

Flexible Few-Shot Learning of Contextual Similarity

Mengye Ren*

MREN@CS.TORONTO.EDU

University of Toronto; Vector Institute

Eleni Triantafillou*

ELENI@CS.TORONTO.EDU

University of Toronto; Vector Institute

Kuan-Chieh Wang*

WANGKUA1@CS.TORONTO.EDU

University of Toronto; Vector Institute

James Lucas*

JLUCAS@CS.TORONTO.EDU

University of Toronto; Vector Institute

Jake Snell

JSNELL@CS.TORONTO.EDU

University of Toronto; Vector Institute

Xaq Pitkow

XAQ@RICE.EDU

Rice University; Baylor College of Medicine

Andreas S. Tolias

ASTOLIAS@BCM.EDU

Baylor College of Medicine; Rice University

Richard Zemel

ZEMEL@CS.TORONTO.EDU

University of Toronto; Vector Institute; Canadian Institute for Advanced Research

Abstract

Existing approaches to few-shot learning deal with tasks that have persistent, rigid notions of classes. Typically, the learner observes data only from a fixed number of classes at training time and is asked to generalize to a new set of classes at test time. Two examples from the same class would always be assigned the same labels in any episode. In this work, we consider a realistic setting where the similarities between examples can change from episode to episode depending on the task context, which is not given to the learner. We define new benchmark datasets for this flexible few-shot scenario, where the tasks are based on images of faces (Celeb-A), shoes (Zappos50K), and general objects (ImageNet-with-Attributes). While classification baselines and episodic approaches learn representations that work well for standard few-shot learning, they suffer in our flexible tasks as novel similarity definitions arise during testing. We propose to build upon recent contrastive unsupervised learning techniques and use a combination of instance and class invariance learning, aiming to obtain general and flexible features. We find that our approach performs strongly on our new flexible few-shot learning benchmarks, demonstrating that unsupervised learning obtains more generalizable representations.

*Equal contribution

1. Introduction

Following the success of machine learning applied to fully-supervised settings, there has been a surge of interest in machine learning within more realistic, natural learning scenarios. Among these, few-shot learning (Lake et al., 2011) (FSL) has emerged as an exciting alternative paradigm. In the few-shot learning setting, the learner is presented with episodes of new learning tasks, where in each episode the learner must identify patterns in a small labeled support set and apply them to make predictions for an unlabeled query set. Since its inception, there has been significant progress on FSL benchmarks. However, standard supervised baselines are often shown to perform as well as carefully designed solutions (Chen et al., 2019; Tian et al., 2020). In this work, we argue that this observation is due in part to the rigidity in which FSL episodes are designed.

In a typical few-shot classification setting, each episode consists of a few examples belonging to one of N classes. Across different training episodes, different images are sampled from the classes in the training set but they will always be given the same class label: an elephant is always an elephant. But many real-world judgments are contextual—they depend on the task at hand and frame-of-reference. A rock is similar to a chair when the aim is to sit, but similar to a club if the aim is to hit. Few-shot learning is especially appropriate in contextual judgments, as people are able to adapt readily to new contexts and make appropriate judgments. So an important question is how to incorporate context into few-shot classification?

In this work, we define a new flexible few-shot learning (FFSL) paradigm. Instead of building episodes from classes, each episode is a binary classification problem that is constructed with some context that is hidden from the learner. In this way, the same data point may be given different labels across multiple episodes. For example, elephants and tables may belong to the same class if the context is “has legs”, but not when the context is “has ears”. Importantly, the learner is not given direct access to the context and must infer it from the examples present in the episode. Also, the contexts are novel at test-time.

Our FFSL problem is significantly more challenging than the standard setup, as it requires a flexible learner that is able to perform well in different contexts. We study generalization issues that occur under supervised representation learning for the flexible few-shot tasks, and we show that these approaches tend to overfit to the training attributes, even when given direct access to the attributes that determine the context. We analyze a simplified version of the problem in order to elucidate one possible cause of this failure.

We contribute new benchmark datasets for this flexible few-shot scenario. The tasks are based on images of faces (Celeb-A) (Liu et al., 2015), shoes (Zappos50K) (Yu and Grauman, 2014), and general objects (ImageNet-with-Attributes) (Deng et al., 2009). We provide a thorough empirical evaluation of existing methods on these tasks. We find that successful approaches in the standard FSL setting fall short on the flexible few-shot tasks. Further, while supervised classification baselines can learn good representation in the standard FSL setting, they suffer in FFSL. Finally, we propose to use a combination of instance and class invariance learning, aiming to obtain general and flexible features. We find that our approach performs strongly on our new flexible few-shot learning benchmarks, demonstrating unsupervised learning obtains more generalizable representations.

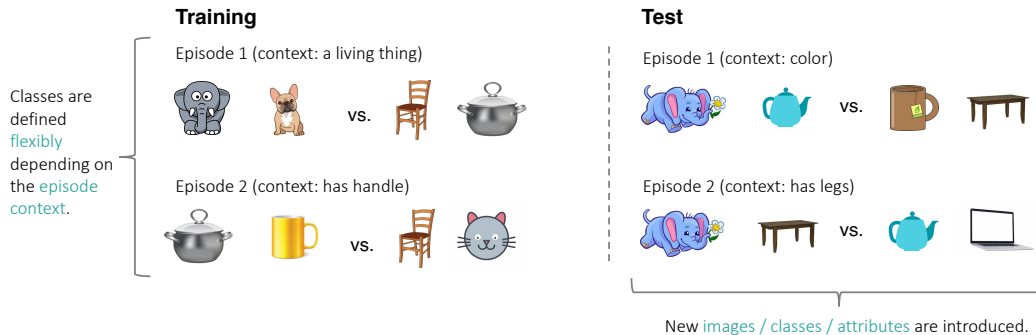


Figure 1: **Illustration of the flexible few-shot learning tasks.** Instead of having a fixed semantic class, each example may belong to different classes flexibly depending on the context of each episode. A context defines a positive class, based on a set of attribute values; all examples that do not match those values belong to the negative class. At test time contexts are defined based on different attributes than during training.

Paradigm	Test time task	Task specification
ZSL (Lampert et al., 2014)	Novel semantic classes	Labeled attributes
FSL (Lake et al., 2011)	Novel semantic classes	Support examples
FFSL (Ours)	Novel classes defined by compositions of unlabeled attributes	Support examples

Table 1: **Differences between zero-shot learning (ZSL), few-shot learning (FSL), and our newly proposed flexible few-shot learning (FFSL).** Our task requires the model to generalize to both new classes and new attributes.

2. Related Work

Few-shot learning: Few-shot learning (FSL) (Fei-Fei et al., 2006; Lake et al., 2011) entails learning new tasks with only a few examples. With an abundance of training data, FSL is closely related to the general meta-learning or learning to learn paradigm, as a few-shot learning algorithm can be developed on training tasks and run on novel tasks at test time. The development of standard few-shot image classification benchmarks such as Omniglot (Lake et al., 2011) and mini-ImageNet (Vinyals et al., 2016), in which each task involves assigning an image to a single semantic class, led to several now-standard methods for few-shot image classification, including MAML (Finn et al., 2017), Matching Network (Vinyals et al., 2016), and Prototypical Network (Snell et al., 2017). Although MAML can adapt based on the context of an FSL episode, it is not empirically better than simpler methods such as Prototypical Networks. To strike a balance between flexibility and simplicity, TADAM (Oreshkin et al., 2018) proposed adapting the network using the FiLM layer (Perez et al., 2018), a generalization of conditional normalization. Various other task conditioning techniques have also been explored in metric-based (Yoon et al., 2019), gradient-based (Ren et al., 2019; Rusu et al., 2019; Simon et al., 2020; Zintgraf et al., 2019), and

parameter prediction-based (Gidaris and Komodakis, 2018; Wang et al., 2019c; Zhao et al., 2018) meta-learners.

The standard few-shot classification task has been extended in various ways, including semi-supervised (Ren et al., 2018), domain shift (Guo et al., 2020; Triantafillou et al., 2019), continual (Antoniou et al., 2020), and online contextualized (Ren et al., 2020) settings. Probabilistic MAML (Finn et al., 2018) investigated the possibility of having ambiguous tasks; in their work, natural image classification was used as an illustrative experiment for comparison to standard MAML, rather than a well developed benchmark. In the same spirit, we extend the study of few-shot learning by introducing our FFSL benchmarks, and show that this task requires consideration of other algorithms.

Zero-shot learning: In zero-shot learning (ZSL) (Akata et al., 2013, 2015; Farhadi et al., 2009; Lampert et al., 2014; Romera-Paredes and Torr, 2015; Xian et al., 2019), a model is asked to recognize classes not present in the training set, supervised only by some auxiliary description (Ba et al., 2015) or attribute values (Farhadi et al., 2009) (see Wang et al. (2019a) for a survey). Lampert et al. (2014) studied the *direct attribute prediction* method, similar to the Supervised Attributes baselines described below in Section 5. The motivation behind our FFSL task can be seen as complementary to ZSL: sometimes a new concept cannot easily be described, but coming up with a small set of representative examples is easier, e.g. “shoes that I like”. An important distinction between ZSL and FFSL is that ZSL uses the same set of attributes for both training and testing (Farhadi et al., 2009; Lampert et al., 2014; Romera-Paredes and Torr, 2015); by contrast, our FFSL task asks the model to learn novel classes defined by attributes for which there are no labels during training. We summarize the relationships between ZSL, FSL and our FFSL in Table 1.

