (2022)
44. Prabhu, A., Torr, P.H., Dokania, P.K.: GDumb: A simple approach that questions our progress in continual learning. In: ECCV. pp. 524–540 (2020)
45. Ravi, N., Reizenstein, J., Novotny, D., Gordon, T., Lo, W.Y., Johnson, J., Gkioxari, G.: Accelerating 3d deep learning with pytorch3d. arXiv:2007.08501 (2020)
46. Rebuffi, S.A., Kolesnikov, A., Sperl, G., Lampert, C.H.: iCaRL: Incremental classifier and representation learning. In: CVPR. pp. 5533–5542 (2017)
47. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al.: Imagenet large scale visual recognition challenge. *International journal of computer vision* **115**, 211–252 (2015)
48. Shin, H., Lee, J.K., Kim, J., Kim, J.: Continual learning with deep generative replay. In: NeurIPS. pp. 2990–2999 (2017)
49. Simon, C., Koniusz, P., Harandi, M.: On learning the geodesic path for incremental learning. In: CVPR. pp. 1591–1600 (2021)
50. Tao, X., Chang, X., Hong, X., Wei, X., Gong, Y.: Topology-preserving class-incremental learning. In: ECCV. pp. 254–270 (2020)
51. Tulsiani, S., Malik, J.: Viewpoints and keypoints. In: CVPR (June 2015)
52. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. *NeurIPS* **30** (2017)
53. Wang, A., Kortylewski, A., Yuille, A.: NeMo: Neural mesh models of contrastive features for robust 3d pose estimation. *ICLR* (2021)
54. Wang, A., Ma, W., Yuille, A., Kortylewski, A.: Neural textured deformable meshes for robust analysis-by-synthesis. In: WACV. pp. 3108–3117 (2024)
55. Wang, A., Mei, S., Yuille, A.L., Kortylewski, A.: Neural view synthesis and matching for semi-supervised few-shot learning of 3d pose. *NeurIPS* **34**, 7207–7219 (2021)
56. Wang, A., Sun, Y., Kortylewski, A., Yuille, A.L.: Robust object detection under occlusion with context-aware compositionalnets. In: CVPR. pp. 12645–12654 (2020)
57. Wang, A., Wang, P., Sun, J., Kortylewski, A., Yuille, A.: Voge: a differentiable volume renderer using gaussian ellipsoids for analysis-by-synthesis. In: ICLR (2022)
58. Wang, F.Y., Zhou, D.W., Ye, H.J., Zhan, D.C.: Foster: Feature boosting and compression for class-incremental learning. In: ECCV (2022)
59. Wu, C., Herranz, L., Liu, X., Van De Weijer, J., Raducanu, B., et al.: Memory replay gans: Learning to generate new categories without forgetting. *NeurIPS* **31** (2018)
60. Wu, Y., Chen, Y., Wang, L., Ye, Y., Liu, Z., Guo, Y., Fu, Y.: Large scale incremental learning. In: CVPR. pp. 374–382 (2019)
61. Xiang, Y., Kim, W., Chen, W., Ji, J., Choy, C., Su, H., Mottaghi, R., Guibas, L., Savarese, S.: Objectnet3d: A large scale database for 3d object recognition. In: ECCV (2016)

- 62. Xiang, Y., Mottaghi, R., Savarese, S.: Beyond pascal: A benchmark for 3d object detection in the wild. In: WACV. pp. 75–82. IEEE (2014)
- 63. Xiang, Y., Mottaghi, R., Savarese, S.: Beyond pascal: A benchmark for 3d object detection in the wild. In: WACV (2014)
- 64. Yan, S., Xie, J., He, X.: Der: Dynamically expandable representation for class incremental learning. In: CVPR. pp. 3014–3023 (2021)
- 65. Zhao, B., Yu, S., Ma, W., Yu, M., Mei, S., Wang, A., He, J., Yuille, A., Kortylewski, A.: Ood-cv: A benchmark for robustness to individual nuisances in real-world out-of-distribution shifts. In: ECCV (2022)
- 66. Zhou, D., Wang, F., Ye, H., Zhan, D.: Pycil: a python toolbox for class-incremental learning. *Sci. China Inf. Sci.* **66**(9) (2023)
- 67. Zhou, X., Karpur, A., Luo, L., Huang, Q.: Starmap for category-agnostic keypoint and viewpoint estimation. In: ECCV. pp. 318–334 (2018)

Supplementary Material for iNeMo: Incremental Neural Mesh Models for Robust Class-Incremental Learning

In the following, we provide further details and ablation studies for our paper. In the first section we define the conventions. We then provide the non-incremental performance of both considered representations (R50 and NeMo) as a reference. Afterwards, we show the advantage of considering uncertainty for the classification in Section C and then give a conclusive ablation study over all components of our method in Section D. Since NeMo is trained with additional pose labels that were not available to baselines, we provide an additional study in Section F where we show that pose labels do not improve the baselines. Finally, we conclude with additional details on the background model and pose estimation, as well as all the considered hyperparameters in our method.

