

DRESS: Disentangled Representation-based Self-Supervised Meta-Learning for Diverse Tasks

Wei Cui Tongzi Wu Jesse C. Cresswell Yi Sui Keyvan Golestan

Layer 6 AI, Toronto, Canada

{wei, tongzi, jesse, amy, keyvan}@layer6.ai

Abstract

Meta-learning represents a strong class of approaches for solving few-shot learning tasks. Nonetheless, recent research suggests that simply pre-training a generic encoder can potentially surpass meta-learning algorithms. In this paper, we first discuss the reasons why meta-learning fails to stand out in these few-shot learning experiments, and hypothesize that it is due to the few-shot learning tasks lacking diversity. We propose DRESS, a task-agnostic Disentangled REpresentation-based Self-Supervised meta-learning approach that enables fast model adaptation on highly diversified few-shot learning tasks. Specifically, DRESS utilizes disentangled representation learning to create self-supervised tasks that can fuel the meta-training process. Furthermore, we also propose a class-partition based metric for quantifying the task diversity directly on the input space. We validate the effectiveness of DRESS through experiments on datasets with multiple factors of variation and varying complexity. The results suggest that DRESS is able to outperform competing methods on the majority of the datasets and task setups. Through this paper, we advocate for a re-examination of proper setups for task adaptation studies, and aim to reignite interest in the potential of meta-learning for solving few-shot learning tasks via disentangled representations.

1. Introduction

Few-shot learning [43] emphasizes the ability to quickly learn and adapt to new tasks, and is regarded as one of the trademarks of human intelligence. In the pursuit of few-shot learning, meta-learning approaches have been widely explored [9, 33, 37], as they allow models to *learn-to-learn*. However, multiple recent studies [8, 34, 35, 40] suggest that a simple *pre-training and fine-tuning* approach is sufficient to support highly competitive performance in few-shot learning tasks. Specifically, an encoder can be trained on a unified dataset that aggregates samples from all available training tasks. During inference, a linear layer is added on

top of the encoder and is fine-tuned using the few-shot support samples to adapt to new tasks.

Compared with meta-learning approaches, pre-training and fine-tuning neglects two crucial sources of information: identities of individual meta-training tasks and distinctions between them; and labels in meta-training tasks.¹ Despite completely ignoring these two sources of information, pre-training and fine-tuning has been shown to achieve better results than meta-learning [8, 35, 40]. This finding is unexpected, and perhaps even puzzling, as it implies that information about training tasks and their labels may be irrelevant to achieving high few-shot learning performance.

We hypothesize that this observation can be attributed to the lack of *task diversity* in many popular few-shot learning benchmarks. For instance, in canonical few-shot learning datasets such as Omniglot [23], miniImageNet [41], and CIFAR-FS [4], the distinct tasks differ solely in that their targets belong to distinct, non-overlapping sets of object classes. In essence, these few-shot learning tasks all share the same nature: main object classification. Hence, there is one degenerate strategy for solving all these tasks simultaneously without the need for individually identifying each task or relying on meta-training labels: compare the main object in the query image to the main objects in the few-shot support images, and assign the class label based on similarity to support images. This strategy can be easily achieved through pre-training with contrastive learning using common image augmentations like rotation and cropping which preserve the semantics of the main object, while discarding other factors such as orientation and background [3]. Given the shared nature of tasks on these specific benchmarks, it is not surprising that a single pre-trained encoder can perform competitively against meta-learning methods.

To rigorously challenge a model’s adaptation ability achieved by either meta-learning or pre-training and fine-tuning, we advocate for the establishment of few-shot learn-

¹When the pre-training stage utilizes a self-supervised loss to train the encoder, all the labels from meta-training tasks are discarded.

ing benchmarks that include tasks with fundamentally distinctive natures. Specifically, we consider tasks beyond main object classification, such as identifying object orientation, background color, ambient lighting, or attributes of secondary objects in the image. At the same time, models should be *agnostic* to the nature of the evaluation tasks. Such setups can reveal the model’s true capacity to learn strictly from the few-shot samples, with *task identification* expected to be an essential component.

Furthermore, we highlight a key consequence of high task diversity: when meta-testing tasks differ significantly in nature from meta-training tasks, the labels in the meta-training tasks may provide misleading guidance to the model. While the motivation of most research on unsupervised and self-supervised meta-learning is to avoid the cost of acquiring labels, the potential for labels to provide misleading guidance under high task diversity has not been discussed. Recognizing this issue, we reaffirm the preference of *self-supervised* meta-training over supervised meta-learning, as self-supervised meta-training can prevent premature fixation on a narrow perspective of the input data.

For effective meta-learning under high task diversities, we bridge the idea of disentangled representation learning with self-supervised meta-learning in a single framework we call *DRESS* — *task-agnostic Disentangled REpresentation-based Self-Supervised meta-learning*. Specifically, we employ an encoder trained to compute disentangled representations, and use it to extract latent encodings of the inputs. We then semantically align these latent representations across all inputs. Within this aligned latent space, we perform clustering independently on each disentangled latent dimension, and use the resultant cluster identities to define pseudo-classes of the inputs. Finally, we construct a set of self-supervised few-shot classification tasks based on these pseudo-classes from each latent dimension. With the disentangled latent dimensions representing distinct attributes and factors of variation within the input images, the constructed few-shot learning tasks are highly diversified. Using these tasks for meta-training, the model can learn to quickly adapt to various unseen tasks, regardless of the task nature.

In addition, to better investigate task diversity, we propose a quantitative task diversity metric based on class partitions. Our metric is directly defined on the input space instead of any learned embedding space, therefore allowing fair and independent comparisons between tasks of varying semantic natures.

We conduct extensive experiments on image datasets containing multiple factors of variation, beyond the main object’s class, and spanning different levels of complexity and realism. To ground our results, we establish three supervised meta-learning baselines that have differing levels of ground-truth information. These supervised baselines not

only serve as upper bounds on performance, but also expose the negative effects of learning from labels when the natures of tasks are mismatched. Our results suggest that DRESS enables few-shot learning performance that can surpass existing methods, and approaches the upper bound of supervised baselines under many experimental setups. Specifically, out of eight testing scenarios, DRESS outperforms all the competing methods on six of them, and outperforms the supervised meta-learning baselines on seven.

Our main contributions can be summarized as follows:

- We identify the lack of task diversity in few-shot learning benchmarks, explaining why pre-training and fine-tuning can outperform meta-learning. We develop more diverse benchmarks for rigorous evaluation.
- We propose DRESS², a self-supervised meta-learning method using disentangled representations for fast adaptation to diverse tasks.
- We introduce a task diversity metric based on class partitions, directly computed in the input space.

2. Related Works

Meta-Learning vs. Pre-training and Fine-tuning There has been a large volume of meta-learning research aiming to solve the general few-shot learning problem [9, 20, 26, 33, 37, 38]. Among them, one canonical and flexible meta-learning approach that stands out is MAML [9], which optimizes the initialization values for a neural network as the meta-parameters.

Researchers have also been interested in the possibilities of unsupervised or self-supervised meta-learning [2, 13, 15, 17, 18, 25, 32, 45]. Notably, CACTUS [13] proposes an unsupervised task construction approach using an encoding-then-clustering procedure. Meta-GMVAE [25] models the dataset using a variational auto-encoder [21] with a mixture of Gaussians as prior, and matches the latent modalities with class concepts. Meanwhile, research including [15, 17] rely on image augmentations and contrastive learning to create samples for pseudo classes, which are used to meta-train the model. Although promising results are obtained on standard few-shot learning benchmarks, these works do not explicitly address the issue of task diversity, nor its effect on fast adaptation performance. The approach we propose uses clustering for task construction, similar to CACTUS. However, to explicitly train the model to learn from diverse input attributes and to adapt to distinct tasks, we leverage the expressiveness of disentangled representations and create tasks focusing on different facets of the input dataset.

Recent studies [8, 34, 35, 40] state that the simple approach of pre-training a generic encoder followed by fine-tuning a projection layer on top can show superior perfor-

²The implementation of DRESS is available at: <https://github.com/layer6ai-labs/DRESS>.

mance compared to meta-learning. Specifically, the input samples from all available meta-training tasks are aggregated into a large dataset, with task identities completely ignored. An encoder is then trained on this large dataset using supervised or self-supervised training techniques (e.g., contrastive learning [7]). When adapting to any meta-testing task, a linear classification layer is added and fine-tuned on top of the encoder over the support samples.