Context dependent similarity: The idea of context dependent similarity of features explored in the present work takes one step towards a more human-like decision maker. In the *contrast model* presented by Tversky (1977), object similarities are expressed by linear combinations of shared and disjoint features of the objects, where the weights depend on the context. As an example, Cuba is similar to Jamaica when speaking about geographic proximity, but similar to Russia when speaking of political viewpoints. Conditional similarity networks (Veit et al., 2017) proposed learning a feature mask for different pre-defined contexts such as colors or styles, using the triplet objective (Wang et al., 2014). Similarly, Wang et al. (2016) proposed to learn a different linear matrix for each context. Kim et al. (2018) explored both context encoding and annotator prior. Context dependent similarity is also an important theme in work on fashion compatibility prediction (Cucurull et al., 2019; Tan et al., 2019; Vasileva et al., 2018). In this paper, we study learning *novel* contextual similarities from only a few examples, and we propose to use a linear classifier at test time with a sparse L1 regularizer to encourage feature selection.

Cold start in recommender systems: Our FFSL tasks share properties of the cold start problem in recommender systems (Gope and Jain, 2017; Lam et al., 2008), in which a new user or item is added to the system with little or no information. As data on the user is being collected, the system must quickly learn to generate good recommendations. The similarity of meta-learning and cold-start recommendation has been explored before (Vartak et al., 2017). Arguably our flexible few-shot tasks are more analogous to cold-start recommendation than the standard FSL setting, as each new

user can be considered as having their own context, consisting of positive examples of items they like and negative examples they dislike.

Compositional learning: Compositional ZSL aims at learning classes (Misra et al., 2017; Purushwalkam et al., 2019; Wang et al., 2019b,c; Yang et al., 2020) defined by a novel composition of labeled attributes. An example of a compositional learning benchmark is the Visual IQ test (Barrett et al., 2018; Zhang et al., 2019), in which novel concepts must be learned from examples, displayed in the form of Raven’s Progressive Matrices. Our FFSL tasks are similar, in that the test-time classes involve novel combinations of attributes. A key distinction is that whereas the novel test classes in compositional ZSL composes known attributes, in our FFSL tasks, the attributes at test time were not labeled during training, and they may not be relevant to any of the classes available at training time.

3. Background

3.1. Standard FSL

The vast majority of standard few-shot classification datasets are constructed as follows. First, a standard supervised classification dataset is obtained (e.g. MNIST). Some number of the classes are designated as training classes (e.g. digits 0-4), and the dataset is partitioned so that all images belonging to the training classes are placed into the training set. The remaining classes are used for validation/testing.

At training time, the learner is given episodes (\mathcal{E}) to learn from. The episode is divided into a labeled *support set* (\mathcal{E}_S) and an unlabeled *query set* (\mathcal{E}_Q). An episode is said to be N -way when it contains data points from only N classes. Additionally, the episode is k -shot when there are k labeled data points from each of the N classes in the support set. Given an episode, the learner must successfully predict the class identity of data points in the query set, given the small amount of labeled information in the support set. Throughout, we use \mathbf{x} to denote input data and y the corresponding class label for this input.

3.2. Unsupervised representation learning

Learning good representation for downstream applications has always been a fundamental aim of deep learning. Hinton and Salakhutdinov (2006) proposed to pretrain subsequent layers of autoencoders for representation learning, and showed good performance for dimensionality reduction, and downstream classification. Following the development of variational autoencoders (VAEs) (Kingma and Welling, 2013), many extensions have been proposed to encourage learning of “disentangled” representations (Higgins et al., 2017; Kim and Mnih, 2018; Liu et al., 2018).

In contrast to traditional generative modeling where the objective is focused on uncovering the data distribution, self-supervised learning has recently emerged as a promising approach for representation learning. These include learning to predict rotations (Kolesnikov et al., 2019), maximizing mutual information between the input and representation (Belghazi et al., 2018; van den Oord et al., 2018), and contrastive learning approaches (Chen et al., 2020; He et al., 2019; Tian et al., 2019;

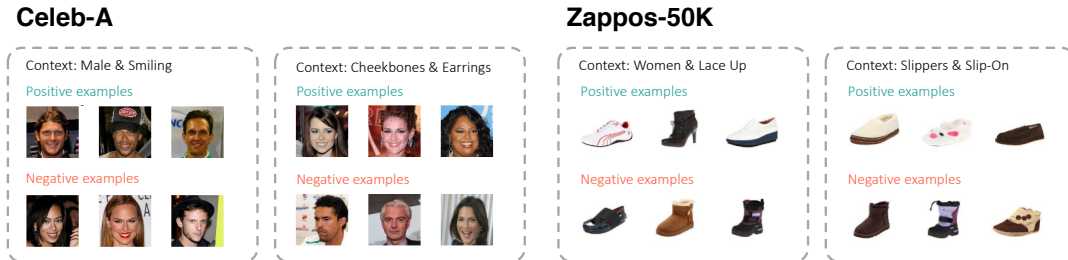


Figure 2: **Sample FFSL episodes using Celeb-A (left) and Zappos-50K (right).** Positive and negative examples are sampled according to the context attributes, but the context information is not revealed to the model at test time.

van den Oord et al., 2018; Xiong et al., 2020). They have shown promise in learning semantic aware representations, almost closing the gap with supervised representation training on the challenging ImageNet benchmark.

SIMCLR (Chen et al., 2020), one of the state-of-the-art methods for self-supervised learning, sends a pair of augmented versions of the same image to the input, and treats them as a positive pair. The pair of hidden representations is further passed into a decoder, producing unit-norm vectors. The network is trained end-to-end to minimize the InfoNCE loss (van den Oord et al., 2018), which distinguishes the positive pair from the rest by encouraging the inner product between the positive pair to gain a higher value than negative pairs.

4. FFSL: Flexible Few-Shot Learning

In this section, we define our FFSL paradigm, and then introduce our two new benchmark datasets for this new FFSL setting.

4.1. Flexible FSL

As in the standard few-shot classification setting, our learner is presented with episodes of data. However, the episodes are not constrained to contain data points from only N classes. Instead, each data point is given either a positive or negative label depending on some criteria that is not known to the learner.

Figure 1 shows some examples of different episodes in our FFSL setting. Each episode contains an image of a pot, but the class identity of the pot varies according to the hidden context. In Episode 1, the pot and the chair are given the same labels whereas in Episode 2 they belong to different classes. Moreover, at test time brand new concepts (e.g. tables) or criteria (e.g. color) may be introduced.

Conceptually, each data point $\mathbf{x} \in \mathcal{X}$ represents some combination of hidden attributes $\mathbf{z} \in \mathcal{Z}$. And each context is an injective function, $f : \mathcal{Z} \rightarrow \{0, 1\}$, that labels each of the data points

depending on their hidden attributes. In this work, we consider contexts that compute conjunctions of binary attributes. At test time, we will see new contexts that we did not see during training.

In order to solve the FFSL task, the learner must correctly find a mapping from the data domain \mathcal{X} to the correct labels. Just like in zero-shot learning, one natural way to solve this problem would be to first find a mapping $h : \mathcal{X} \rightarrow \mathcal{Z}$, that is persistent across episodes, and then estimate the context in each episode. However, we do not limit our exploration to methods that use this approach, since FFSL allows different partitions of the \mathcal{Z} space for training and testing.

Next we describe how we generate the FFSL datasets using existing image datasets with attributes, Celeb-A faces (Liu et al., 2015) and UT Zappos-50K shoes (Yu and Grauman, 2014). Sample episodes from each dataset are shown in Figure 2.

Celeb-A: The Celeb-A dataset (Liu et al., 2015) contains around 200K images of celebrities’ faces. We split half to training, and a quarter to validation and testing each. Each image is annotated with 40 binary attributes, detailing hair color, facial expressions, and other descriptors. We picked 27 salient attributes and split 14 for training and 13 for both val and test. There is no overlap between training or test attributes but they may sometimes belong to a common category, e.g. blond hair is in training and brown hair is in test. Split details are included in the supplementary materials.

Zappos-50K: The UT Zappos-50K dataset (Yu and Grauman, 2014) contains just under 50K images of shoes annotated with attribute values, out of which we kept a total of 76 that we considered salient. We construct an image-level split that assigns 80% of the images to the training set, 10% to the validation and 10% to the test set. We additionally split the set of attribute values into two disjoint sets that are used to form the training and held-out FFSL tasks, respectively.

FFSL episode construction: For each episode, we randomly select one or two attributes (two for Zappos-50K) and look for positive example belonging to these attributes simultaneously. And we also sample an equal number of negative examples that don’t belong to one or neither of the selected attributes. This will construct a *support set* of positive and negative samples, and then we repeat the same process for the corresponding *query set* as well.

5. Learning Novel Contextual Similarity

In this section, we present different baselines to solve FFSL tasks, and our proposed Unsupervised with Fine-Tuning (UFT) model. Overall, we separate learning into two stages: *representation learning* and *few-shot learning*. The aim of the representation learning stage is to utilize the training set towards obtaining a backbone that can successfully solve flexible few-shot tasks. Then, in the FSL stage, a test episode with a few examples is presented, and the learner utilizes the trained backbone network and performs additional learning on top.

Traditionally, the problem of few-shot classification was addressed via episodic models, as described in Section 3.1, so we adopt a number of representative models from this category as baselines: Prototypical Networks (Snell et al., 2017), Matching Networks (Vinyals et al., 2016), a MAML (Finn et al., 2017) variant that only optimizes the top-most layer in the inner loop (Raghu et al., 2020), and two competitive FSL methods that were specifically designed to be context-aware:

TADAM (Oreshkin et al., 2018) and TAFENet (Wang et al., 2019c). TADAM (Oreshkin et al., 2018) uses a task encoding vector to predict the β, γ terms in the batch normalization layers in the convolutional blocks, whereas TAFENet (Wang et al., 2019c) predicts weights in the fully connected feature layers. We used the prototype vector of the positive class as the task encoding for these two models. We train all episodic models on flexible few-shot learning tasks that are derived from the training set of attributes, and we refer to this approach as Flexible Few-Shot Episodic (FFSE) learning.

Following the success of performing supervised pretraining (Chen et al., 2019; Tian et al., 2020), instead of using episodic training as mentioned above, we propose to adopt a two-stage approach to address this problem. The first stage is representation learning, which can be obtained from either supervised or unsupervised pretraining. Then, the second stage is few-shot learning, where we solve a few-shot episode by reading out the representation into binary classes.

