

A Universal Representation Transformer Layer for Few-Shot Image Classification

Lu Liu^{1,2*}, William Hamilton^{1,3†}, Guodong Long², Jing Jiang², Hugo Larochelle^{1,4†}

¹ Mila, ² Australian AI Institute, UTS, ³ McGill University, ⁴ Google Research, Brain Team
Correspondence to lu.liu.cs@icloud.com

Abstract

Few-shot classification aims to recognize unseen classes when presented with only a small number of samples. We consider the problem of *multi-domain* few-shot image classification, where unseen classes and examples come from diverse data sources. This problem has seen growing interest and has inspired the development of benchmarks such as Meta-Dataset. A key challenge in this multi-domain setting is to effectively integrate the feature representations from the diverse set of training domains. Here, we propose a Universal Representation Transformer (URT) layer, that meta-learns to leverage universal features for few-shot classification by dynamically re-weighting and composing the most appropriate domain-specific representations. In experiments, we show that URT sets a new state-of-the-art result on Meta-Dataset. Specifically, it outperforms the best previous model on three data sources or performs the same in others. We analyze variants of URT and present a visualization of the attention score heatmaps that sheds light on how the model performs cross-domain generalization. Our code is available at <https://github.com/liulu112601/URT>

1 Introduction

Learning tasks from small data remains a challenge for machine learning systems, which show a noticeable gap compared to the ability of humans to understand new concepts from few examples. A promising direction to address this challenge is developing methods that are capable of performing transfer learning across the collective data of many tasks. Since machine learning systems generally improve with the availability of more data, a natural assumption is that few-shot learning systems should benefit from leveraging data *across many different tasks and domains*—even if each individual task has limited training data available.

This research direction is well captured by the problem of *multi-domain few-shot classification*. In this setting, training and test data spans a number of different domains, each represented by a different source dataset. A successful approach in this multi-domain setting must not only address the regular challenge of few-shot classification—i.e., the challenge of having only a handful of examples per class. It must also discover how to leverage (or ignore) what is learned from different domains, achieving generalization and avoiding cross-domain interference.

Recently, Triantafillou et al. [20] proposed a benchmark for multi-domain few-shot classification, Meta-Dataset, and highlighted some of the challenges that current methods face when training data is heterogeneous. Crucially, they found that methods which trained on all available domains would normally obtain improved performance on some domains at the expense of others. Following on their work, progress has been made, which includes the design of adapted hyper-parameter optimization

*This work was done while Lu Liu was a research intern with Mila.

†Canada CIFAR AI Chair

strategies [17] and more flexible meta-learning algorithms [16]. Most notable is SUR (Selecting Universal Representation) [5], a method that relies on a so-called universal representation, extracting from a collection of pre-trained and domain-specific neural network backbones. SUR prescribes a hand-crafted feature-selection procedure to infer how to weight each backbone for each task at hand, and produces an adapted representation for each task. This was shown to lead to state-of-the-art performance on Meta-Dataset.

In SUR, the classification procedure for each task is fixed and not learned. Thus, except for the underlying universal representation, there is no transfer learning performed with regards to how classification rules are inferred across tasks and domains. Yet, cross-domain generalization might be beneficial in that area as well, in particular when tasks have only few examples per class.

Present work. To explore this question, we propose the Universal Representation Transformer (URT) layer, which can effectively learn to transform a universal representation into task-adapted representations. The URT layer is inspired from Transformer networks [21] and uses an attention mechanism to learn to retrieve or blend the appropriate backbones to use for each task. By training this layer across few-shot tasks from many domains, it can support transfer across these tasks.

We show that our URT layer on top of a universal representation’s pre-trained backbones sets a new state-of-the-art performance on Meta-Dataset. It succeeds at outperforming SUR on 3 dataset sources without impairing accuracy on the others. To interpret the strategy that URT learns to weigh the backbones from different domains, we visualize the attention scores for both seen and unseen domains and find that our model generates meaningful weights for the pre-trained domains. A comprehensive analysis on variants and ablations of the URT layer is provided to show the importance of various components of URT, notably the number of attention heads.

2 Few-Shot Classification

2.1 Problem Setting

In this section, we will introduce the problem setting for few-shot classification and the formulation of meta-learning for few-shot classification. Few-shot classification aims to classify samples where only few examples are available for each class. We describe a few-shot learning classification task as the pair of examples, comprising of a support set S to define the classification task and the query set Q of samples to be classified.

Meta-learning is a technique that aims to model the problem of few-shot classification as learning to learn from instances of few-shot classification tasks. The most popular way to train a meta-learning model is with episodic training. Here, tasks $T = (Q, S)$ are sampled from a larger dataset by taking subsets of the dataset to build a support set S and a query set Q for the task. A common approach is to sample N -way- K -shot tasks, each time selecting a random subset of N classes from the original dataset and choosing only K examples for each class to add to the support set S .