Task Diversity The diversity among tasks used in the meta-training stage is imperative for the model’s fast adaptation ability, and has been considered previously [1, 22, 29, 39, 48]. One obstacle to investigating task diversity lies in the difficulty of quantifying it, which has been approached through metrics or proxies for quantifying the diversity between different tasks [1, 22, 39]. However, definitions of task diversity have either relied on a shared input embedding space [1, 22], or been defined through projection mappings from the input to output spaces [39].

Recently, [30] conducted thorough experiments suggesting that existing meta-learning methods show very slight improvements over the pre-training and fine-tuning approach on tasks with higher *Task2Vec* diversity coefficients [1, 29]. Nonetheless, to the best of our knowledge, the intuition behind the link between task diversity and the performance of few-shot learning has yet to be discussed. Similarly, no meta-learning approach has explicitly exploited the idea of diversifying meta-training tasks for boosting the fast adaptation ability of a model.

Disentangled Representation Learning Disentangled representation learning has been mainly investigated in the context of generative modeling [12, 14, 16, 19, 36, 44, 46, 47], with the objective of learning representations that capture independent factors of variation within the input distribution. For complex images, factors of variations include the main object identity, as well as object orientation, background, ambient lighting, view angle, and so on. We leverage disentangled representations when creating tasks for self-supervised meta-learning.

3. Problem Formulation

3.1. Definition of Few-Shot Learning

Consider an N -way K -shot classification task T , with inputs belonging to a set of classes $\{c_1, c_2, \dots, c_N\}$. For model adaptation, there are K labeled *support samples* available for each class: $T = \bigcup_{n=1}^N \{(x_i^s, y_i^s = c_n)_{i=1}^K\}$. At test time, the model needs to classify K_q unlabeled *query samples* $\{x_i^q\}_{i=1}^{K_q}$ drawn from the same distribution as the support samples.

In general, a *set* of few-shot learning tasks are available for both training $\{T_j^{\text{train}}\}_{j=1}^{N_{\text{train}}}$ and testing $\{T_j^{\text{test}}\}_{j=1}^{N_{\text{test}}}$, where N_{train} and N_{test} stand for the number of few-shot learning tasks for training and testing respectively.

3.2. Meta-Learning

In meta-learning, individual tasks serve a similar role as individual data points in conventional machine learning. Specifically, $\{T_j^{\text{train}}\}_{j=1}^{N_{\text{train}}}$ and $\{T_j^{\text{test}}\}_{j=1}^{N_{\text{test}}}$ are commonly referred to as the *meta-training tasks* and *meta-testing tasks*, respectively. Now, consider a parameterized machine learning model $f(\phi, \theta)$, with ϕ being the collection of meta-parameters shared across different tasks, and θ being the task-specific parameters to be adapted on each individual task. The general meta-training process optimizes the model’s meta-parameters over the meta-training tasks:

$$\phi^* = \arg \min_{\phi} \sum_{j=1}^{N_{\text{train}}} \mathcal{L}_{T_j^{\text{train}}}(f(\phi, \theta)), \quad (1)$$

where \mathcal{L}_T denotes the loss function of the model on task T (i.e., for a few-shot learning task, the loss function can be the cross-entropy loss on the query sample predictions).

After the meta-training stage, the model can be adapted to meta-testing tasks by only tuning the task-specific parameters. Specifically, for any given task T_j^{test} :

$$\theta_{T_j^{\text{test}}}^* = \arg \min_{\theta} \mathcal{L}_{T_j^{\text{test}}}(f(\phi^*, \theta)). \quad (2)$$

After adaptation, the model’s performance is computed as:

$$\sum_{j=1}^{N_{\text{test}}} \mathcal{L}_{T_j^{\text{test}}}(f(\phi^*, \theta_{T_j^{\text{test}}}^*)). \quad (3)$$

Meta learning has been one of the main forces for solving few-shot learning problems as described in Sec. 3.1. In one of the most popular meta-learning approaches, MAML [9], ϕ is the collection of initialization values for all the parameters in the model; while θ represents the set of finalized values of these model parameters after being optimized on a given task.

3.3. Pre-training and Fine-tuning

A more brute force, yet still general scheme for few-shot learning is pre-training and fine-tuning. It replaces the meta-training stage by *pre-training*, where all the data points from $\{T_j^{\text{train}}\}$ are aggregated into a large dataset, hence ignoring individual task identities. The model $f(\phi, \theta)$ is first trained on this aggregated dataset to optimize ϕ , and is then fine-tuned to adapt to testing tasks through optimizing θ . The conventional choices of ϕ and θ in the pre-training and fine-tuning scheme are often different from meta-learning. For example, in classification tasks ϕ is often the collection of parameters of an encoder that maps the inputs to a general-purpose embedding space, while θ is the collection of parameters in an adaptation layer that maps the embedding space to the outputs.

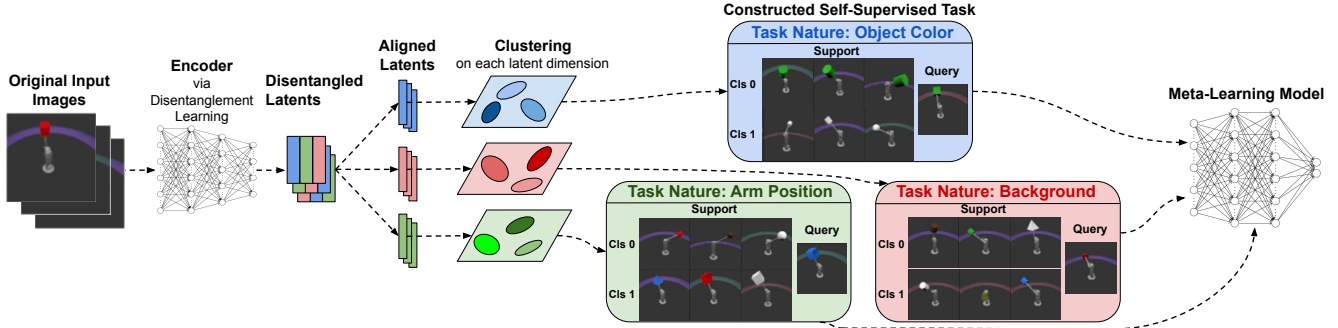


Figure 1. DRESS creates diversified self-supervised meta-training tasks through disentanglement learning. Images are first encoded into disentangled latent representations. The latent representations are then semantically aligned across the dataset so that sets of clusters can be formed on each latent dimension individually. Each set of clusters acts as pseudo-classes for a distinct self-supervised classification task. Hence, each disentangled latent dimension creates a meta-learning task with its own unique nature.

4. Methodology

We introduce DRESS, our task-agnostic Disentangled REpresentation-based Self-Supervised meta-learning approach. DRESS leverages disentangled latent representations of input images to construct self-supervised few-shot learning tasks that power the meta-training process. Following DRESS, the trained model is able to adapt to unseen tasks of diverse natures that require distinct learning rules to solve. The multi-stage flow diagram of DRESS is illustrated in Fig. 1, with details of each stage as follows:

1. **Encoding Disentangled Representations:** First, all images available for meta-training are collected, and used to train a general purpose encoder with the objective of producing disentangled representations (e.g., a factorized diffusion autoencoder (FDAE) [44], or latent slot diffusion model (LSD) [16]). We then use the trained encoder to encode the images to obtain their disentangled latent representations.
2. **Aligning Latent Dimensions:** After collecting the disentangled representations for all the training images, we align the latent dimensions of representations across images, corresponding to the semantic meaning of each latent dimension. For instance, some encoders [16, 28, 36] disentangle attributes by applying multiple attention masks over each input image, so we can align the resultant features by aligning the attention masks. After alignment, a given dimension in the latent space conveys the same semantic information across all images (e.g., main object color, background color, lighting color).
3. **Clustering Along Disentangled Latent Dimensions:** We perform clustering on each dimension over the latent values. Since dimensions are disentangled and aligned, clustering each latent dimension produces a distinct partition of the entire set of inputs that corresponds to one specific semantic property.
4. **Forming Diverse Self-Supervised Tasks:** Finally, we construct self-supervised learning tasks using cluster identities as the *pseudo-class* labels. We create a large

number of few-shot classification tasks under each disentangled latent dimension by first sampling a subset of cluster identities as classes, and then sampling a subset of images under each class as the few-shot support samples and query samples.