5.1. Stage 1: Representation learning

Supervised Attribute prediction (SA): Since attributes are the building blocks from which the “classes” are ultimately defined in our flexible episodes, it is natural to consider a training objective that explicitly trains the base backbone to recognize the training attributes, for the purpose of representation learning. For this, we add a classification layer on top of the feature extractor and use it for the task of attribute prediction. We use a sigmoid activation, effectively solving independent binary tasks to predict the presence or absence of each attribute. This approach is reminiscent of the non-episodic training proposed in Chen et al. (2019) for standard few-shot classification, except that we train to predict attributes instead of object classes as is usually done. We refer to this approach as Supervised Attribute prediction (SA). As an oracle, we also consider a variant that performs attribute prediction for the full set of both training and held-out attributes, referred to as (SA*). The performance of this approach should be thought of as an oracle for this problem, since the test attribute labels are also used at training time.

Unsupervised learning (U): Our proposed flexible few-shot learning tasks differ from the standard setting in that they require a flexible learner that is capable of adapting to different contexts. We hypothesized that learning general-purpose features that capture varying aspects of objects would be helpful to enable this desired flexibility across episodes. We therefore also considered a self-supervised approach, for its ability to learn general semantic features. We chose SIMCLR as a representative from this category due to its empirical success. We refer to this unsupervised approach as (U).

Unsupervised Learning with Fine-tuning (UFT) Finally, we explored whether we can combine the merits of unsupervised representation learning and supervised attribute classification, by fine-tuning SIMCLR’s unsupervised representation for the task of attribute prediction discussed above (SA). To prevent SA from completely overriding the unsupervised features, we add another classifier decoder MLP before the sigmoid classification layer (see Figure 3-B). Empirically, finetuning on SA is found to be beneficial, but early stopping is needed to prevent optimizing too much towards training attributes.

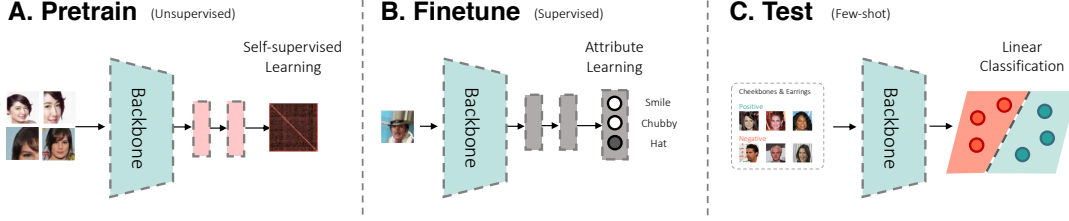


Figure 3: **Our proposed UFT method for FFSL.** **A:** We first pretrain the network with unsupervised contrastive objective to learn general features. **B:** Then we finetune the network to classify the set of training attributes. Both stages employ a different decoder header so that the representation remains general. **C:** Finally at test time we use a linear classifier with a sparsity regularizer.

5.2. Stage 2: Few-shot learning

Once a representation is learned, it remains to be decided how to use the small support set of each given test episode in order to make predictions for the associated query set. MatchingNet (Vinyals et al., 2016) uses a nearest neighbor classifier, whereas ProtoNet (Snell et al., 2017) uses the nearest centroid. Following Chen et al. (2019), we propose to directly learn a linear classifier on top of the representation. This approach learns a weight coefficient for each feature dimension, thus performing some level of feature selection, unlike the nearest-centroid and nearest-neighbor variants. Still, the weights need to be properly regularized to encourage high-fidelity selection. For this, we apply an L1 regularizer on the weights to encourage sparsity. The overall objective of the classifier is:

$$\arg \min_{\mathbf{w}, b} -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y}) + \lambda \|\mathbf{w}\|_1, \quad (1)$$

where $\hat{y} = \sigma(\mathbf{w}^\top \mathbf{h} + b)$, and \mathbf{h} is the representation vector extracted from the CNN backbone. The learning of a classifier is essentially done at the same time as the selection of feature dimensions. We refer to this approach as Logistic Regression (LR), and compare it below to the MatchingNet (nearest-neighbor) and ProtoNet (nearest-centroid) approaches.

6. Experiments

In this section we present our experimental evaluations on our FFSL benchmarks using the different representation learning and few-shot learning methods described in the previous section. For the Celeb-A dataset, we included an additional representation learning method: **ID**, where the backbone was trained for the objective of solving the auxiliary task of face identity classification. In addition to the **SA*** oracle, we provided another oracle **GT-LR**, where the representations used in the test episodes are the ground-truth binary attribute values of the given examples, and the readout is performed by training a linear classifier.

Experimental details: Images were cropped and resized to $84 \times 84 \times 3$. We used ResNet-12 (He et al., 2016; Oreshkin et al., 2018) with 64, 128, 256, 512 channels in each ResBlock. The decoder

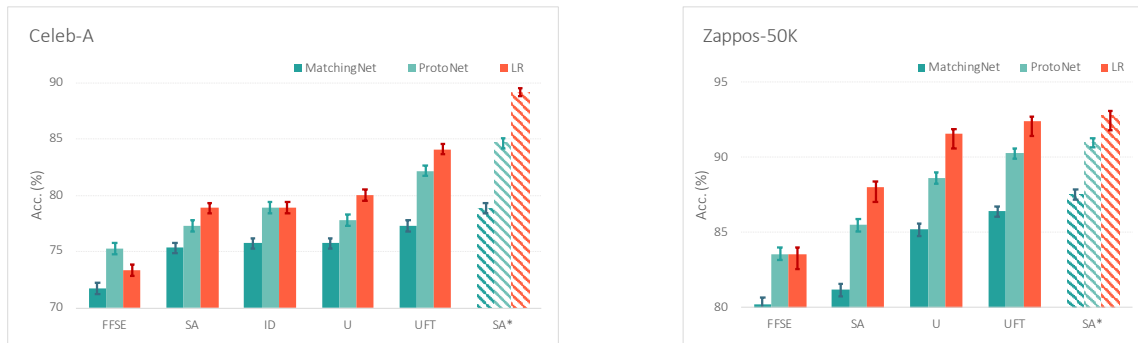


Figure 4: **20-shot FFSL results comparing different representation learning and FSL stage combinations.** **FFSE**: Episodic learning using the flexible few-shot episodes. **SA**: Supervised attribute classification. **ID**: Auxiliary representation learning (for Celeb-A this is face ID classification). **U**: Unsupervised contrastive learning. **UFT**: Our proposed U pretraining followed by SA finetuning. **SA***: Supervised attribute binary classification on **train+test** attributes, which serves as an oracle (striped bars). A set of few-shot learners are evaluated: 1) MatchingNet; 2) ProtoNet; 3) Logistic Regression (LR). For FFSE, LR is trained like MAML (Finn et al., 2017; Raghu et al., 2020) with an inner loop to solve the classification problem. UFT with LR achieves the best performance on both benchmarks. Chance is 50%.

Method	Type	Celeb-A		Zappos-50K	
		5-shot	20-shot	5-shot	20-shot
Chance	-	50.00±0.00	50.00±0.00	50.00±0.00	50.00±0.00
MatchingNet (Vinyals et al., 2016)	FFSE	68.30±0.76	71.73±0.52	77.26±0.60	80.19±0.49
MAML/ANIL (Raghu et al., 2020)	FFSE	71.24±0.74	73.35±0.53	77.82±0.50	83.58±0.41
TAFENet (Wang et al., 2019c)	FFSE	69.10±0.76	72.11±0.54	79.45±0.55	83.29±0.46
ProtoNet (Snell et al., 2017)	FFSE	72.12±0.75	75.27±0.51	78.83±0.51	83.57±0.42
TADAM (Oreshkin et al., 2018)	FFSE	73.54±0.70	76.06±0.53	77.25±0.49	78.63±0.47
SA-LR	Pretrain	72.91±0.74	78.86±0.48	82.22±0.49	88.02±0.38
UFT-LR	Pretrain	78.18±0.68	84.09±0.48	86.03±0.42	92.44±0.30
Test Attribute Oracles					
SA*-LR	Pretrain	84.74±0.60	89.15±0.38	88.21±0.39	92.80±0.27
GT-LR	-	91.07±0.49	98.16±0.17	97.70±0.18	99.59±0.05

Table 2: Flexible 5- and 20-shot learning results on Celeb-A and Zappos-50K, using the same ResNet-12 backbone network.

network for contrastive learning has two 512-d layers and outputs 128-d vectors. The classifier decoder network has two 512-d layers and outputs a 512-d vector. We trained SIMCLR using random crop areas of 0.08 – 1.0, color augmentation 0.5, and InfoNCE temperature 0.5, for 1000 epochs using LARS (You et al., 2017) and cosine schedule with batch size 512 and peak learning rate 2.0. SA finetuning lasts for another 2k steps with batch size 128 and learning rate 0.1 for the