The meta-learning problem can then be formulated by the following optimization:

$$\min_{\Theta} \mathbb{E}_{(S,Q) \sim p(T)} [\mathcal{L}(S, Q, \Theta)], \quad \mathcal{L}(S, Q, \Theta) = \frac{1}{|Q|} \sum_{(\mathbf{x}, y) \sim Q} -\log p(y|\mathbf{x}, S; \Theta) + \lambda \Omega(\Theta), \quad (1)$$

where $p(T)$ is the distribution of tasks, Θ are the parameters of the model and $p(y|\mathbf{x}, S; \Theta)$ is the probability assigned by the model to label y of query example \mathbf{x} (given the support set S), and $\Omega(\Theta)$ is an optional regularization term on the model parameters with factor λ .

Conventional few-shot classification targets the setting of N -way- K -shot, where the number of classes and examples are fixed in each episode. Popular benchmarks that follow this approach include Omniglot [8] or benchmarks made of subsets of ImageNet, such as *mini*ImageNet [22] and *tiered*ImageNet [15]. In such benchmarks, the tasks used for training cover a set of classes that is disjoint from the classes found in the test set of tasks. However, with the training and test sets tasks coming from a single dataset/domain, the distribution of tasks found in either sets is similar and lacks variability, which may be unrealistic in practice.

It is in this context that Triantafillou et al. [20] proposed Meta-Dataset, as a further step towards large-scale, multi-domain few shot classification. Meta-Dataset includes ten datasets (domains), with eight of them available for training. Additionally, each task sampled in the benchmark varies in the

number of classes N , with each class also varying in the number of shots K . As in all few-shot learning benchmarks, the classes used for training and testing do not overlap.

2.2 Background and Related Work

Transfer by fine-tuning A simple and effective method for few-shot classification is to perform transfer learning by first learning a neural network classifier on all data available for training and using its representation to initialize and then fine-tune neural networks on the few-shot classification tasks found at test time [2, 20, 4, 17]. Specifically, Saikia et al. [17] have shown that competitive performance can be reached using a strong hyper-parameter optimization method applied on a carefully designed validation metric appropriate for few-shot learning.

Meta-Learning Another approach is to use meta-learning to more directly train a model to learn to perform few-shot classification, in an end-to-end way. A large variety of such models have been explored, inspired by memory-augmented networks [18], LSTMs [13] and metric-based classifiers [22]. The two most popular methods are Prototypical Networks [19] and Model Agnostic Meta-Learning (MAML) [6]. Prototypical Networks assume that every class can be represented as a prototype in a learned embedding space (represented by a neural network). Prototypes are calculated as the average of the representations of the support examples of each class. A 1-nearest centroid classifier is then adopted for classification and the neural network representation is trained to facilitate classification in few-shot tasks directly. MAML models the procedure of learning to learn as a bilevel optimization, where an outer loop backpropagates loss gradients through an optimization-based inner loop to learn its initialization. Triantafillou et al. [20] showed that prototypical networks and MAML could be combined by leveraging prototypes for the initialization of the output weights value in the inner loop. Requeima et al. [16] also proposed Conditional Neural Adaptive Processes (CNAPs) for few-shot classification, which can be seen as extending prototypical networks with a more sophisticated architecture that allows for improved task adaptation.

Universal Representations In contrast, our work instead builds on that of Dvornik et al. [5] and their method SUR (Selecting from Universal Representations). Bilen and Vedaldi [1] introduced the term *universal representation* to refer to a representation that supports good performance in multiple domains. One proposal towards such a representation is to train different neural networks backbones separately on the data of each available domain, then simply to concatenate the representation learned by each. Another is to introduce some parameter sharing between the backbones, by having a single network conditioned on the domain of the provenance of each batch of training data [14], e.g. using Feature-wise Linear Modulate (FiLM) [12]. SUR proposes to leverage a universal representation in few-shot learning tasks with a feature selection procedure that assigns different weights to each of the domain-specific subvectors of the universal representation. The objective is to assign high weights only to the domain-specific representations that are specifically useful for each few-shot task at hand. The weights are inferred by optimizing a loss on the support set that encourages high accuracy of a nearest-centroid classifier. As such, the method does not involve any meta-learning—a choice motivated by the concern that meta-learning may struggle in generalizing to domains that are dissimilar to the training domains. SUR achieved state-of-the-art performance on Meta-Dataset. However, a contribution of our work is to provide evidence that meta-learning can actually be used to replace SUR’s hand-designed inference procedure and improve performance further.