As different dimensions within the disentangled representation depict distinct aspects of the input data, the sets of self-supervised tasks constructed from disentangled dimensions are naturally diversified, requiring distinct decision rules to solve. When using these tasks for meta-training, the model can digest each factor of variation within the data, and therefore learns to adapt to unseen few-shot tasks regardless of their contexts, natures, and meanings. See Fig. 5 in Appendix A.1 for a general illustration of the meta-learning pipeline.

4.1. Selection of the Meta-Learning Algorithm

As a method to construct self-supervised tasks, DRESS is compatible with any conventional meta-learning algorithm for model training. However, not all meta-learning algorithms are well-suited to the highly diversified tasks DRESS generates. In this paper, we pair DRESS with the optimization-based adaptation approach MAML [9] because of its simplicity and ubiquity in meta-learning benchmarks. Discussions of other popular meta-learning algorithms are provided in Appendix A.2.

Nonetheless, we do not claim that MAML is necessarily the optimal choice for DRESS. We acknowledge the importance of providing different bases for meta-learning models to learn distinct semantics within disentangled latent dimensions. Therefore, we conjecture that *ensemble learning*, such as *mixture of experts*, holds significant potential for meta-learning in high task diversity setups and could be seamlessly integrated with DRESS to form a high-performing self-supervised learning pipeline. However, exploring ensemble learning as a meta-learning algorithm is beyond the scope of this paper.

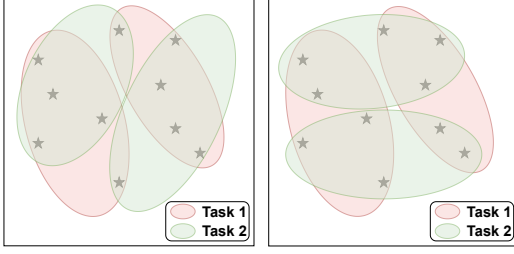


Figure 2. Illustration of class-partition based task diversity. A binary classification task is defined by two ellipses of the same color on an input space. **Left:** Two similar tasks where classes have high overlap in terms of data points. **Right:** Two dissimilar tasks, with less overlap between the class partitions.

4.2. Task Diversity based on Class Partitions

In DRESS, different encoders with different embedding spaces could be used to construct tasks. Correspondingly, we advocate for a task diversity metric that is not tied to any specific embedding space, but is instead linked to the original input space. Specifically, we introduce a task diversity metric based on the task’s class partitions. Consider two classification tasks defined on the same inputs as in Figure 2. Each task partitions the dataset based on class identities. The similarity between the two tasks can be measured by the similarity between their respective partitions.

The mathematical definition of our class-partition-based task diversity metric is as follows: consider an input dataset of K data points, $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^K$, and two (potentially multi-class) classification tasks, T_1 and T_2 , defined on \mathcal{D} . Assume T_1 and T_2 both have N classes which can be mapped to $\{c_j\}_{j=1}^N$ (if one task has fewer classes, we treat the missing classes as having zero samples). T_1 and T_2 can be described by two sets of class labels $\{y_i^1\}_{i=1}^K$ and $\{y_i^2\}_{i=1}^K$, respectively, with one label for each input in \mathcal{D} . Equivalently, each task can be represented by a class-based partition of \mathcal{D} . For T_1 , the class partition is denoted as $\mathcal{P}^1 = \{\mathcal{P}_{c_j}^1\}_{j=1}^N$, where $\mathcal{P}_{c_j}^1 = \{x_i \mid y_i^1 = c_j\}$. Similarly, \mathcal{P}^2 represents the class partition for T_2 .

Our task diversity metric is computed using these class-based partitions. First, we match subsets between \mathcal{P}^1 and \mathcal{P}^2 to maximize the pairwise overlaps, which can be achieved using methods such as bipartite matching. For each matched pair of subsets, we then compute the intersection-over-union (IoU) ratio. Finally, we calculate the average IoU value across all subset pairs across the two partitions. A low average IoU indicates that \mathcal{P}^1 and \mathcal{P}^2 differ significantly, suggesting that T_1 and T_2 are relatively diverse tasks. We note that during the step of relabeling the classes, the semantic information of the class concepts in each task is lost. Therefore, the proposed metric only quantifies task diversity from the function mapping perspective. Nonetheless, learning to jointly solve tasks that are diversi-

fied in their function mappings (as measured by our metric) has been shown to promote the development of better adaptation capacity [39].

5. Experimental Setup

5.1. Datasets

To study the effect of highly diversified tasks in different scenarios, we consider well-curated datasets with controlled factors of variations, as well as a complex real-world dataset. Specifically, for curated datasets, we consider SmallNORB [24], Shapes3D [5], Causal3D [42], and MPI3D [10], covering a data-complexity spectrum from easy to hard. These datasets all include labels for multiple independently varying factors, and are often used in research on generative modeling and disentangled representation learning. Furthermore, for a complex real-world dataset, we explore the popular CelebA dataset [27]. We provide the full details of all the factors of variation for each of these datasets in Appendix B. We emphasize that we do not include experiments on popular few-shot learning benchmarks such as Omniglot, *mini*ImageNet, and CIFAR-FS, due to their lack of task diversity among meta-training and meta-testing stages, as mentioned in Sec. 1.

5.2. Baseline Methods

We compare DRESS to several popular baseline methods on few-shot learning, categorized as follows:

Supervised Meta-Learning: We implement three variations of supervised meta-learning baselines with increasingly relevant information about ground-truth factors:

- *Supervised-Original:* Only use the ground-truth factors that do not define meta-testing tasks to create supervised meta-training tasks.
- *Supervised-All:* Use all the ground-truth factors from a dataset to create supervised meta-training tasks.
- *Supervised-Oracle:* Only use the ground-truth factors that define the meta-testing tasks to create supervised meta-training tasks.

These methods progressively increase the relevancy of information available to the model, but are increasingly unrealistic for practical settings. Supervised-Original must learn to generalize from a limited set of ground-truth factors to unknown factors at meta-testing time. Supervised-All has the most information, but needs to identify the task natures and relevant factors for successful adaptation, therefore representing the performance upper bound when the evaluation tasks are agnostic. Supervised-Oracle has perfect knowledge of factors utilized in meta-testing tasks, and represents the ultimate performance upper bound.

Few-Shot Direct Adaptation: This represents the lower bound of performance when a model is directly optimized on the support samples from each meta-testing task.

Table 1. Few-shot learning classification accuracies, with each trial conducted over 1000 meta-testing few-shot learning tasks (FSDA: Few-shot Direct Adaptation, PTFT: Pre-Training & Fine-Tuning).

Method	SmallNORB	Shapes3D	Causal3D	MPI3D-Easy	MPI3D-Hard
Supervised-Original	61.87% \pm 0.80%	62.03% \pm 1.55%	52.11% \pm 0.28%	57.75% \pm 0.46%	63.27% \pm 1.25%
Supervised-All	79.56% \pm 0.27%	99.93% \pm 0.02%	88.77% \pm 1.00%	99.29% \pm 0.29%	91.03% \pm 1.70%
Supervised-Oracle	80.22% \pm 0.41%	99.97% \pm 0.02%	93.47% \pm 0.17%	100.00% \pm 0.00%	99.42% \pm 0.11%
FSDA	73.93% \pm 0.91%	65.70% \pm 2.05%	66.92% \pm 0.86%	60.59% \pm 0.29%	62.27% \pm 0.28%
PTFT	57.97% \pm 1.87%	57.88% \pm 2.19%	55.61% \pm 0.21%	92.93% \pm 0.48%	79.50% \pm 0.76%
Meta-GMVAE	68.60% \pm 0.73%	59.10% \pm 1.73%	59.18% \pm 0.79%	99.39% \pm 0.13%	50.02% \pm 0.26%
PsCo	74.18% \pm 0.42%	97.62% \pm 0.58%	70.77% \pm 0.50%	83.52% \pm 2.01%	79.52% \pm 0.74%
CACTUS-DeepCluster	75.78% \pm 0.36%	86.81% \pm 0.68%	65.66% \pm 0.41%	84.95% \pm 0.56%	72.77% \pm 0.97%
CACTUS-DINOv2	62.76% \pm 0.76%	80.62% \pm 0.25%	53.87% \pm 0.54%	94.39% \pm 0.44%	81.92% \pm 0.39%
DRESS	78.13% \pm 0.36%	93.05% \pm 0.18%	76.42% \pm 0.38%	99.94% \pm 0.03%	84.95% \pm 0.50%

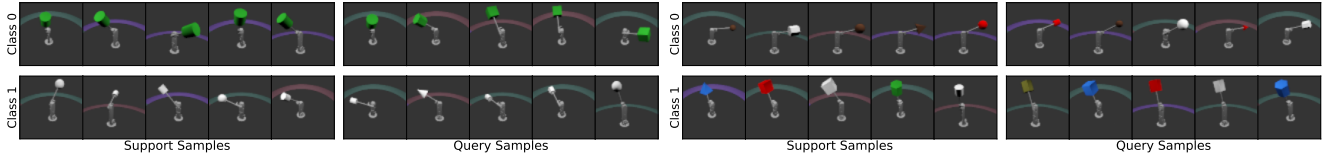


Figure 3. Two self-supervised tasks constructed by DRESS on MPI3D. **Left:** The task focuses on classifying the object color. **Right:** The task focuses on identifying the robot arm angle.