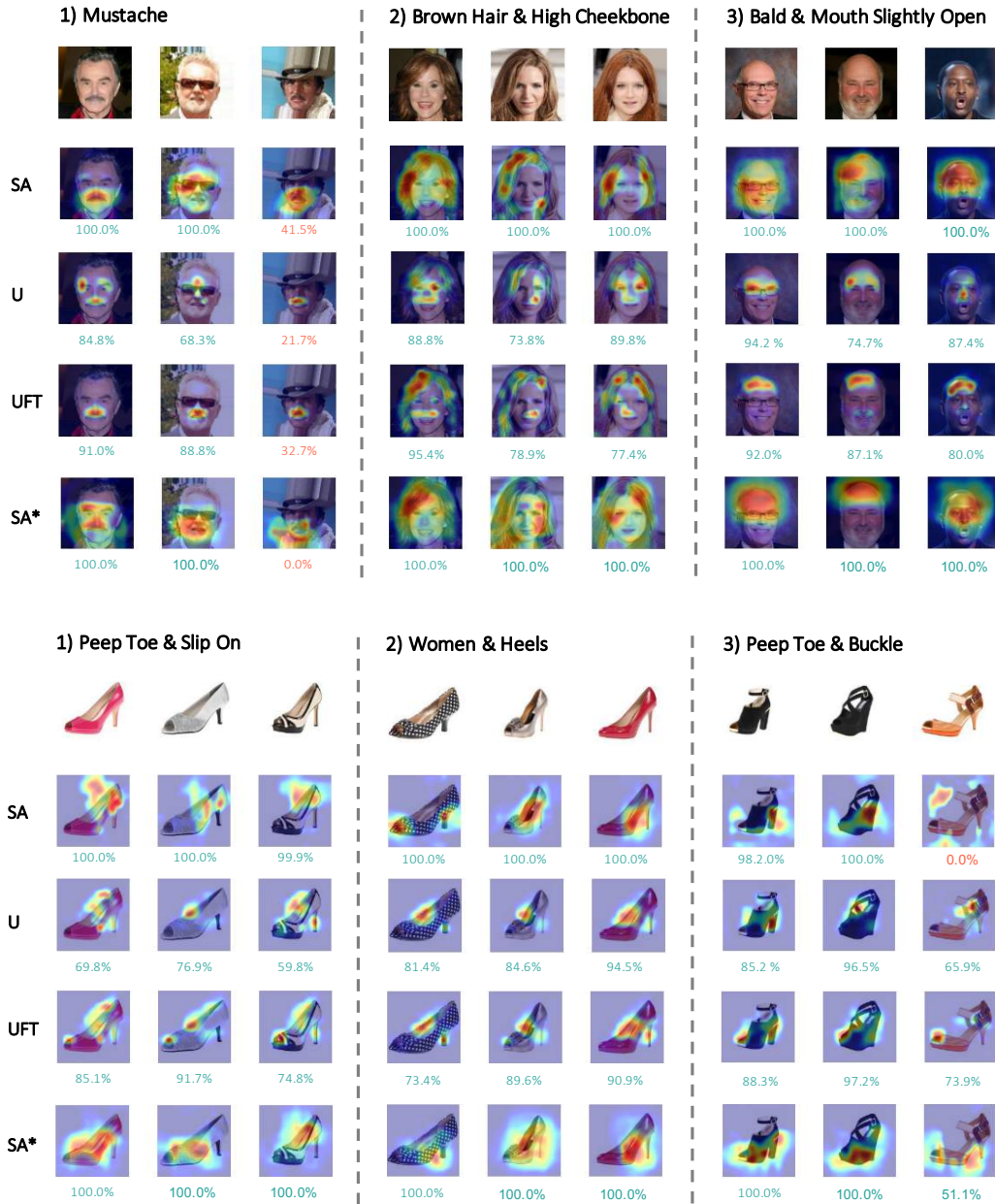


Figure 5: **Visualization of few-shot classifiers using CAM (Zhou et al., 2016), on top of different representations.** Left: Celeb-A; Right: Zappos-50K. Context attributes that define the episode are shown above and images are from the query set of the positive class at test time. Classifier sigmoid confidence scores are shown at the bottom. Red numbers denote wrong classification (below 50% confidence). UFT shows accurate and localized classification explanation.

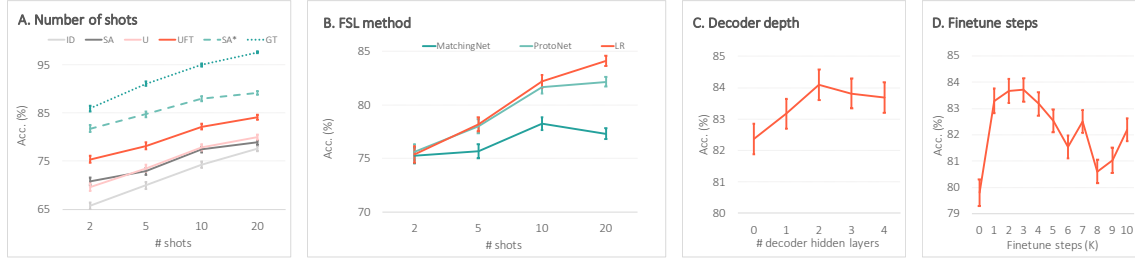


Figure 6: **Additional results on the Celeb-A dataset. A: How many examples are needed for FFSL?** Performance increases with number of shots, even when given the binary ground-truth attribute vector (GT), suggesting that there is greater natural ambiguity in the task than in standard FSL. **B: Comparison of few-shot learning methods on different number of shots.** LR improves with more shots. **C: Effect of the number of decoder layers during finetuning.** Adding a decoder keeps the representation general and reduces overfitting to the training attributes. **D: Effect of the number of finetuning steps.** A small amount of finetuning on the training attribute is beneficial, but eventually hurts.

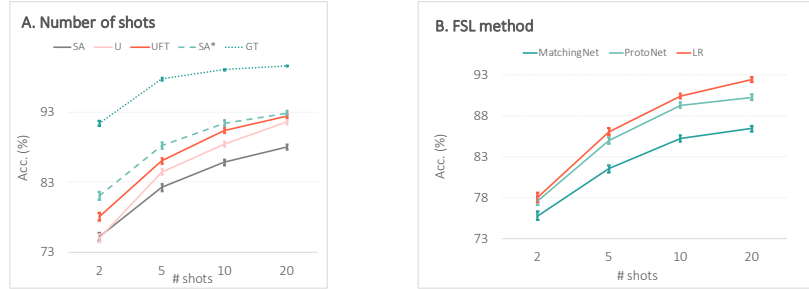


Figure 7: **Additional results on the Zappos-50K dataset. A: How many examples are needed for FFSL?** Performance increases with number of shots, even when given the binary ground-truth attribute vector (GT), suggesting that there is greater natural ambiguity in the task than in standard FSL. **B: Comparison of few-shot learning methods on different number of shots.** LR improves with more shots.

decoder and 0.01 for the backbone and momentum 0.9. ID, SA and SA* use batch size 256 with a learning rate 0.1 for 30k steps, with 0.1x learning rate decay at 20k and 25k steps, and momentum 0.9. Features were normalized before LR classifiers (Tian et al., 2020), and we used cosine similarity for ProtoNet with a learned temperature term (Oreshkin et al., 2018). For the MAML/ANIL variant, we used 10 inner loop steps to adapt the classification layer to each task.

At test time, we report few-shot classification results of running 600 episodes. We report both average accuracy and standard error. LR classifiers are solved with batch gradient descent with 1000 epochs, with a learning rate of 1.0 on Celeb-A and 0.1 on Zappos-50K.

6.1. Results and discussion

Main results: Figure 4 shows our main results on Celeb-A and Zappos-50K with 20-shot FFSL episodes. Results on 5- and 20-shot and other competitive few-shot learning methods are tabulated in Table 2. On both benchmarks, training on flexible few-shot episodes based on training attributes (FFSE) performed worst. Similarly, supervised attribute (SA) learning and learning via the auxiliary task of class facial identification (ID) were not helpful for representation learning either. Interestingly, U attained relatively better test performance, suggesting that the training objective in contrastive learning indeed preserves more general features—not just shown for semantic classification tasks in prior literature, but also for the flexible class definitions present here. Our proposed UFT approach contributed further gains in performance, suggesting that a combination of unsupervised features with some supervised attribute information is indeed beneficial for this task. We also tried to finetune SIMCLR’s representation using FFSE but this did not perform well. We conclude that episodic learning may not help learn higher-level features about the FFSL task itself. Lastly, we confirmed that UFT is able to reduce the generalization gap between SA and SA*, in fact almost closing it entirely in the case of Zappos-50K. These results were consistent across our benchmarks.

Discussion & Analysis We hypothesize that the weak performance of FFSE and SA on our benchmarks is due to the fact that their training objectives essentially encourage ignoring features that aren’t useful at training time, but may still be useful at test time, due to the shift in similarity contexts between the training and testing phases. In Appendix E, we study a toy FFSL problem which further illustrates these generalization issues. We explore training a ProtoNet model on data from a linear generative model, where each episode presents significant ambiguity in resolving the correct context. We show that in this setting, the prototypical network is forced to discard information on the test attributes in order to solve the training tasks effectively, and thus fails to generalize. However, in the equivalent FSL problem no such information destruction occurs as the ProtoNet need not resolve ambiguous context. Intuitively, this is because the distance comparison is made only along the selected classes’ features in each episode. Thus, the same model applied to FSL achieves near-perfect train and test accuracy.

Across our benchmarks, we found that UFT was the most effective representation learning algorithm we explored for FFSL. Interestingly, this result contrasts with standard FSL literature, where unsupervised representation learning still lags behind supervised pretraining (Medina et al., 2020). On the other hand, our flexible few-shot learning results confirms a significant and complementary gain brought by unsupervised representation learning.

Visualizing few-shot classifiers: To understand and interpret the decision made by few-shot linear classifiers, we visualize the classifier weights by using class activation mapping (CAM) (Zhou et al., 2016), and plot the heatmap over the 11×11 spatial feature map in Figure 5. SA sometimes shows incorrect localization as it is not trained to classify those novel test attributes. As expected, SA* oracle performs well, and the area of responsibility is bigger since the training objective encourages the propagation of attribute information spatially across the full feature map. Models pretrained with unsupervised objectives, especially UFT, show surprisingly accurate and localized heatmaps that pinpoint the location of the attributes (e.g. mustache or baldness), while not being exposed to any labeled information of these attributes. This is understandable, as during the contrastive learning stage, local features for mustache or baldness can be good descriptors that match

	SA	U	UFT	SA*	GT
LR	77.4	79.2	83.1	87.1	95.8
+L1 (1e-4)	77.6 (+0.2)	79.4 (+1.2)	83.2 (+0.1)	87.4 (+0.3)	96.1 (+0.3)
+L1 (1e-3)	78.2 (+0.8)	80.2 (+1.0)	83.8 (+0.7)	88.4 (+1.3)	97.1 (+1.3)
+L1 (1e-2)	75.7 (-1.7)	78.3 (-0.9)	79.5 (-3.6)	87.6 (+0.5)	98.2 (+2.4)

Table 3: Effect of the L1 regularizer on different representations for the validation set of Celeb-A.

different views of the same instance together. The fact that the SA models are able to sometimes get correct answers without showing evidence of paying attention to local regions (e.g. in the 3rd column), suggest that they are performing general feature matching rather than reasoning about context dependent similarities. Moreover, U and UFT seem to preserve more uncertainty in the classification outputs.