A Conventions

In the tables of the main paper, we followed previous work [27, 32, 46] and reported the average of the testing accuracies over all tasks \mathcal{T}^i with Φ_i , which we denoted as $acc(1 : i)$. In the supplemental material, we deviate from this setting and **report the final accuracy with $\Phi_{N_{\text{task}}}$ on the whole test dataset after integrating all tasks**, as it determines the final performance loss that is usually most significant. We denote the final accuracy on all seen classes after training on the final task $\mathcal{T}^{N_{\text{task}}}$ as $\overline{acc}(1 : N_{\text{task}})$.

B Non-Incremental Upper Bounds

To determine an upper bound for the performance of ResNet50 and NeMo approaches, we train on all classes jointly and provide the results in Table 3. While both approaches are able to achieve similar performance for classification on P3D, NeMo significantly outperforms the ResNet50 on O3D. We suspect that the reason for this is that O3D contains a large number of occluded and truncated objects. NeMo generally outperforms the ResNet50 for pose estimation implemented following [67].

C Considering Uncertainty in Classification

We proposed an extension to the classification strategy introduced by Jesslen *et al.* [20] in Equation 14 of the main paper which was motivated by the observation that classes sharing visually similar features were confused more often when training the model in an incremental setting. We believe that when training on all classes jointly, the contrastive loss between all features of different classes is sufficient to ensure that all parts of different objects have distinct feature representations. However, such disentanglement is significantly more challenging in an incremental setting. The results from Table 4 show that our proposed strategy to exclude pixels that ambiguously relate to meshes of multiple possible classes (i.e. uncertain pixels) brings a significant improvement.

Table 3: We trained the RestNet50 and NeMo approaches jointly on all classes to determine an upper bound for their performance. All networks were initialized with weights from DINOv1 [5], which itself was trained in an unsupervised fashion. For joint training of NeMo, we follow the training protocol of Jesslen *et al.* [20] with the exception of using pre-trained weights as mentioned.

Metric	Type	Repr.	<i>P3D</i>	<i>O3D</i>
Classification	Joint	R50	98.32	76.23
$\overline{acc}(1 : N_{\text{task}})$ in % \uparrow	Joint	NeMo	99.27	85.28

Metric	Type	Repr.	<i>P3D</i>
Pose $\pi/6$	Joint	R50	74.6
$\overline{acc}(1 : N_{\text{task}})$ in % \uparrow	Joint	NeMo	87.25
Pose $\pi/18$	Joint	R50	36.5
$\overline{acc}(1 : N_{\text{task}})$ in % \uparrow	Joint	NeMo	65.81

Table 4: We compare our inference approach to the one proposed by Jesslen *et al.* [20]. When training jointly, the performance is nearly identical. However, when training incrementally, disentangling visually similar features becomes more challenging and our proposed strategy significantly improves the result.

Metric	Type	Inference	<i>P3D</i>	<i>O3D</i>
			<i>B0 + 3</i>	<i>B0 + 20</i>
Classification	Joint	[20]	99.28	85.28
	Joint	Ours	99.27	85.28
$\overline{acc}(1 : N_{\text{task}})$ in % \uparrow	Incremental	[20]	95.06	75.8
	Incremental	Ours	96.41	80.17

Table 5: Top: We provide an ablation study for the 2D ResNet50 and NeMo with the traditional class-incremental techniques LwF and iCaRL. As visible, traditional techniques work less well on NeMo. **Bottom:** We provide an ablation study of our model components and show that all of our additions increase the performance. We indicate the used exemplar selection strategy in the column "Replay", where H denotes the herding strategy [46] and PA our pose-aware exemplar selection strategy. Note that we used our improved inference method from Section C for all methods.

Metric	Method	Repr.	Replay	Init	Pos.	KD	<i>P3D</i>	<i>O3D</i>
							<i>B0 + 3</i>	<i>B0 + 20</i>
Classification	Finetune	NeMo	-				17.47	17.81
	LwF	R50	-			✓	83.75	56.44
	LwF	NeMo	-			✓	17.45	17.72
	iCaRL	R50	H			✓	91.79	61.75
	iCaRL	NeMo	H			✓	93.72	68.87
$\overline{acc}(1 : N_{\text{task}})$ in % \uparrow	Ours	NeMo	H				93.60	69.01
	Ours	NeMo	PA				94.70	70.32
	Ours	NeMo	PA	✓			94.78	71.67
	Ours	NeMo	PA	✓	✓		94.98	72.09
	Ours	NeMo	PA	✓	✓	✓	96.41	80.17

D Ablation

In the main paper, we have shown that our novel class-incremental learning strategy with neural meshes significantly outperforms the 2D baselines. In the following, we provide an analysis of how much the individual parts of our model contribute to this result.