Pre-training and Fine-tuning: We implement the encoder-based pre-training and fine-tuning method as described in [40], using SimCLR [7] with its standard image augmentation pipeline for pre-training, with details in Appendix C. **Unsupervised & Self-Supervised Meta-Learning:** For unsupervised meta-learning baselines, we explore CACTUS [13] with two encoders: DeepCluster [6] trained from scratch on each dataset, and off-the-shelf DINOv2 [31]. We refer to these two baselines as *CACTUS-DeepCluster* and *CACTUS-DINOv2*, with detailed settings in Tab. 11 in Appendix D. Additionally, we experiment with two recent unsupervised and self-supervised meta-learning approaches: *Meta-GMVAE* [25] and *PsCo* [15].

We unify the model architecture, meta-training, and meta-testing setups for these methods across all experiments, as detailed in Appendix D.

5.3. Task Setups for Meta-Training & Meta-Testing

For meta-testing, we construct supervised few-shot learning tasks based on the selection of a *subset* of the attributes with ground-truth labels from each dataset. Consequently, the natures and levels of difficulty of the tasks are determined by this subset of attributes. Given the subset of attributes selected, the meta-testing tasks are created using the ground-truth labels, following a procedure similar to [13]. First, we randomly pick a small number of attributes from the attribute subset, and define two distinct value combinations on the picked attributes. Images whose attributes match the first value combination are assigned to the positive class, while those matching the second combination are assigned to the negative class. To ensure a high

number of distinct tasks given limited number of attributes and image samples, we focus on 2-way 5-shot tasks, though our method and analysis are not restricted to this specific setup. Additionally, we include 5 query samples per class. Importantly, the model does not receive prior information about the class distribution of these query samples. For the three supervised meta-training baselines, we also create supervised meta-training tasks following the same procedure. The details on the subsets of attributes for supervised meta-training tasks and meta-testing tasks are provided in the corresponding tables in Appendix B for each dataset.

To thoroughly evaluate each method’s adaptation across different task setups, we create multiple configurations for the two most complex datasets, MPI3D and CelebA, by varying how attributes are grouped. For MPI3D, we define two few-shot learning setups:

- *MPI3D-Easy*: meta-testing tasks focusing on identifying the background and camera height attributes.
- *MPI3D-Hard*: meta-testing tasks focusing on predicting the horizontal and vertical robot arm angular positions (with an extra layer of complexity introduced from visual disturbance caused by varying camera heights).

For CelebA, we define three few-shot learning setups focusing on different task natures:

- *CelebA-Hair*: meta-testing tasks focusing on all attributes relevant to the person’s hair.
- *CelebA-Primary*: meta-testing tasks focusing on primary facial attributes or features.
- *CelebA-Random*: meta-testing tasks constructed from a randomly selected subset of attributes.

For the self-supervised meta-training tasks generated by

Table 2. Few-shot learning classification accuracies on the more complex and realistic CelebA dataset, with each trial conducted over 1000 meta-testing few-shot learning tasks (FSDA: Few-shot Direct Adaptation, PTFT: Pre-Training & Fine-Tuning).

Method	CelebA-Hair	CelebA-Primary	CelebA-Random
Supervised-Original	68.89% \pm 0.56%	76.98% \pm 1.10%	81.90% \pm 0.21%
Supervised-All	79.10% \pm 0.23%	88.07% \pm 0.25%	85.64% \pm 0.20%
Supervised-Oracle	87.84% \pm 0.30%	91.20% \pm 0.12%	90.73% \pm 0.15%
FSDA	63.28% \pm 0.25%	69.31% \pm 0.49%	57.74% \pm 0.44%
PTFT	59.57% \pm 0.26%	67.07% \pm 0.28%	65.12% \pm 0.34%
Meta-GMVAE	64.18% \pm 0.20%	67.87% \pm 0.29%	64.94% \pm 0.18%
PsCo	66.24% \pm 0.29%	65.96% \pm 0.55%	60.49% \pm 0.38%
CACTUS-DeepCluster	67.40% \pm 0.96%	71.42% \pm 0.14%	62.16% \pm 1.09%
CACTUS-DINOv2	69.37% \pm 0.19%	77.00% \pm 0.30%	74.39% \pm 0.27%
DRESS	73.83% \pm 0.14%	77.41% \pm 0.12%	68.28% \pm 0.51%

DRESS, as well as the applicable meta-learning baselines, we use the same task setup format: 2-way 5-shot classification, with 5 query samples per class. For the remaining baselines, we strictly follow their own training procedures.

5.4. Implementation Details of DRESS

Curated Datasets: For our experiments on SmallNORB, Shapes3D, Causal3D and MPI3D, we adopt the FDAE architecture [44] for the encoder. We train a FDAE model from scratch on each dataset and use it to encode the images into disentangled representations. The FDAE encoder computes a pair of codes for each visual concept, the content code and the mask code. We regard this pair of codes as two independent latent dimensions.³

When using FDAE as the encoder, no explicit computation is required for the latent alignment stage in DRESS. Since FDAE employs deterministic convolutional neural networks, each output head of the network computes a fixed semantic mapping from input images. As a result, the latent dimensions are inherently organized in a consistent semantic order. This allows us to proceed directly to clustering after encoding all images. For each latent dimension (*i.e.*, a vector representation per image), we first apply PCA for dimensionality reduction, then perform K-Means clustering with 200 clusters, which define the pseudo-classes for self-supervised meta-training tasks.

Real-World Dataset: For our CelebA experiments, we adopt the LSD encoder [16] (trained from scratch) instead of FDAE due to the latter’s capacity limitations in capturing detailed facial features. This also demonstrates DRESS’s flexibility in obtaining disentangled representations from various encoder architectures. The LSD encoder utilizes slot attention [28] to learn disentangled latent representations by computing visual *slots*, with each slot attending to different regions of the image through a learned attention

mask. However, because of the stochastic nature of slot attention, the order of the slots varies across images, requiring explicit latent alignment before clustering. To align the slots, we gather a batch of attention masks, cluster them with K-Means (with the number of clusters equal to the number of slots per image), and reorder the attention slots based on cluster identities of their corresponding attention masks.

6. Results & Analysis

6.1. Experimental Results on Curated Datasets

We present the few-shot classification accuracies in Tab. 1 for all the curated datasets.⁴ On these datasets, DRESS consistently achieves the best few-shot adaptation performance among unsupervised methods, with an exception on the Shapes3D dataset. The performance of Supervised-Original is unimpressive, indicating that meta-training targets could mislead a supervised model when adapting to highly diversified tasks as we discussed in Sec. 4. In contrast to [40], pre-training and fine-tuning is not on par with meta-learning approaches, due to the more challenging and diverse tasks we benchmark on. CACTUS shows varying results across datasets with different encoders, reflecting the importance of the latent representations in task construction. As DRESS uses disentangled representation learning to construct diversified pre-training tasks, it obtains superior results across these datasets and task setups. We provide visualizations of two tasks constructed by DRESS in Fig. 3, and visualizations of more tasks constructed by DRESS focusing on other factors within MPI3D in Appendix E.

6.2. Experimental Results on Real-World Dataset

We report few-shot classification accuracies on the three CelebA setups from Sec. 5.3 in Tab. 2. DRESS outperforms

³Our notion of *latent dimension* is based on the semantic meaning. For FDAE encodings, each code for a visual concept is a vector. In this case, one latent dimension corresponds to a vector space.

⁴The reported results show the mean and standard deviation over 4 trials with different seeds, a procedure we follow throughout the experiments in this paper.

Table 3. Ablation on Disentangled Representations, Latent Dimension Alignment, and Individual Dimension Clustering.