Number of shots: Since we have a flexible definition of classes in each episode, it could be the case that the support examples are ambiguous. For example, by presenting only an elephant and a cat in the support set, it is unclear whether the positive set is determined by "animals" or "4 legs". Figures 6-A and 7-A show several approaches evaluated using LR with varying numbers of support examples per class in Celeb-A and Zappos-50K FFSL episodes, respectively. GT-LR gradually approached 100% accuracy as the number of shots approached 20. This demonstrates that FFSL tasks potentially require more support examples than standard FSL to resolve ambiguity. Again here, UFT consistently outperformed U, SA, and ID baselines across different number of shots. Figures 6-B and 7-B plot the performance of different FSL methods, using a common UFT representation. LR performs better than MatchingNet and ProtoNet with more support examples.

Effect of decoder depth: Figure 6-C studies the effect of a decoder for attribute classification finetuning. Adding an MLP decoder was found to be beneficial for unsupervised representation learning in prior literature (Chen et al., 2020). Here we found that adding a decoder is also important for SA finetuning, contributing to over 2% improvement.

Effect of SA finetuning: Figure 6-D plots the validation accuracy on FFSL tasks during finetuning for a total of 10k steps. It is found that the accuracy grows from 80% and peaks at 2k steps with over 84%, and then drops. This suggests that a little finetuning on supervised attributes is beneficial, but prolonged finetuning eventually makes the representation less generalizable.

Effect of L1 regularization: Table 3 studies the effect of adding the L1 regularizer on LR. In standard FSL (Chen et al., 2019), typically no regularization is required; however, we found that adding L1 regularization on the weights is beneficial for FFSL. This is especially noticeable on SA* and GT, since it allows the few-shot learner to have a sparse selection of feature dimensions based on the support set context.

Attribute readout: We can also study the usefulness of different representations in terms of their ability to classify attributes. For this task, we trained a new fully connected layer with sigmoid activation on top of the representations, to predict all 40 attributes in the Celeb-A dataset. We used the Adam optimizer and trained for 100 epochs with learning rate 1e-3. The results are reported in

Mean AUC	RND	FFSE	ID	SA	U	UFT	SA*
All (40 attributes)	79.18	88.80	91.29	90.27	92.80	93.33	94.46
Train+Test (27 attributes)	82.27	93.38	94.31	94.23	95.78	96.52	97.18
Train (14 attributes)	84.40	96.04	95.34	95.72	96.43	97.23	97.50
Test (13 attributes)	79.96	90.52	93.19	92.63	95.08	95.76	96.84

Table 4: Celeb-A attribute readout binary prediction performance of different representations, measured in mean AUC. RND denotes using a randomly initialized CNN; FFSE uses a ProtoNet with episodic training.

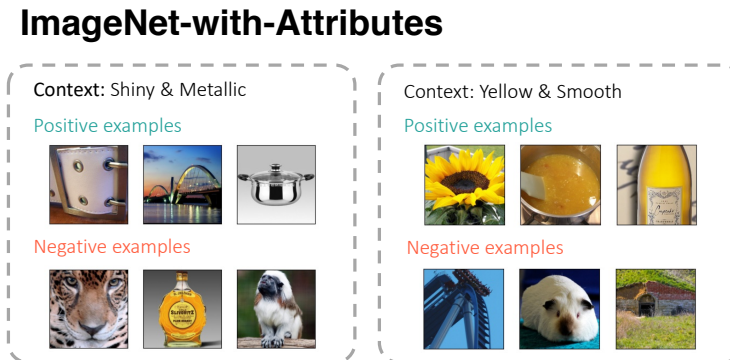


Figure 8: **Sample FFSL episodes of ImageNet-with-Attributes dataset.** Positive and negative examples are sampled according to the context attributes, but the context information is not revealed to the model at test time.

Table 4. UFT is shown to reduce the gap between SA and SA*, demonstrating that unsupervised pretraining allows us to generalize well to new attributes. sectionFlexible Few-Shot Learning on ImageNet Although SIMCLR is shown to be effective, its training requires a significant amount of computing resource, prolonged training time, and a large training set. This motivated us to look into using off-the-shelf SIMCLR models for the FFSL problem. Fortunately, there are publicly released versions of SIMCLR models based on ResNet50 pretrained on ImageNet-1k. In this section, we provide some initial results on using pretrained SIMCLR models on a version of our FFSL task constructed using a subset of the ImageNet dataset.

6.2. ImageNet-with-Attributes

A small subset of the ImageNet dataset comes with attribute annotations. The set has 9.6k images, and each image contains annotations on 25 attributes. To construct a FFSL benchmark using this set, we first split the examples into 5k, 2k, and 2.6k examples for training, validation, and test splits respectively. Four of the 25 attributes (long, square, round, rectangular) are not considered, since they are often too ambiguous to recognize. We split the remaining 21 attributes into 11 training and

10 test attributes. We include attribute split details in Appendix D. The task episodes are constructed in the exact same style as our CelebA benchmark. Figure 8 shows two example episodes.

6.3. Experimental details

Similar to our other results, we compare the pretrained SIMCLR model (U), with UFT, SA, and SA*. The pretrained models U are downloaded from the official source. We used the Hub module¹ from Google’s GitHub repository². All models use a ResNet50 as the encoder. One key difference is that here the U model is pretrained on the full ImageNet-1k training set, a much larger training set than the 5k examples in our ImageNet-with-Attributes training set. UFT is finetuned on the 5k ImageNet-with-Attributes training set. A 2-layer MLP decoder with hidden size 512 was used. Finetuning was done using the Adam optimizer with $1e-4$ learning rate for the decoder and $1e-5$ for the backbone. SA and SA* are trained from scratch on this 5k ImageNet-with-Attributes training set using the same hyperparameters as described for the other two benchmarks. Here we did not include FFSE models because they do not scale to ResNet50 with high-resolution images.

6.4. Results

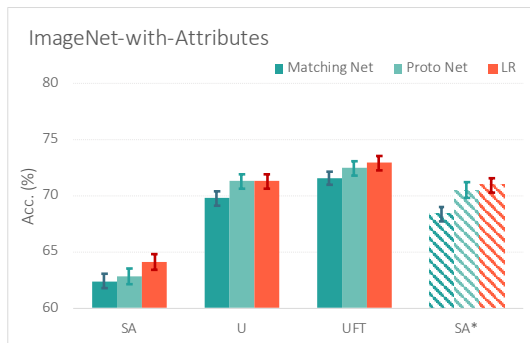


Figure 9: ImageNet-with-Attributes 20-shot FFSL accuracy. UFT improves over U. Both U and UFT outperforms SA* because they were pretrained with unsupervised learning on a much larger training set (ImageNet1k).

In Figure 9, on the 20-shot FFSL task, UFT outperformed U. Unlike our previous results, both U and UFT outperformed SA* here because they were pretrained on a much larger training set. SA and SA* were trained on the ImageNet-with-Attributes 5k training set, and they did not generalize as well to unseen validation and test examples. This result suggests that the off-the-shelf pretrained SIMCLR model provides representations that generalize well, and thus is a promising approach to problems at scale. The fact that UFT improves over U suggest that finetuning the model to focus on attributes will improve the performance on the FFSL task.

1. [gs://simclr-checkpoints/simclrv2/pretrained/r50_1x_sk0/hub](https://github.com/google-research/simclr/blob/master/pretrained/r50_1x_sk0/hub)

2. <https://github.com/google-research/simclr>

Mean AUC	SA	SA*	U	UFT
All (25 attributes)	72.01	73.02	81.08	82.49
Train+Test (21 attributes)	73.43	78.98	80.14	82.37
Train (11 attributes)	72.69	75.86	80.63	82.43
Test (10 attributes)	72.01	74.98	81.08	83.30

Table 5: ImageNet-with-Attributes attribute readout binary prediction performance of different representations, measured in mean AUC.

For further analysis, we looked at how much information on the attributes is stored in the features of each method. On top of a fixed backbone, we trained a MLP classifier to predict all the attributes (both seen and unseen) using the training examples. Table 5 shows the binary attribute prediction performance measured using AUC. Performance on this binary attribute prediction task is highly correlated with the performance on the FFSL. SA has the least amount of information on the unseen attributes, and U has the most. Again, SA* here performs worse on these test examples compared to U/UFT because of the limited number of training examples. In sum, this section shows that using off-the-shelf SIMCLR models pretrained on a large unsupervised dataset can be a promising future direction.

7. Conclusion

The notion of a class often changes depending on the context, yet existing few-shot classification relies on a fixed semantic class definition. In this paper, we propose a flexible few-shot learning paradigm where the similarity criteria change based on the episode context. We proposed benchmarks using the Celeb-A and Zappos-50K datasets to create flexible definitions with existing attribute labels. We explored various ways to perform representation learning for this new task. Unlike in standard FSL, we found that supervised representation learning generalizes poorly on the test set, due to the shift in similarity contexts between the training and testing phases. Unsupervised contrastive learning on the other hand preserved more generalizable features, and further finetuning on supervised attribute classification yielded the best results. Interestingly, unsupervised learning is also thought to play an important role in human and animal learning. Finally, while a sparse regularizer is found to be helpful during readout, in future work we hope it could also influence the representation learning process and hence produce more disentangled representations for FFSL.