Transformer Networks Our meta-learning approach to leverage universal representations is inspired directly from Transformer networks [21]. Transformer networks are neural network architectures characterized by the use self-attention mechanisms to solve tasks. Our model structure is inspired by the structure of the dot-product self-attention in the Transformer, which we adapted here to multi-domain few-shot learning by designing appropriate parametrizations for queries, keys and values. Self-attention was explored in the single-domain training regime by Ye et al. [23], however for a different purpose, where each representation of individual examples in a task support set is influenced by all other examples. Such a strategy is also employed by CNAPs, but with the latter using FiLM as the conditioning mechanism, instead of self-attention. Regardless, the aim of this paper is to propose a different strategy. Rather than using self-attention between individual examples in the support set, our model uses self-attention to select between different domain-specific backbones.

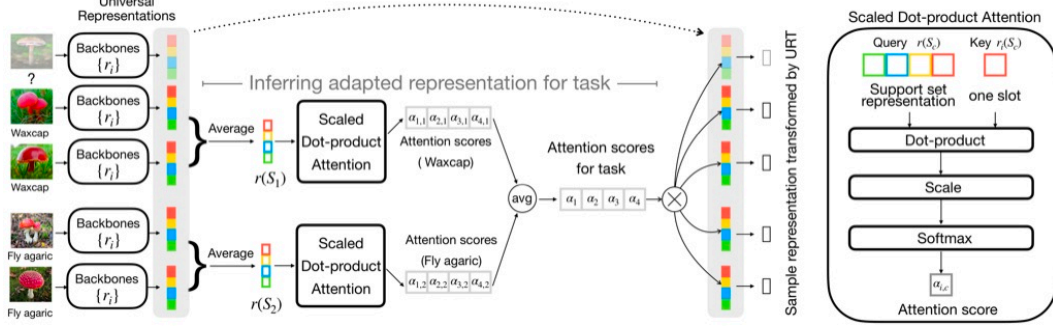


Figure 1: Illustration of how a single-head URT layer uses a universal representation to produce a task-specific representation. This example assumes the use of four backbones, with each color illustrating their domain-specific sub-vector representation in the universal representation.

3 Universal Representation Transformer Layer

In this section, we describe our proposed URT layer, which uses meta-learning episodic training to learn how to combine the domain-specific backbones of a universal representation for any given few-shot learning classification task.

Conceptually, the proposed model views the support set S of a task as providing information on how to query and retrieve from the set $\{r_i\}$ of m pre-trained backbones the most appropriate backbone to build an adapted representation ϕ for the task.

We would like the model to support a variety of strategies on how to retrieve backbones. For example, it might be beneficial for the model to retrieve a single backbone from the set, especially if the domain of the given task matches perfectly that of a domain found in the training set. Alternatively, if some of the training domains benefit from much more training data than others, a better strategy might be to attempt some cross-domain generalization towards the few-shot learning task by blending many backbones together, even if none matches the domain of the task perfectly.

This motivates us to use dot-product self-attention, inspired by layers of Transformer networks [21]. For this reason, we refer to our model as a Universal Representation Transformer (URT) layer. Additionally, since each class of the support set might require a different strategy, we perform attention separately for each class and their support set $S_c = \{x | (x, y) \in S \text{ and } y = c\}$.

3.1 Single-Head URT Layer

We start by describing an URT layer consisting of a single attention head. An illustration of a single-head URT layer is shown in Figure 1. Let $r_i(x)$ be the output vector of the backbone for domain i . We then write the universal representation as

$$r(x) = \text{concat}(r_1(x), \dots, r_m(x)). \quad (2)$$

This representation provides a natural starting point to obtain a representation of a support set class. Specifically, we will note

$$r(S_c) = \frac{1}{|S_c|} \sum_{x \in S_c} r(x) \quad (3)$$

as the representation for the set S_c . From this, we can describe the URT layer by defining the queries³, keys, the attention mechanism and output of the layer:

Queries q_c : For each class c , we obtain a query through $q_c = W^q r(S_c) + b^q$, where we have a learnable query linear transformation represented by matrix W^q and bias b^q .

³Unable to avoid the unfortunate double usage of the term "query" due to conflicting conventions, we highlight the difference between the *query* sets Q of few-shot tasks and the *queries* q_c of an attention mechanism.