Method	Causal3D	CelebA-Hair	CelebA-Primary
DRESS	76.42% \pm 0.38%	73.83% \pm 0.14%	77.41% \pm 0.12%
DRESS w/o Disent. Represent.	54.02% \pm 0.37%	-	-
DRESS w/o Latent Dim. Align.	-	73.02% \pm 0.18%	-
DRESS w/o Ind. Dim. Cluster.	-	-	74.22% \pm 0.26%

Table 4. Task Diversity Score on Each Dataset. Higher score indicates greater task diversity.

Method	Shapes3D	MPI3D-Hard	SmallNORB	Causal3D	CelebA-Hair
Supervised-Original	0.97	0.95	0.95	0.99	0.88
Supervised-All	0.99	0.99	0.98	0.99	0.88
Supervised-Oracle	0.99	0.98	0.99	0.99	0.85
CACTUS-DeepCluster	0.80	0.79	0.68	0.88	0.92
CACTUS-DINOv2	0.61	0.58	0.58	0.63	0.74
DRESS	0.90	0.92	0.70	0.73	0.98



Figure 4. Two self-supervised tasks constructed by DRESS on CelebA. **Left:** The task focuses on identifying if the person wears eyeglasses or not. **Right:** The task focuses on identifying if the person has hair bangs (hair curtain covering the forehead) or not.

all unsupervised methods on CelebA-Hair, excelling at capturing secondary features (i.e. hair features) beyond primary facial attributes. It also ranks first on CelebA-Primary, slightly ahead of CACTUS-DINOv2. It must be noted that as DINOv2 is the state-of-art high capacity vision encoder, it is expected to capture the information from the main objects (i.e., the faces). On CelebA-Random, DRESS falls behind CACTUS-DINOv2 but remains superior to other baselines. This drop likely stems from the fact that disentangled representations struggle to model fine details like *bags under eyes* and *bushy eyebrows*, which is a common phenomenon we have observed after experimenting with different encoders trained with disentangled representation learning. Similar to Sec. 6.1, Supervised-Original performs poorly, showing that labels can misguide adaptation to unseen tasks. We provide visualizations of two tasks constructed by DRESS on the CelebA dataset in Fig. 4. More task visualizations of diverse facial attributes by DRESS on CelebA are provided in Appendix E.

6.3. Ablation Studies

We conduct ablation studies on each key design decision of DRESS. To control the number of experiments, each ablation study is conducted under one task setup from Sec. 5.3. See Tab. 3 for all the ablation results.

Disentangled Representations: We replace the disentanglement learning encoder (i.e., FDAE or LSD) with the state-of-the-art DINOv2 encoder, which does not focus on learning disentangled representations. After extracting la-

tent representations from DINOv2, we follow the remaining steps of DRESS. Similar to FDAE, DINOv2 produces latent representations with a consistent semantic ordering across inputs, making latent alignment unnecessary. We perform this ablation study on the Causal3D setup where we find disentangled representations are absolutely crucial for strong performance. Without disentanglement, the constructed self-supervised tasks do not correspond to well-defined features of the data, making them much less useful for meta-learning.

Latent Dimension Alignment: Here we combine DRESS with the LSD encoder, as it does not produce aligned latent dimensions naturally. For the ablation, we skip the process of clustering the attention masks and re-ordering the attention slots. We conduct this ablation study on the CelebA-Hair setup, where there is a small but noticeable drop in performance. Without alignment, the same feature may express different semantic concepts on different datapoints, again making the generated tasks not correspond neatly to individual features of the data.

Clustering within each Disentangled Latent Dimension: Instead of performing independent clustering on each dimension separately, we directly cluster the entire latent space to generate the partitions. We then apply the final stage of DRESS to create self-supervised tasks from the obtained partitions. This ablation study is conducted on the CelebA-Primary setup, where performance is significantly reduced. When clustering all dimensions together, it is possible that the clustering naturally relies on individual dis-

entangled features, but is not guaranteed. When multiple features are used, the generated tasks will no longer cleanly define separate factors of variation in the data.

6.4. Quantitative Results on Task Diversity

We compute the class-partition based task diversity, as proposed in Sec. 4.2, for DRESS and applicable baselines. Details for computing this metric are provided in Appendix F. Tab. 4 reports the task diversity scores for each dataset, showing that DRESS consistently produces more diverse tasks than the baseline task construction method CACTUS. For supervised meta-learning methods, the task diversity scores are computed on the partitions constructed as specified in Sec. 5.3. They serve as upper bounds on the task diversity from each dataset, as they leverage the knowledge of the ground-truth attributes or factors of variations, which are intrinsically diverse by definition.

7. Conclusion

We surfaced an issue in popular few-shot learning benchmarks: tasks are not diverse enough to truly test model adaptation ability. Instead, tasks with distinct natures can serve as more informative benchmarks. We proposed a self-supervised meta-learning approach that harnesses the expressiveness of disentangled representations to construct self-supervised tasks. Our approach enables models to acquire broad knowledge on underlying factors in a dataset, and quickly adapt to unseen tasks. Experimental results validate that our approach empowers the model in fast adaptation under high task diversity benchmarks.

References

- [1] Alessandro Achille, Michael Lam, Rahul Tewari, Avinash Ravichandran, Subhansu Maji, Charles Fowlkes, Stefano Soatto, and Pietro Perona. Task2vec: Task embedding for meta-learning. In *ICCV*, 2019. 3
- [2] Antreas Antoniou and Amos Storkey. Assume, augment and learn: Unsupervised few-shot meta-learning via random labels and data augmentation. *arXiv preprint arXiv:1902.09884*, 2019. 2
- [3] Randall Balestriero, Mark Ibrahim, Vlad Sobal, Ari Morcos, Shashank Shekhar, Tom Goldstein, Florian Bordes, Adrien Bardes, Gregoire Mialon, Yuandong Tian, Avi Schwarzschild, Andrew Gordon Wilson, Jonas Geiping, Quentin Garrido, Pierre Fernandez, Amir Bar, Hamed Pirsiavash, Yann LeCun, and Micah Goldblum. A cookbook of self-supervised learning. *arXiv:2304.12210*, 2023. 1
- [4] Luca Bertinetto, Joao F. Henriques, Philip Torr, and Andrea Vedaldi. Meta-learning with differentiable closed-form solvers. In *ICLR*, 2019. 1
- [5] Chris Burgess and Hyunjik Kim. 3D shapes dataset. <https://github.com/deepmind/3dshapes-dataset>, 2018. 5
- [6] Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep clustering for unsupervised learning of visual features. In *ECCV*, 2018. 6
- [7] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *Proceedings of the 37th International Conference on Machine Learning*, pages 1597–1607, 2020. 3, 6
- [8] Vincent Dumoulin, Neil Houlsby, Utku Evci, Xiaohua Zhai, Ross Goroshin, Sylvain Gelly, and Hugo Larochelle. A unified few-shot classification benchmark to compare transfer and meta learning approaches. In *NeurIPS*, 2021. 1, 2
- [9] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *Proceedings of the 34th International Conference on Machine Learning*, pages 1126–1135, 2017. 1, 2, 3, 4, 11, 14
- [10] Muhammad Waleed Gondal, Manuel Wuthrich, Djordje Miladinovic, Francesco Locatello, Martin Breidt, Valentin Volchkov, Joel Akpo, Olivier Bachem, Bernhard Schölkopf, and Stefan Bauer. On the transfer of inductive bias from simulation to the real world: a new disentanglement dataset. In *NeurIPS*, 2019. 5
- [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 14
- [12] Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-VAE: Learning basic visual concepts with a constrained variational framework. In *ICLR*, 2017. 3
- [13] Kyle Hsu, Sergey Levine, and Chelsea Finn. Unsupervised learning via meta-learning. In *ICLR*, 2019. 2, 6
- [14] Kyle Hsu, Jubayer Ibn Hamid, Kaylee Burns, Chelsea Finn, and Jiajun Wu. Tripod: Three complementary inductive biases for disentangled representation learning. In *Proceedings of the 41st International Conference on Machine Learning*, pages 19101–19122, 2024. 3
- [15] Huiwon Jang, Hankook Lee, and Jinwoo Shin. Unsupervised meta-learning via few-shot pseudo-supervised contrastive learning. In *ICLR*, 2023. 2, 6
- [16] Jindong Jiang, Fei Deng, Gautam Singh, and Sungjin Ahn. Object-centric slot diffusion. In *NeurIPS*, 2023. 3, 4, 7
- [17] Siavash Khodadadeh, Ladislav Bölöni, and Mubarak Shah. Unsupervised meta-learning for few-shot image classification. In *NeurIPS*, 2019. 2
- [18] Siavash Khodadadeh, Sharare Zehtabian, Saeed Vahidian, Weijia Wang, Bill Lin, and Ladislav Boloni. Unsupervised meta-learning through latent-space interpolation in generative models. In *ICLR*, 2021. 2
- [19] Hyunjik Kim and Andriy Mnih. Disentangling by factorising. In *Proceedings of the 35th International Conference on Machine Learning*, pages 2649–2658, 2018. 3
- [20] Minyoung Kim and Timothy M. Hospedales. A Hierarchical Bayesian Model for Few-Shot Meta Learning. In *ICLR*, 2024. 2
- [21] Diederik P Kingma and Max Welling. Auto-encoding variational Bayes. In *ICLR*, 2014. 2
- [22] Ramnath Kumar, Tristan Deleu, and Yoshua Bengio. The effect of diversity in meta-learning. In *AAAI*, 2022. 3