Contribution Statement

All authors contributed to the high-level idea and writing of the paper. MR contributed to the code base for running attribute-based few-shot learning experiments, discovered that unsupervised learning plus finetuning is beneficial, performed experiments on Celeb-A, and created most of the figures and graphics. ET helped with figure creation, implemented the flexible few-shot version of Zappos-50K and ran the experiments on that dataset. KCW contributed to the FFSL task definition, implemented the ImageNet-with-Attributes FFSL benchmark, and the associated code for using the

off-the-shelf models. JL designed and implemented early experiments in the FFSL setting, provided the formal description of FFSL, and analyzed the linear toy problem presented in Appendix E. JS contributed to the analysis of early FFSL experiments. XP and AT contributed ideas about the underlying question and its possible solutions, and helped interpret results. RZ contributed many ideas behind the underlying question studied here and the problem formulation, and led the team’s brainstorming about how to test the hypotheses, the datasets and benchmarks, and modeling approaches and visualizations.

Acknowledgments

We would like to thank Claudio Michaelis for several helpful discussions about the problem formulation, Alireza Makhzani for related generative modeling ideas, and Mike Mozer for discussions about contextual similarity. Resources used in preparing this research were provided, in part, by the Province of Ontario, the Government of Canada through CIFAR, and companies sponsoring the Vector Institute (www.vectorinstitute.ai/#partners). This project is supported by NSERC and the Intelligence Advanced Research Projects Activity (IARPA) via Department of Interior/Interior Business Center (DoI/IBC) contract number D16PC00003. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright annotation thereon. Disclaimer: The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of IARPA, DoI/IBC, or the U.S. Government.

References

- Zeynep Akata, Florent Perronnin, Zaïd Harchaoui, and Cordelia Schmid. Label-embedding for attribute-based classification. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2013.
- Zeynep Akata, Scott E. Reed, Daniel Walter, Honglak Lee, and Bernt Schiele. Evaluation of output embeddings for fine-grained image classification. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2015.
- Antreas Antoniou, Massimiliano Patacchiola, Mateusz Ochal, and Amos J. Storkey. Defining benchmarks for continual few-shot learning. *CoRR*, abs/2004.11967, 2020.
- Lei Jimmy Ba, Kevin Swersky, Sanja Fidler, and Ruslan Salakhutdinov. Predicting deep zero-shot convolutional neural networks using textual descriptions. In *IEEE International Conference on Computer Vision, ICCV*, 2015.
- David GT Barrett, Felix Hill, Adam Santoro, Ari S Morcos, and Timothy Lillicrap. Measuring abstract reasoning in neural networks. *arXiv preprint arXiv:1807.04225*, 2018.
- Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeswar, Sherjil Ozair, Yoshua Bengio, Aaron Courville, and R Devon Hjelm. Mine: mutual information neural estimation. *arXiv preprint arXiv:1801.04062*, 2018.

- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey E. Hinton. A simple framework for contrastive learning of visual representations. *CoRR*, abs/2002.05709, 2020.
- Wei-Yu Chen, Yen-Cheng Liu, Zsolt Kira, Yu-Chiang Frank Wang, and Jia-Bin Huang. A closer look at few-shot classification. In *Proceedings of the 7th International Conference on Learning Representations, ICLR*, 2019.
- Guillem Cucurull, Perouz Taslakian, and David Vázquez. Context-aware visual compatibility prediction. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2019.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2009.
- Ali Farhadi, Ian Endres, Derek Hoiem, and David A. Forsyth. Describing objects by their attributes. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition, CVPR*, 2009.
- Li Fei-Fei, Rob Fergus, and Pietro Perona. One-shot learning of object categories. *IEEE TPAMI*, 28(4):594–611, 2006.
- 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, ICML*, 2017.
- Chelsea Finn, Kelvin Xu, and Sergey Levine. Probabilistic model-agnostic meta-learning. In *Advances in Neural Information Processing Systems 31, NeurIPS*, 2018.
- Spyros Gidaris and Nikos Komodakis. Dynamic few-shot visual learning without forgetting. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2018.
- Jyotirmoy Gope and Sanjay Kumar Jain. A survey on solving cold start problem in recommender systems. In *2017 International Conference on Computing, Communication and Automation (ICCCA)*, pages 133–138. IEEE, 2017.
- Yunhui Guo, Noel C. F. Codella, Leonid Karlinsky, John R. Smith, Tajana Rosing, and Rogério Schmidt Feris. A broader study of cross-domain few-shot learning. In *14th European Conference on Computer Vision, ECCV*, 2020.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2016.
- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross B. Girshick. Momentum contrast for unsupervised visual representation learning. *CoRR*, abs/1911.05722, 2019.
- Irina Higgins, Loïc 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 *5th International Conference on Learning Representations, ICLR*, 2017.

- Geoffrey E Hinton and Ruslan R Salakhutdinov. Reducing the dimensionality of data with neural networks. *science*, 313(5786):504–507, 2006.
- Hyunjik Kim and Andriy Mnih. Disentangling by factorising. In *Proceedings of the 35th International Conference on Machine Learning, ICML*, 2018.
- Kun Ho Kim, Oisin Mac Aodha, and Pietro Perona. Context embedding networks. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2018.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Alexander Kolesnikov, Xiaohua Zhai, and Lucas Beyer. Revisiting self-supervised visual representation learning. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 1920–1929, 2019.
- Brenden M. Lake, Ruslan Salakhutdinov, Jason Gross, and Joshua B. Tenenbaum. One shot learning of simple visual concepts. In *Proceedings of the 33th Annual Meeting of the Cognitive Science Society, CogSci*, 2011.
- Xuan Nhat Lam, Thuc Vu, Trong Duc Le, and Anh Duc Duong. Addressing cold-start problem in recommendation systems. In *Proceedings of the 2nd international conference on Ubiquitous information management and communication*, pages 208–211, 2008.
- Christoph H. Lampert, Hannes Nickisch, and Stefan Harmeling. Attribute-based classification for zero-shot visual object categorization. *IEEE Trans. Pattern Anal. Mach. Intell.*, 36(3):453–465, 2014.
- Yu Liu, Fangyin Wei, Jing Shao, Lu Sheng, Junjie Yan, and Xiaogang Wang. Exploring disentangled feature representation beyond face identification. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2018.
- Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision, ICCV*, 2015.
- Carlos Medina, Arnout Devos, and Matthias Grossglauser. Self-supervised prototypical transfer learning for few-shot classification. *CoRR*, abs/2006.11325, 2020.
- Ishan Misra, Abhinav Gupta, and Martial Hebert. From red wine to red tomato: Composition with context. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2017.
- Boris N. Oreshkin, Pau Rodríguez López, and Alexandre Lacoste. TADAM: task dependent adaptive metric for improved few-shot learning. In *Advances in Neural Information Processing Systems 31, NeurIPS*, 2018.
- Ethan Perez, Florian Strub, Harm de Vries, Vincent Dumoulin, and Aaron C. Courville. Film: Visual reasoning with a general conditioning layer. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, AAAI*, 2018.

- Senthil Purushwalkam, Maximilian Nickel, Abhinav Gupta, and Marc’Aurelio Ranzato. Task-driven modular networks for zero-shot compositional learning. In *IEEE/CVF International Conference on Computer Vision, ICCV*, 2019.
- Aniruddh Raghu, Maithra Raghu, Samy Bengio, and Oriol Vinyals. Rapid learning or feature reuse? towards understanding the effectiveness of MAML. In *8th International Conference on Learning Representations, ICLR*, 2020.
- Mengye Ren, Eleni Triantafillou, Sachin Ravi, Jake Snell, Kevin Swersky, Joshua B. Tenenbaum, Hugo Larochelle, and Richard S. Zemel. Meta-learning for semi-supervised few-shot classification. In *Proceedings of 6th International Conference on Learning Representations, ICLR*, 2018.
- Mengye Ren, Renjie Liao, Ethan Fetaya, and Richard S. Zemel. Incremental few-shot learning with attention attractor networks. In *Advances in Neural Information Processing Systems 32, NeurIPS*, 2019.
- Mengye Ren, Michael L. Iuzzolino, Michael C. Mozer, and Richard S. Zemel. Wandering within a world: Online contextualized few-shot learning. *CoRR*, abs/2007.04546, 2020.
- Bernardino Romera-Paredes and Philip H. S. Torr. An embarrassingly simple approach to zero-shot learning. In *Proceedings of the 32nd International Conference on Machine Learning, ICML*, 2015.
- Andrei A. Rusu, Dushyant Rao, Jakub Sygnowski, Oriol Vinyals, Razvan Pascanu, Simon Osindero, and Raia Hadsell. Meta-learning with latent embedding optimization. In *7th International Conference on Learning Representations, ICLR*, 2019.
- Christian Simon, Piotr Koniusz, Richard Nock, and Mehrtash Harandi. On modulating the gradient for meta-learning. In *16th European Conference on Computer Vision, ECCV*, 2020.
- Jake Snell, Kevin Swersky, and Richard S. Zemel. Prototypical networks for few-shot learning. In *Advances in Neural Information Processing Systems 30, NIPS*, 2017.
- Reuben Tan, Mariya I. Vasileva, Kate Saenko, and Bryan A. Plummer. Learning similarity conditions without explicit supervision. In *IEEE/CVF International Conference on Computer Vision, ICCV*, 2019.
- Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. *CoRR*, abs/1906.05849, 2019.
- Yonglong Tian, Yue Wang, Dilip Krishnan, Joshua B. Tenenbaum, and Phillip Isola. Rethinking few-shot image classification: a good embedding is all you need? *CoRR*, abs/2003.11539, 2020.
- Eleni Triantafillou, Tyler Zhu, Vincent Dumoulin, Pascal Lamblin, Kelvin Xu, Ross Goroshin, Carles Gelada, Kevin Swersky, Pierre-Antoine Manzagol, and Hugo Larochelle. Meta-dataset: A dataset of datasets for learning to learn from few examples. In *8th International Conference on Learning Representations, ICLR*, 2019.
- Amos Tversky. Features of similarity. *Psychological Review*, 84(4):327–352, 1977. ISSN 19391471. doi: 10.1037/0033-295X.84.4.327.