Table 5 shows the contribution of each of our model components. We start with the most naive extension of NeMo to the class-incremental setting: in each task, we initialize the required number of meshes and fine-tune the feature extractor on each new task dataset. As expected, this leads to bad results. Next, we extend the models by the traditional distillation [27] (LwF) and herding exemplar [46] (iCaRL) strategies. The latter brings significant improvements. This shows that overall exemplar replay is necessary to retain knowledge while training Neural Mesh Models incrementally. We also compare applying LwF and iCaRL to either the 2D ResNet50 or NeMo and find that those strategies in most cases work better for the 2D setting, hence not being simply transferable to NeMo.

We finally demonstrate that our additions to maintain a structured latent space provide the best results by introducing the latent space initialization, positional regularization, and adding knowledge distillation. The results indicate that knowledge distillation has little effect (row 9), while adding the pose-aware replay (row 5 to row 6) has the largest impact on the result. This shows that the pose-aware exemplar selection strategy is critical and all other additions further improve the performance.

Table 6: Final task accuracy on Pascal3D with decreasing number of exemplars per class. Even with few exemplars, iNeMo retains good accuracy.

Metric	Exemplars	<i>P3D</i>
		<i>B0 + 3</i>
Classification	20	96.41
$\overline{acc}(1 : N_{\text{task}})$ in % \uparrow	10	93.36
	5	82.59

E Training with less Replay Memory

Memory replay is essential when training iNeMo, as it allows updating neural meshes from previous tasks. However, storing too many samples per class in memory can become quite expensive and as such it is crucial for methods to be effective in utilizing replay with fewer samples. We show in Table 6 that iNeMo can adapt to lower memory sizes, but is optimal for the chosen 20 exemplars.

F Enhancing 2D Classifiers with Pose Annotations

Neural Mesh Models leverage meshes to host 3D consistent features and consequently, their training requires camera pose annotations. However, such pose annotations were not used in the 2D baselines, which could in principle give the Neural Mesh Models an advantage. To this end, we evaluate if using the pose annotation could improve the results of the 2D baselines. We extend the ResNet50 model with a second classifier head to predict the pose following [67] and use the following combined loss:

$$\mathcal{L}_{pe-cl} = \mathcal{L}_{iCaRL} + \lambda_{pe}\mathcal{L}_{Starmap}. \quad (15)$$

We then train the models in a class-incremental fashion with the iCaRL [46] strategy. The results are provided in Table 7 and show that the additional pose supervision introduced in this way does not help to improve the classification accuracy. When increasing the weight of the pose loss λ_{pe} , the performance consistently decreases with the best performing model being the default iCaRL network with $\lambda_{pe} = 0$.

G Additional Implementational Details

In this section, we provide the full implementation details about our method.

Data Preparation. NeMo [53] was originally proposed for 3D pose estimation, which means that the degrees of freedom to the camera pose are azimuth, elevation, and roll angle. This implies that the objects are scaled accordingly and

Table 7: Adding an additional pose estimation head and providing additional supervision does not lead to better representation learning. It is not obvious how conventional classifiers could leverage the additional camera pose annotation.

Metric	λ_{pe}	$P3D$
		$B0 + 3$
	0.00	91.79
Classification	0.33	91.42
$\overline{acc}(1 : N_{task})$ in % \uparrow	0.66	90.56
	1.00	89.86

centered in the images. We follow this procedure and use the publicly available code of NeMo. To make the sizes of all input images consistent, we further pad all images to the size of 640×800 , where we fill the padded region with random textures from the Describable Textures Dataset [10].

Obtaining the 3D Cuboid Mesh is possible, since P3D [62] and O3D [61] provide a selection of 3D CAD models for each object category. For our 3D cuboid mesh, we consider the average bounding box of those models. We then sample vertices uniformly on its surface, leading to roughly 1,100 vertices per mesh.

Annotations at training time are computed with PyTorch3D’s [45] mesh rasterizer. Concretely, we render the neural meshes with ground-truth camera poses to compute their projection and binary object masks. Additionally, we compute the projected coordinates $\pi(k)$ of each vertex V^k and its binary visibility o^k . Given the class label, we render the corresponding mesh at $\frac{1}{8}$ of the original image resolution (same size as the output of the feature extractor Φ). At each pixel, we determine vertex visibility by considering the closest face using the returned z-buffer. To parameterize the rasterizer, we use a relatively simple camera model with a focal length of 1. As there is no viewport specified for neither the P3D or OOD-CV [65] dataset, we follow previous work [20, 53, 57] and use a viewport of 3000/8. For the O3D [61] dataset we use their specified viewport of 2000/8.