- [23] Brenden M. Lake, Ruslan Salakhutdinov, Jason Gross, and Joshua B. Tenenbaum. One shot learning of simple visual concepts. In *Conference of the Cognitive Science Society*, 2011. 1
- [24] Yann Lecun, Fufei Huang, and Leon Bottou. Learning methods for generic object recognition with invariance to pose and lighting. In *CVPR*, 2004. 5
- [25] Dong Bok Lee, Dongchan Min, Seanie Lee, and Sung Ju Hwang. Meta-GMVAE: Mixture of Gaussian VAEs for unsupervised meta-learning. In *ICLR*, 2021. 2, 6
- [26] Hae Beom Lee, Hayeon Lee, Jaewoong Shin, Eunho Yang, Timothy M. Hospedales, and Sung Ju Hwang. Online hyperparameter meta-learning with hypergradient distillation. In *ICLR*, 2022. 2
- [27] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *ICCV*, 2015. 5
- [28] Francesco Locatello, Dirk Weissenborn, Thomas Unterthiner, Aravindh Mahendran, Georg Heigold, Jakob Uszkoreit, Alexey Dosovitskiy, and Thomas Kipf. Object-centric learning with slot attention. In *NeurIPS*, 2020. 4, 7
- [29] Brando Miranda, Patrick Yu, Yu-Xiong Wang, and Sanmi Koyejo. The curse of low task diversity: On the failure of transfer learning to outperform maml and their empirical equivalence. *arXiv preprint arXiv:2208.01545*, 2022. 3
- [30] Brando Miranda, Patrick Yu, Saumya Goyal, Yu-Xiong Wang, and Sanmi Koyejo. Is pre-training truly better than meta-learning? *arXiv preprint arXiv:2306.13841*, 2023. 3
- [31] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel HAZIZA, Francisco Massa, Alaaeldin El-Nouby, Mido Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Herve Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. DINOv2: Learning robust visual features without supervision. *Transactions on Machine Learning Research*, 2024. 6
- [32] Eva Pachetti, Sotirios Tsaftaris, and Sara Colantonio. Boosting few-shot learning with disentangled self-supervised learning and meta-learning for medical image classification. *arXiv preprint arXiv:2403.17530*, 2024. 2
- [33] Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. In *ICLR*, 2017. 1, 2, 11
- [34] Zhiqiang Shen, Zechun Liu, Jie Qin, Marios Savvides, and Kwang-Ting Cheng. Partial is better than all: Revisiting fine-tuning strategy for few-shot learning. In *AAAI*, 2021. 1, 2
- [35] Yang Shu, Zhangjie Cao, Jinghan Gao, Jianmin Wang, Philip S. Yu, and Mingsheng Long. Omni-training: Bridging pre-training and meta-training for few-shot learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45:15275–15291, 2023. 1, 2
- [36] Gautam Singh, Fei Deng, and Sungjin Ahn. Illiterate DALL-E learns to compose. In *ICLR*, 2022. 3, 4
- [37] Jake Snell, Kevin Swersky, and Richard S. Zemel. Prototypical networks for few-shot learning. In *NeurIPS*, 2017. 1, 2, 11
- [38] Xiaozhuang Song, Shun Zheng, Wei Cao, James J.Q. Yu, and Jiang Bian. Efficient and effective multi-task grouping via meta learning on task combinations. In *NeurIPS*, 2022. 2
- [39] Yi Sui, Tongzi Wu, Jesse C. Cresswell, Ga Wu, George Stein, Xiao Shi Huang, Xiaochen Zhang, and Maksims Volkovs. Self-supervised representation learning from random data projectors. In *ICLR*, 2024. 3, 5
- [40] Yonglong Tian, Yue Wang, Dilip Krishnan, Joshua B. Tenenbaum, and Phillip Isola. Rethinking few-shot image classification: a good embedding is all you need? In *ECCV*, 2020. 1, 2, 6, 7
- [41] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Koray Kavukcuoglu, and Daan Wierstra. Matching networks for one shot learning. In *NeurIPS*, 2016. 1
- [42] Julius von Kügelgen, Yash Sharma, Luigi Gresele, Wieland Brendel, Bernhard Schölkopf, Michel Besserve, and Francesco Locatello. Self-supervised learning with data augmentations provably isolates content from style. In *NeurIPS*, 2021. 5
- [43] Yaqing Wang, Quanming Yao, Joern-Henrik Jacobsen, David Duvenaud, Mohammad Norouzi, and Kevin Swersky. Generalizing from a few examples: A survey on few-shot learning. *ACM Computing Surveys (CSUR)*, 53(3):1–34, 2020. 1
- [44] Ancong Wu and Wei-Shi Zheng. Factorized diffusion autoencoder for unsupervised disentangled representation learning. In *AAAI*, 2024. 3, 4, 7
- [45] Hui Xu, Jiaxing Wang, Hao Li, Deqiang Ouyang, and Jie Shao. Unsupervised meta-learning for few-shot learning. *Pattern Recognition*, 2021. 2
- [46] Tao Yang, Yuwang Wang, Yan Lu, and Nanning Zheng. Dis-Diff: Unsupervised disentanglement of diffusion probabilistic models. In *NeurIPS*, 2023. 3
- [47] Zhongqi Yue, Jiankun Wang, Qianru Sun, Lei Ji, Eric I-Chao Chang, and Hanwang Zhang. Exploring diffusion time-steps for unsupervised representation learning. In *ICLR*, 2024. 3
- [48] Amir R. Zamir, Alexander Sax, William Shen, Leonidas Guibas, Jitendra Malik, and Silvio Savarese. Taskonomy: Disentangling task transfer learning. In *CVPR*, 2018. 3

A. Visualization and Discussions of Meta-Learning Algorithms

A.1. Meta-Learning on Few-Shot Learning Pipeline

We provide the visualization for the general pipeline on applying meta-learning to solve few-shot learning tasks in Fig. 5.

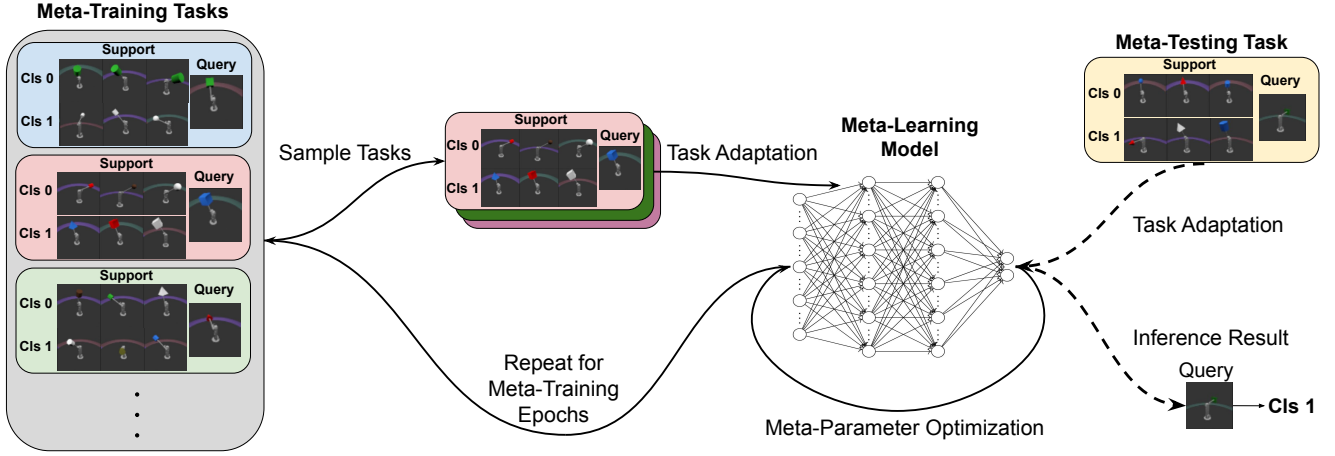


Figure 5. During the meta-training stage, the model adapts on batches of sampled tasks. The model’s performance is optimized for meta-parameter optimization. After meta-training, the model can be quickly adapted to meta-testing tasks and perform few-shot inference.