- Aäron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *CoRR*, abs/1807.03748, 2018.
- Manasi Vartak, Arvind Thiagarajan, Conrado Miranda, Jeshua Bratman, and Hugo Larochelle. A meta-learning perspective on cold-start recommendations for items. In *Advances in neural information processing systems, NIPS*, 2017.
- Mariya I. Vasileva, Bryan A. Plummer, Krishna Dusad, Shreya Rajpal, Ranjitha Kumar, and David A. Forsyth. Learning type-aware embeddings for fashion compatibility. In *15th European Conference on Computer Vision, ECCV*, 2018.
- Andreas Veit, Serge J. Belongie, and Theofanis Karaletsos. Conditional similarity networks. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2017.
- Oriol Vinyals, Charles Blundell, Tim Lillicrap, Koray Kavukcuoglu, and Daan Wierstra. Matching networks for one shot learning. In *Advances in Neural Information Processing Systems 29, NIPS*, 2016.
- Jiang Wang, Yang Song, Thomas Leung, Chuck Rosenberg, Jingbin Wang, James Philbin, Bo Chen, and Ying Wu. Learning fine-grained image similarity with deep ranking. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2014.
- Wei Wang, Vincent W Zheng, Han Yu, and Chunyan Miao. A survey of zero-shot learning: Settings, methods, and applications. *ACM Transactions on Intelligent Systems and Technology, TIST*, 10 (2):1–37, 2019a.
- Xiaofang Wang, Kris M. Kitani, and Martial Hebert. Contextual visual similarity. *CoRR*, abs/1612.02534, 2016.
- Xin Wang, Fisher Yu, Trevor Darrell, and Joseph E. Gonzalez. Task-aware feature generation for zero-shot compositional learning. *CoRR*, abs/1906.04854, 2019b.
- Xin Wang, Fisher Yu, Ruth Wang, Trevor Darrell, and Joseph E. Gonzalez. Tafe-net: Task-aware feature embeddings for low shot learning. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2019c.
- Yongqin Xian, Christoph H. Lampert, Bernt Schiele, and Zeynep Akata. Zero-shot learning - A comprehensive evaluation of the good, the bad and the ugly. *IEEE Trans. Pattern Anal. Mach. Intell.*, 41(9):2251–2265, 2019.
- Yuwen Xiong, Mengye Ren, and Raquel Urtasun. Loco: Local contrastive representation learning. *CoRR*, abs/2008.01342, 2020.
- Muli Yang, Cheng Deng, Junchi Yan, Xianglong Liu, and Dacheng Tao. Learning unseen concepts via hierarchical decomposition and composition. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*, 2020.
- Sung Whan Yoon, Jun Seo, and Jaekyun Moon. Tapnet: Neural network augmented with task-adaptive projection for few-shot learning. In *Proceedings of the 36th International Conference on Machine Learning, ICML*, 2019.

- Yang You, Igor Gitman, and Boris Ginsburg. Scaling SGD batch size to 32k for imagenet training. *CoRR*, abs/1708.03888, 2017.
- Aron Yu and Kristen Grauman. Fine-grained visual comparisons with local learning. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2014.
- Chi Zhang, Feng Gao, Baoxiong Jia, Yixin Zhu, and Song-Chun Zhu. RAVEN: A dataset for relational and analogical visual reasoning. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2019.
- Fang Zhao, Jian Zhao, Shuicheng Yan, and Jiashi Feng. Dynamic conditional networks for few-shot learning. In *15th European Conference on Computer Vision, ECCV*, 2018.
- Bolei Zhou, Aditya Khosla, Àgata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2016.
- Luisa M. Zintgraf, Kyriacos Shiarlis, Vitaly Kurin, Katja Hofmann, and Shimon Whiteson. Fast context adaptation via meta-learning. In *Proceedings of the 36th International Conference on Machine Learning, ICML*, 2019.

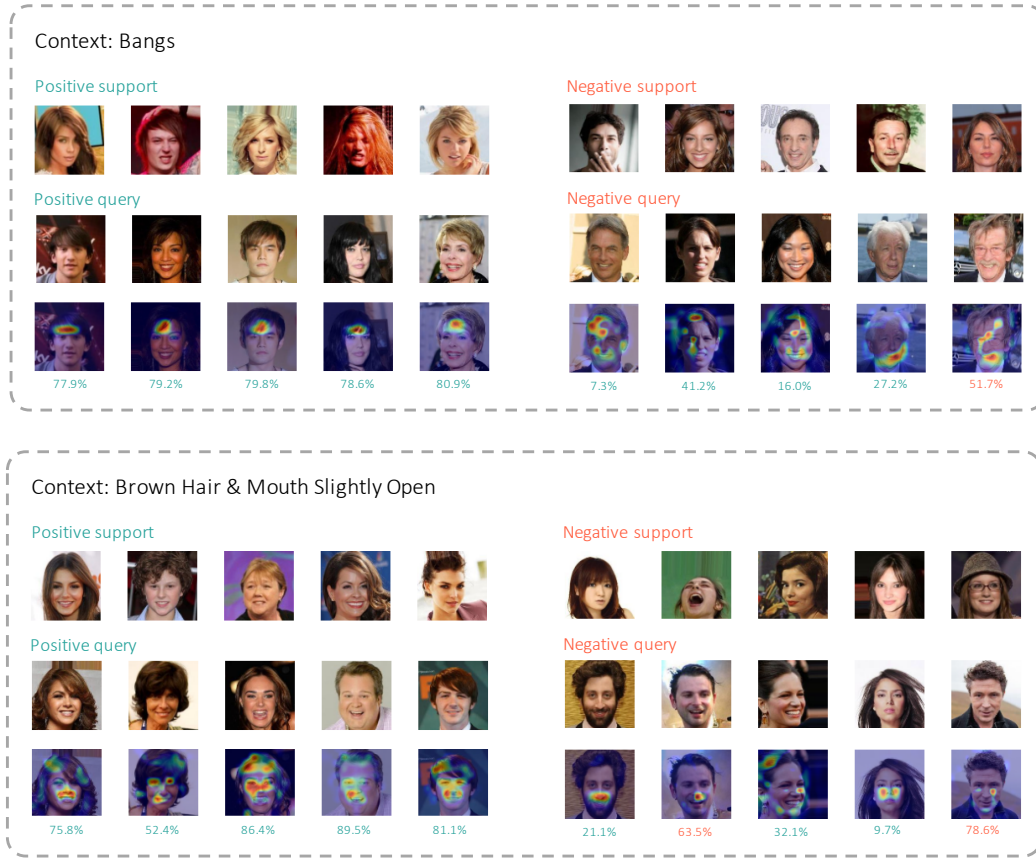


Figure 10: **Visualization of Celeb-A 20-shot LR classifiers using CAM on top of UFT representations.** Context attributes that define the episode are shown above. Classifier sigmoid confidence scores are shown at the bottom. Red numbers denote wrong classification and green denote correct.

Appendix A. Additional heatmap visualization

We provide additional visualization results in Figure 10, 11, and we plot the heat map to visualize the LR classifier weights. Some of the examples are ambiguous in nature and they are challenging to get it correct.

Appendix B. Attribute splits of Celeb-A

We include the attribute split for Celeb-A in Table 6. There are 14 attributes in training and 13 attributes in val/test. We discarded the rest of the 13 attributes in the original datasets since they were not very obvious.



Figure 11: **Visualization of Zappos-50K 20-shot LR classifiers using CAM on top of UFT representations.** Context attributes that define the episode are shown above. Classifier sigmoid confidence scores are shown at the bottom. Red numbers denote wrong classification and green denote correct.

Train	5_o_Clock_Shadow	Black_Hair	Blond_Hair	Chubby
	Double_Chin	Eyeglasses	Goatee	Gray_Hair
	Male	No_Beard	Pale_Skin	Receding_Hairline
	Rosy_Cheeks	Smiling		
Val/Test	Bald	Bangs	Brown_Hair	Heavy_Makeup
	High_Cheekbones	Mouth_Slightly_Open	Mustache	Narrow_Eyes
	Sideburns	Wearing_Earrings	Wearing_Hat	Wearing_Lipstick
	Wearing_Necktie			

Table 6: **Attribute Splits for Celeb-A**

Train	Category-Shoes	Category-Sandals	SubCategory-Oxfords	SubCategory-Heel
	SubCategory-Boot	SubCategory-Slipper Flats	SubCategory-Short heel	SubCategory-Flats
	SubCategory-Slipper Heels	SubCategory-Athletic	SubCategory-Knee High	SubCategory-Crib Shoes
	SubCategory-Over the Knee	HeelHeight-High heel	Closure-Pull-on	Closure-Ankle Strap
	Closure-Zipper	Closure-Elastic Gore	Closure-Sling Back	Closure-Toggle
	Closure-Snap	Closure-T-Strap	Closure-Spat Strap	Gender-Men
	Gender-Boys	Material-Rubber	Material-Wool	Material-Silk
	Material-Aluminum	Material-Plastic	Toestyle-Capped Toe	Toestyle-Square Toe
	Toestyle-Snub Toe	Toestyle-Bicycle Toe	Toestyle-Open Toe	Toestyle-Pointed Toe
	Toestyle-Almond	Toestyle-Apron Toe	Toestyle-Snip Toe	Toestyle-Medallion
	Category-Boots	Category-Slippers	SubCategory-Mid-Calf	SubCategory-Ankle
	SubCategory-Loafers	SubCategory-Boat Shoes	SubCategory-Clogs and Mules	SubCategory-Sneakers and Athletic Shoes
	SubCategory-Heels	SubCategory-Prewalker	SubCategory-Prewalker Boots	SubCategory-Firstwalker
Val/Test	HeelHeight-Short heel	Closure-Lace up	Closure-Buckle	Closure-Hook and Loop
	Closure-Slip-On	Closure-Ankle Wrap	Closure-Bungee	Closure-Adjustable
	Closure-Button Loop	Closure-Monk Strap	Closure-Belt	Gender-Women
	Gender-Girls	Material-Suede	Material-Snakeskin	Material-Corduroy
	Material-Horse Hair	Material-Stingray	Toestyle-Round Toe	Toestyle-Closed Toe
	Toestyle-Moc Toe	Toestyle-Wingtip	Toestyle-Center Seam	Toestyle-Algonquin
	Toestyle-Bump Toe	Toestyle-Wide Toe Box	Toestyle-Peep Toe	

Table 7: Attribute splits for Zappos-50K

Appendix C. Attribute splits of Zappos-50K

The Zappos-50K dataset annotates images with different values relating to the following aspects of shoes: ‘Category’, ‘Subcategory’, ‘HeelHeight’, ‘Insole’, ‘Closure’, ‘Gender’, ‘Material’ and ‘Toestyle’.