H Pose Estimation

For the pose estimation, we follow previous work [20, 53]. For completeness, we also provide a brief explanation here on how one can estimate the 3D object pose of an object of class c , given the trained feature extractor Φ and the neural mesh \mathfrak{M}_c .

3D Pose Estimation. During inference, we do not have access to the camera pose and corresponding perspective transformation. Since the camera intrinsics

and distance to the object are assumed to be known a-priori, we need to optimize for the unknown camera pose Q^{pred} . We define the foreground \mathcal{F} in the same way as we did for the classification in Section 4.4 of the main part of the paper. However, in addition to pixels that have been recognized as background, we also remove pixel positions that fall outside the projection of the cuboid, leading to $\mathcal{F}_{\pi^{\text{pred}}} = \mathcal{F} \cap \{f_{\pi^{\text{pred}}(k)} | \forall V_k \in \mathcal{V}_c \wedge o_c^k = 1\}$.

Finding Q^{pred} is done via a render-and-compare approach. We do so by initializing a rough estimate and optimizing iteratively. Given the current camera pose and its incurred perspective transform π^{pred} , we maximize the current object likelihood:

$$\max_{Q^{\text{pred}}} \prod_{f_{\pi^{\text{pred}}(k)} \in \mathcal{F}_{\pi^{\text{pred}}}} P(f_{\pi^{\text{pred}}(k)} | \theta_c^k). \quad (16)$$

By considering the vMF distribution, we optimize the initial camera pose using PyTorch3D’s [45] differentiable rasterizer by minimizing the negative log likelihood:

$$\mathcal{L}_{RC}(Q^{\text{pred}}) = \sum_{f_{\pi^{\text{pred}}(k)} \in \mathcal{F}_{\pi^{\text{pred}}}} -f_{\pi^{\text{pred}}(k)}^\top \cdot \theta_c^k. \quad (17)$$

Efficient Pose Estimation via Template Matching. The convergence of the above process is highly reliant on the provided initial pose, making it prohibitively slow in a worst case scenario. Wang *et al.* [55] proposed to speed it up by pre-rendering all neural meshes from 144 distinct viewing angles. Before the render-and-compare process, the output of the feature extractor is compared to each of these pre-rendered maps and the camera pose Q^{pred} is initialized with the pose that maximized the object likelihood. This simple procedure is remarkably effective, giving a speed-up of approximately 80% [20] over the original approach [53].

I Modelling the Background

For both classification and pose estimation, we leverage a set of features \mathcal{B} . This approach of disentangling foreground and background features into separate sets was introduced by Bai *et al.* [3]. Although we do not have a combined foreground set (but rather separate, 3D consistent meshes), we adopt their handling of the background model.

Learning the Background Model. We maintain a set \mathcal{B} of N_{bg} features that are sampled from positions in the feature map that fall outside the cuboid projection. From each sample in a training batch of size b , we sample N_{bgupdate} new background features. Consequently, we need to remove $b \cdot N_{\text{bgupdate}}$ from \mathcal{B} to avoid going over the allocated memory limit. Replacement is done, by maintaining a counter for each background feature, that indicates how many update steps it has been alive in \mathcal{B} and prioritizing the oldest features for removal.

Balancing the Background Model. Ideally, the background should contain features from a wide variety of background options (*i.e.* water from boats, sky from airplanes, urban scenes from cars, ...). However, sampling background features from the current task dataset only means that \mathcal{B} would be heavily biased towards background features from the currently considered classes. Therefore, we balance \mathcal{B} after each training phase by sampling background features from the exemplar memory $\mathcal{E}^{1:i}$, which was constructed evenly from all classes and viewing angles.

J Hyperparameter Collection

As there are quite a few hyperparameters involved in our method, we include this brief section that notes down all parameters for our final model.

Opt.	LR	(β_1, β_2)	Task-Epoch	Batch Size	Weight Decay
Adam	1e-5	(0.9, 0.999)	50	16	1e-4

Table 8: Optimization Parameters.

Table 9: Loss Weighting.

κ_1	κ_2	κ_3	λ_{eff}	λ_{kd}
1/0.07	1	0.5	0.1	10.0

Table 10: Mesh- and Background-related Parameters.

d	η	N_{bg}	$N_{bgupdate}$	R
128	0.9	2560	5	48

Table 11: Pose Estimation Parameters.

Opt.	LR	(β_1, β_2)	Epochs
Adam	5e-2	(0.4, 0.6)	30