A.2. Discussions on Suitable Selection of Meta-Learning Algorithm

The majority of meta-learning algorithms can be categorized under one of the three general themes: black-box adaptation [33]; optimization-based adaptation, with MAML [9] being the notable example; and non-parametric adaptation, with ProtoNet [37] being the notable example.

The non-parametric adaptation scheme often relies on a pre-trained latent space, based on which the adaptation to new tasks is achieved (*e.g.*, the computation of the *prototypes* in ProtoNet). However, as we have advocated through the design of DRESS, we need different partitions on a given dataset based on disentangled latent dimensions to correspond to different semantics or nature of diverse tasks. Therefore, the non-parametric adaptation scheme lacks the capacity for fully benefiting from DRESS. We also opt out of the black-box adaptation scheme for its lack of inductive bias in the learning process, due to this same reason.

Optimization-based meta-learning algorithms are suitable to combine with DRESS for learning tasks with diverse natures. This class of algorithms does not impose the assumption that the model should adapt to all the tasks based on any specific latent space, therefore allowing the model the flexibility in learning different fundamental concepts and attributes from the data, and benefiting from the comprehensive set of meta-training tasks provided by DRESS.

B. Dataset Descriptions

B.1. SmallNORB

SmallNORB contains 48,600 images, of which we use 24,300 images for meta-training and 24,300 images for meta-testing, following the pre-defined train-test split convention on the dataset. Each image has a resolution of 96×96 pixels with a single gray-scale color channel. We simply repeat this channel three times to create three-channel images to be compatible with all of the encoders tested (such as the pre-trained DINOv2, which expects three-channel images as inputs off-the-shelf).

The images in the dataset include 5 factors of variations, as detailed in Tab. 5. Note that we ignored the additional factor of *camera ID* in SmallNORB, as we exclusively take images from the first camera.

B.2. Shapes3D

Shapes3D contains 480,000 images, of which we use 400,000 images for meta-training and 50,000 images for meta-testing, following the pre-defined train-test split convention on the dataset. Each image has a resolution of 64×64 pixels with RGB color channels.

The images in the dataset include 6 factors of variations, as detailed in Tab. 6.

Table 5. Factors of Variation in SmallNORB

Attribute Name	Cardinality	Constructed Tasks
Generic Category	5	Meta-Train
Instance ID	5	Meta-Train
Elevation Angle	18	Meta-Test
Azimuth Angle	9	Meta-Test
Lighting	6	Meta-Test

Table 6. Factors of Variation in Shapes3D

Attribute Name	Cardinality	Constructed Tasks
Floor Hue	10	Meta-Test
Wall Hue	10	Meta-Test
Object Hue	10	Meta-Train
Scale	8	Meta-Train
Shape	4	Meta-Train
Orientation	15	Meta-Test

B.3. Causal3D

Causal3D contains 237,600 images, of which we use 216,000 images for meta-training and 21,600 images for meta-testing, following the pre-defined train-test split convention on the dataset. Each image has a resolution of 224×224 pixels with RGB color channels.

The images in the dataset include 7 factors of variations, as detailed in Tab. 7. Each of these factors are continuous values in the original form, which we have quantized to 10 levels. We emphasize that in DRESS and the competing unsupervised methods we experimented with, the models are agnostic to the quantization decision (i.e. there are 10 different values in each latent dimension that we use for creating meta-testing few-shot learning tasks). Note that the original dataset also provides labels for additional factors which we neglected in our experiments, such as rotation angles.

Table 7. Factors of Variation in Causal3D

Attribute Name	Cardinality	Constructed Tasks
X Position	10	Meta-Train
Y Position	10	Meta-Train
Z Position	10	Meta-Train
Object Color	10	Meta-Train
Ground Color	10	Meta-Test
Spotlight Position	10	Meta-Test
Spotlight Color	10	Meta-Test

B.4. MPI3D

MPI3D consists of four dataset variants. We utilize the *MPI3D_toy* dataset containing simplistic rendered images with clear color contrast. Throughout the paper, we refer to this dataset simply as MPI3D. The dataset contains 1,036,800 images, of which we use 1,000,000 images for meta-training and 30,000 images for meta-testing, following the pre-defined train-test split convention on the dataset. Each image has a resolution of 64×64 pixels with RGB color channels.

The images in the dataset include 7 factors of variations, as detailed in Tab. 8. We note that for the two factors *horizontal axis* and *vertical axis*, denoting the robot arm’s angular position, the ground truth labels for each are based on a 40-interval partition of the entire 180-degree angular range, leading to a mere 4.5-degree maximum angle difference for two different factor values. In our experiments, we re-group the partitions into 10 intervals for each of the two axes, leading to an 18-degree

maximum angle difference between two factor values.

Table 8. Factors of Variation in MPI3D under each Task Setup

Attribute Name	Cardinality	MPI3D-Easy Task Setup	MPI3D-Hard Task Setup
Object Color	6	Meta-Train	Meta-Train
Object Shape	6	Meta-Train	Meta-Train
Object Size	2	Meta-Train	Meta-Train
Camera Height	3	Not Used	Meta-Test
Background Color	3	Not Used	Meta-Test
Horizontal Axis	40	Meta-Test	Not Used
Vertical Axis	40	Meta-Test	Not Used

B.5. CelebA

CelebA consists of 202,599 images of celebrity faces, of which we use 162,770 images for meta-training and 19,962 images for meta-testing, leaving 19,687 images for meta-validation for a subset of approaches. Note that this is the conventional split when experimenting with CelebA. Each image has a resolution of 178×218 pixels with RGB color channels. We conduct a cropping around the face regions in these images before feeding them into each model, for both meta-training and meta-testing.

The images in the dataset include 40 binary factors of variations. Instead of listing out all these 40 factors, in Tab. 9, we only list the binary attributes reserved for meta-testing few-shot learning tasks under each attribute split setup. The remaining attributes were used for constructing meta-training tasks exclusively for supervised baselines.

Table 9. Factors of Variation in CelebA under each Task Setup

Task Setup	Attribute Name
CelebA-Hair	Bangs
	Black Hair
	Blond Hair
	Brown Hair
	Gray Hair
	Receding Hairline
	Straight Hair
	Wavy Hair
CelebA-Primary	Bald
	Big Lips
	Big Nose
	Blond Hair
	Eyeglasses
	Pale Skin
	Straight Hair
	Wearing Hat
CelebA-Random	5 o’Clock Shadow
	Bags under Eyes
	Bald
	Blurry
	Bushy Eyebrows
	Double Chin
	Goatee
	Mouth Slightly Open

C. Detailed Setups for Pre-training and Fine-tuning

For pre-training, we use an encoder backbone that shares the same architecture as the ResNet-18 [11] backbone used for FDAE. After pre-training, a trainable linear layer is attached on top of the encoder for the adaptation process on evaluation tasks. The encoder is frozen throughout the adaptation process. We include the details for this approach in Tab. 10. Note that we do not use a supervised loss in pre-training in order to avoid the encoder focusing only on tasks that are irrelevant to the meta-evaluation tasks, as we have discussed in Sec. 4.

Table 10. Pre-Training and Fine-Tuning Setup

Setting	Value
Pre-Training Epochs	10
Tasks in Meta-Evaluation	1000
Gradient Descent Steps in Adaptation	5

Regarding the number of epochs for pre-training, in the pre-training procedure the entire set of meta-training image inputs are fed to the encoder (i.e. 400,000 images for Shapes3D; and 1,000,000 images for MPI3D). Therefore, with 10 epochs over the entire meta-training dataset, the number of forward-backward computations for optimizing the encoder already surpasses the models trained with the meta-learning methods.

D. Additional Setup Details for Meta-Learning Methods

In this section, we further provide more details on the implementation of DRESS as well as meta-learning baselines.