We discarded the ‘Insole’ values, since those refer to the inside part of the shoe which isn’t visible in the images. We also discarded some ‘Material’ values that we deemed hard to recognize visually. We also modified the values of ‘HeelHeight’ which originally was different ranges of cm of the height of the heel of each shoe. Instead, we divided those values into only two groups: ‘short heel’ and ‘high heel’, to avoid having to perform very fine-grained heel height recognition which we deemed was too difficult.

These modifications leave us with a total of 79 values (across all higher-level categories). Not all images are tagged with a value from each category, while some are even tagged with more than one value from the same category (e.g. two different materials used in different parts of the shoe). We split these values into 40 ‘training attributes’ and 39 ‘val/test attributes’. As mentioned in the main paper, the training attributes are used to construct training episodes (when performing episodic training), or to define the classification layer (in the case of the SA and UFT models). The ‘val/test attributes’, on the other hand, are used to construct our flexible evaluation episodes. For example, a particular training episode might define its positive class as the conjunction: ‘Category=Sandals and Material=Plastic’ and a test episode might define its positive class as the conjunction: ‘Category=Boots and Closure=Buckle’.

We include the complete list of attributes in Table 7. The format we use is ‘X-Y’ where X stands for the category (e.g. ‘Material’) and Y stands for the value of that category (e.g. ‘Wool’). We do this to avoid ambiguity, since it may happen that different categories have some value names in common, e.g. ‘Short Heel’ is a value of both ‘SubCategory’ and ‘HeelHeight’.

Train	pink	spotted	wet	blue
	shiny	rough	striped	white
	metallic	wooden	gray	
Val/Test	brown	green	violet	red
	orange	yellow	furry	black
	vegetation	smooth		

Table 8: Attribute Splits for ImageNet-with-Attributes

Appendix D. Attribute splits of ImageNet-with-Attributes

We include the attribute split for ImageNet-with-Attributes in Table 8. There are 11 attributes in training and 10 attributes in val/test. We discarded the rest of the 4 attributes in the original datasets.

Appendix E. Flexible Few-Shot Toy Problem

In this section, we present a toy problem that illustrates the challenges introduced by the flexible few-shot learning setting and the failures of existing approaches on this task. This simple model captures the core elements of our flexible few-shot tasks, including ambiguity, domain shift from training to test time, and the role of learning good representations. The primary limitation of this model is the fact that it is fully linear and the attribute values are independent—in a more realistic FFSL task recovering a good representation from the data is significantly more challenging, and the data points will have a more complex relationship with the attributes as in our benchmark datasets.

Problem setup We define a flexible few-shot learning problem where the data points $\mathbf{x} \in \mathbb{R}^m$ are generated from binary attribute strings, $\mathbf{z} \in \{0, 1\}^d$, with $\mathbf{x} = A\mathbf{z} + \boldsymbol{\zeta}$ for some matrix $A \in \mathbb{R}^{m \times d}$ with full column rank and noise source $\boldsymbol{\zeta}$. Thus, each data point \mathbf{x} is a sum of columns of A with some additive noise.

We consider contexts that classify the examples as positive when two attributes are both 1-valued, and negative otherwise. For the training episodes, the contexts depend only on the first $d_1 < d$ attributes. At test time, the episode contexts depend on the remaining $d - d_1$ attributes. The episodes are generated by sampling a context uniformly, that is, choosing two attributes that represent a class. Then k data points are sampled with positive labels (the two attributes defining the context are 1-valued) and k with negative labels (at least one of the attributes is 0-valued).

Linear prototypical network Now, consider training a prototypical network on this data with a linear embedding network, $g(\mathbf{x}) = W\mathbf{x}$. Within each episode, the prototypical network computes the prototypes for the positive and negative classes,

$$\mathbf{c}_j = \frac{1}{k} \sum_{\mathbf{x}_i \in S_j} g(\mathbf{x}_i) = \frac{1}{k} \sum_{\mathbf{x}_i \in S_j} \sum_{l=1}^d z_{il} W \mathbf{a}_l, \text{ for } j \in \{0, 1\},$$

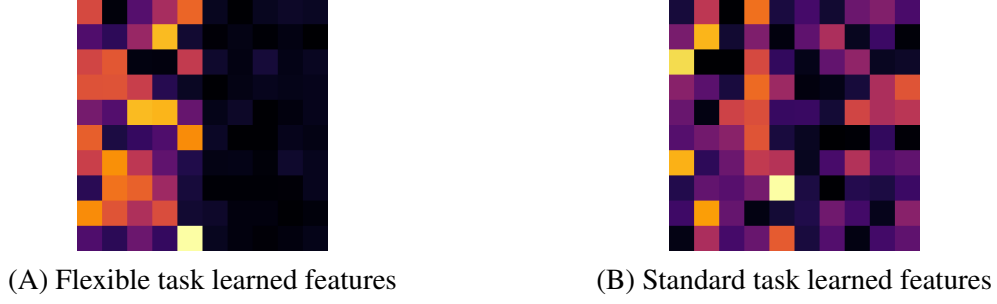


Figure 12: Projecting data features into prototypical network embedding space (WA) for the linear toy problem. Values closer to zero are darker in colour. On the flexible task, the model destroys information from the test attributes to remove ambiguity at training time.

where S_j is the set of data points in the episode with label j , and \mathbf{a}_l is the l^{th} column of the matrix A . Further, the prototypical network likelihood is given by,

$$p(y = 0|\mathbf{x}) = \frac{\exp\{-\|W\mathbf{x} - \mathbf{c}_0\|_2^2\}}{\exp\{-\|W\mathbf{x} - \mathbf{c}_0\|_2^2\} + \exp\{-\|W\mathbf{x} - \mathbf{c}_1\|_2^2\}}.$$

The goal of the prototypical network is thus to learn weights W that lead to small distances between data points in the same class and large distances otherwise. In the flexible few-shot learning tasks, there is an additional challenge in that class boundaries shift between episodes. The context defining the boundary is unknown and must be inferred from the episode. However, with few shots (small k) there is ambiguity in the correct context — with a high probability that several possible contexts provide valid explanations for the observed data.

Fitting the prototypical network Notice that under our generative model, with $\mathbf{x} = W\mathbf{z} + \boldsymbol{\zeta}$ and for $j \in \{0, 1\}$ we have,

$$W\mathbf{x} - \mathbf{c}_j = WA(\mathbf{z} - \frac{1}{k} \sum_{\mathbf{z}_i \in S_j} \mathbf{z}_i) + \frac{1}{k} \sum_i W\boldsymbol{\zeta}_i + W\boldsymbol{\zeta}.$$

Notice that if $\mathbf{v}_j(\mathbf{z}) = A(\mathbf{z} - \frac{1}{k} \sum_{\mathbf{z}_i \in S_j} \mathbf{z}_i) \in \text{Ker}(W)$, the kernel of W , then the entire first term is zero. Further, if $\mathbf{z} \in S_j$ (the same class as the prototype) then there is no contribution from the positive attribute features in this term. Otherwise, this term is guaranteed to have some contribution from the positive attribute features.

Therefore, if W projects to the linear space spanned by the positive attribute features then $W\mathbf{v}_j(\mathbf{z})$ is zero when $\mathbf{z} \in S_j$ and non-zero otherwise. This means that the model will be able to solve the episode without contextual ambiguity. Then the optimal weights are those that project to the set of features used in the training set—destroying all information about the test attributes which would otherwise introduce ambiguity.

We observed this effect empirically in Figure 12, where we have plotted the matrix $\text{abs}(WA)$. Each column of these plots represents a column of A mapped to the prototypical network’s embedding space. The first 5 columns correspond to attributes used at training time, and the remaining 5 to those used at test time.

In the flexible task described above, as our analysis suggests, the learned prototypical feature weights project out the features used at test time (the last 5 columns). As a result, the model achieved 100% training accuracy but only 51% test accuracy (chance is 50%).

We also compared against an equivalent problem set up that resembles the standard few-shot learning setting. In the FSL problem, the binary attribute strings may have only a single non-zero entry and each episode is a binary classification problem where the learner must distinguish between two classes. Now the vector \mathbf{z} is a one-hot encoding and the comparison to the prototypes occurs only over a single feature column of A , thus there is no benefit to projecting out the test features. As expected, the model we learned (Figure 12 B) is not forced to throw away test-time information and achieves 100% training accuracy and 99% test accuracy.

Settings for Figure 12 We use 10 attributes, 5 of which are used for training and 5 for testing. We use a uniformly random sampled $A \in \mathbb{R}^{30 \times 10}$ and the prototypical network learns $W \in \mathbb{R}^{10 \times 30}$. We use additive Gaussian noise when sampling data points with a standard deviation of 0.1. The models are trained with the Adam optimizer using default settings over a total of 30000 random episodes, and evaluated on an additional 1000 test episodes. We used $k = 20$ to produce these plots, but found that the result was consistent over different shot counts.