Firstly, for DRESS, the supervised meta-learning baselines, as well as the two CACTUS baselines, we use MAML [9] as the meta-optimization engine, with a convolutional neural network (CNN) of identical specification as the base learner, for fair comparisons between the methods. The few-shot direct adaptation baseline also uses a CNN of the same specification. For the remaining baselines, we follow the design details as in the original papers.

We summarize in Tab. 11 hyper-parameter values of the meta-learning baselines CACTUS-DeepCluster and CACTUS-DINOv2. We note that for the DeepCluster encoder, PCA is applied on its output to reach the number of latent dimensions as listed.

Table 11. Few-Shot Learning Setup

Setting	CACTUS-DeepCluster	CACTUS-DINOv2
Latent Dimensions	256	384
Randomly Scaled Latent Spaces	50	50
Clusters Over Each Latent Space	300	300

In Tab. 12, we provide meta-training and meta-testing hyper-parameters for DRESS and two meta-learning baselines, CACTUS-DeepCluster and CACTUS-DINOv2.

E. Additional Task Visualizations from DRESS

We provide more visualizations on self-supervised few-shot learning tasks generated by DRESS on MPI3D in Fig. 6, as well as tasks generated by DRESS on CelebA in Fig. 7. As evidenced by these visualizations, the generated tasks have very distinctive natures, covering multiple aspects and factors of variations within the corresponding datasets. When being trained on such diversified tasks, the resulting model naturally acquires the ability to adapt well on unseen tasks, regardless of the semantics that the tasks focus on.

Table 12. Few-Shot Learning Setup for Meta-Learning Methods

Setting	Value
Tasks per Meta-Training Epoch	8
Meta-Training Epochs	30,000
Tasks in Meta-Evaluation	1,000
Gradient Descent Steps in Task Adaptation	5
Adaptation Step Learning Rate	0.05
Meta-Optimization Step Learning Rate	0.001

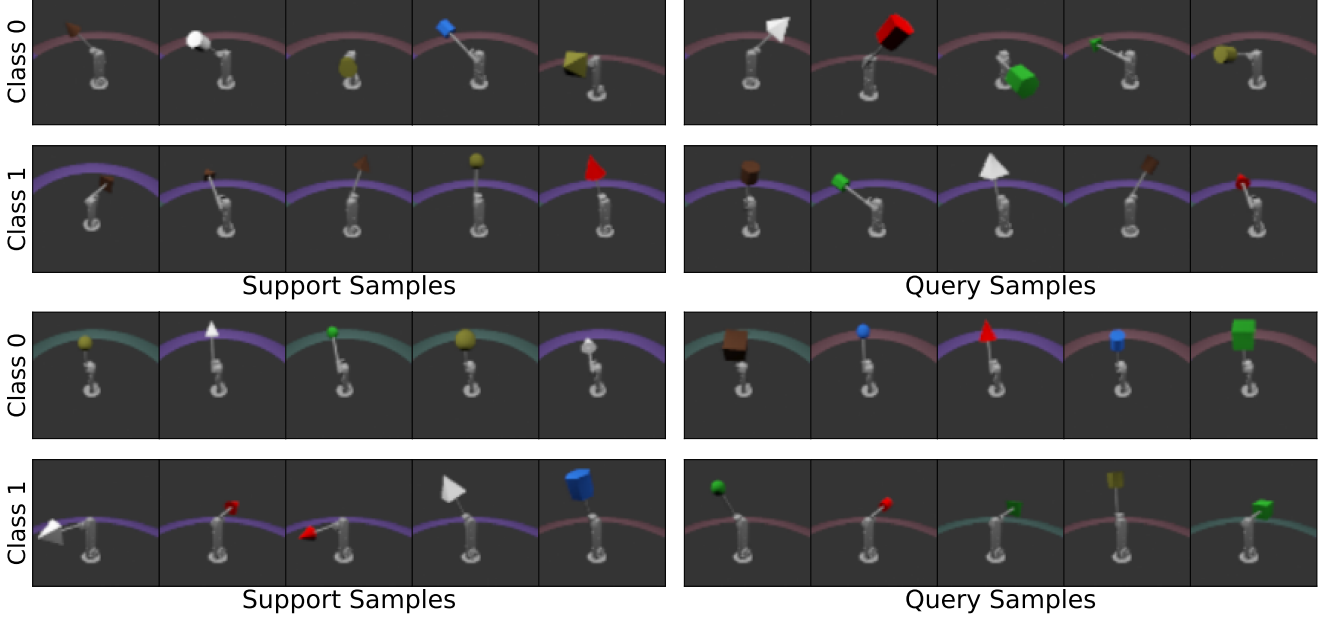


Figure 6. More self-supervised tasks constructed by DRESS on MPI3D. The top task focuses on the background color; while the bottom task focuses on the camera height.

F. Computation Details for Class-Partition based Task Diversity Metric

With the task diversity defined in Sec. 4.2, we aim to compute the intersection-over-union ratio (IoU) over pairs of tasks created by each method. Nonetheless, as we focus on the few-shot learning tasks (five-shot two-way tasks to be specific), the number of input samples on each task is very small. Therefore, if we directly take two such few-shot learning tasks, there is most likely no intersection in the samples they cover.

To address this difficulty, we instead focus on the partitions over the entire dataset. As described in Sec. 4 and Sec. 5.3, for DRESS, supervised meta-learning baselines, as well as the CACTUS-based baselines, the individual tasks are directly sampled from the dataset-level partitions. Therefore, computing the diversity metric over these partitions can give us a good proxy to the evaluation of the task diversity from each method. We now present the procedure for computing the class-partition based task diversity.

Consider two partitions on the same dataset generated by a specific meta-learning method: $\mathcal{P}^1 = \{\mathcal{P}_i^1\}_{i=1}^{K_p}$ and $\mathcal{P}^2 = \{\mathcal{P}_i^2\}_{i=1}^{K_p}$, where \mathcal{P}_i^1 and \mathcal{P}_i^2 denotes the i -th subset in \mathcal{P}^1 and \mathcal{P}^2 respectively, and K_p is the number of subsets in each partition. Note that if one partition has fewer subsets, we can simply regard it having extra empty subsets, such that the total number of subsets reaches K_p . We use these dataset partitions to replace the class partitions within each task, i.e. $\{\mathcal{P}_{c_j}^1\}$ and $\{\mathcal{P}_{c_j}^2\}$.

We summarize our procedure for computing the values on the purposed task diversity metric in Algorithm 1. Note that



Figure 7. More self-supervised tasks constructed by DRESS on CelebA. The top task focuses on the gender of the person; while the bottom task focuses on if the person has mouth open or not.

instead of performing strict bipartite matching for subsets between \mathcal{P}^1 and \mathcal{P}^2 , we match the subsets through a greedy process: going through the subsets one-by-one in \mathcal{P}^1 , and find the best match from the remaining subsets in \mathcal{P}^2 . While this greedy procedure does not strictly guarantee the perfect matches between the two partitions, it provides a decent estimates for our quantitative analysis at a manageable level of computational cost.

Algorithm 1 Task Diversity Metric Computation Procedure

Input: $\mathcal{P}^1 = \{\mathcal{P}_i^1\}_{i=1}^{K_p}$, $\mathcal{P}^2 = \{\mathcal{P}_i^2\}_{i=1}^{K_p}$
 $\text{idx_list} \leftarrow [1, 2, \dots, K_p]$
 $\text{IoU_list} \leftarrow \emptyset$
for $i \leftarrow 1$ to K_p **do**
 $\text{idx_matched} \leftarrow 0$
 $\text{highest_IoU} \leftarrow 0$
 for $j \in \text{idx_list}$ **do**
 $\text{IoU} = \frac{|\mathcal{P}_i^1 \cap \mathcal{P}_j^2|}{|\mathcal{P}_i^1 \cup \mathcal{P}_j^2|}$
 if $\text{IoU} > \text{highest_IoU}$ **then**
 $\text{idx_matched} \leftarrow j$
 $\text{highest_IoU} \leftarrow \text{IoU}$
 end if
 end for
 $\text{IoU_list.append}(\text{highest_IoU})$
 if $\text{idx_matched} > 0$ **then**
 $\text{idx_list.pop}(\text{idx_matched})$
 end if
end for
 $\text{avg_IoU_score} \leftarrow \text{avg}(\text{IoU_list})$
 $\text{diversity_score} \leftarrow 1 - \text{avg_IoU_score}$
Output: diversity_score