Algorithm 1 Training of URT layer

Input: Number of tasks τ_{total} , m pre-trained backbones ;

```
1: for  $\tau \in \{1, \dots, \tau_{total}\}$  do
2:   Sample a few-shot task  $T$  with support set  $S$  and query set  $Q$ ;
3:   # Infer adapted representation for task from  $S$ 
4:   For each class, obtain representation using  $m$  pre-trained backbones as in Eq. (3);
5:   Obtain attention scores using Eq. (4,5) for each head using support set  $S$ ;
6:   # Use adapted representation to predict labels in  $Q$  from support set  $S$ 
7:   Compute adapted representation of examples in  $S$  and  $Q$  as in Eq. (6,7);
8:   Compute probabilities of label of examples in  $Q$  using Prototypical Network as in Eq. (9);
9:   Compute loss as in Eq. (1,8) and perform gradient descent step on URT parameters  $\Theta$ ;
10: end for
```

Keys $\mathbf{k}_{i,c}$: For each domain i and class c , we define keys as $\mathbf{k}_{i,c} = \mathbf{W}^k r_i(S_c) + \mathbf{b}^k$, using a learnable linear transformation \mathbf{W}^k and \mathbf{b}^k and where $r_i(S_c) = 1/|S_c| \sum_{\mathbf{x} \in S_c} r_i(\mathbf{x})$, using a similar notation as for $r(S_c)$.

Attention scores α_i : as for regular Transformer layers, we use scaled dot-product attention

$$\alpha_{i,c} = \frac{\exp(\beta_{i,c})}{\sum_{i'} \exp(\beta_{i',c})}, \beta_{i,c} = \frac{\mathbf{q}_c^\top \mathbf{k}_{i,c}}{\sqrt{l}}, \quad (4)$$

where l is the dimensionality of the keys and queries. Then, these per-class scores are aggregated to obtain scores for the full support set by averaging

$$\alpha_i = \frac{\sum_c \alpha_{i,c}}{N}. \quad (5)$$

Equipped with these attention scores, the URT layer can now produce an adapted representation for the task (for the support and query set examples) by computing

$$\phi(\mathbf{x}) = \sum_i \alpha_i r_i(\mathbf{x}). \quad (6)$$

As we can see, this approach has the flexibility of either selecting a single domain-specific backbone (by assigning $\alpha_i = 1$ for a single domain) or blending different domains together (by having $\alpha_i \gg 0$ for multiple backbones).

3.2 Multi-Head URT Layer

The URT layer described so far can only learn to retrieve a single backbone (or blending of backbones). Yet, it might be beneficial to retrieve multiple different (blended) backbones, especially for a few-shot task that would include many classes of varying complexity.

Thus, to achieve such diversity in the adapted representation, we also consider URT layers with multiple heads, i.e. where each head corresponds to the calculation of Equation 6 and each head has its own set of parameters ($\mathbf{W}^q, \mathbf{b}^q, \mathbf{W}^k, \mathbf{b}^k$). Denoting each head now as ϕ_h , a multi-head URT layer then produces as its output the concatenation of all of its heads:

$$\phi(\mathbf{x}) = \text{concat}(\phi_1(\mathbf{x}), \dots, \phi_H(\mathbf{x})). \quad (7)$$

Empirically we found that the randomness in the initialization of head weights alone did not lead to uniqueness and being complimentary between the heads, so inspired by Lin et al. [10], we add a regularizer to avoid duplication of the attention scores:

$$\Omega(\Theta) = \|(\mathbf{A}\mathbf{A}^\top - \mathbf{I})\|_F^2, \quad (8)$$

where $\|\cdot\|_F$ is the Frobenius norm of a matrix and $\mathbf{A} \in \mathbb{R}^{n \times m}$ is the matrix for attention scores, with \mathbf{A}_h being the vector of all scores α_i for head h . The identity matrix \mathbf{I} regularizes each set of attention scores to be more focused so that multiple heads can attend to different domain-specific backbones.

Table 1: Test performance (mean+CI%95) over 600 few-shot tasks. URT and the most recent methods, which are listed in the first row, are compared on 13 test datasets, which are listed in the first column. The numbers in **bold** have intersecting confidence intervals with the most accurate method.

	ProtoNet[19]	MAML[6]	ProtoMAML[20]	CNAPs[16]	BOHB-E[17]	SUR[5]	URT
ILSVRC	44.5 ± 1.1	37.8 ± 1.0	46.5 ± 1.1	52.3 ± 1.0	55.4 ± 1.1	56.3 ± 1.1	55.7 ± 1.0
Omniglot	79.6 ± 1.1	83.9 ± 1.0	82.7 ± 1.0	88.4 ± 0.7	77.5 ± 1.1	93.1 ± 0.5	94.4 ± 0.4
Aircraft	71.1 ± 0.9	76.4 ± 0.7	75.2 ± 0.8	80.5 ± 0.6	60.9 ± 0.9	85.4 ± 0.7	85.8 ± 0.6
Birds	67.0 ± 1.0	62.4 ± 1.1	69.9 ± 1.0	72.2 ± 0.9	73.6 ± 0.8	71.4 ± 1.0	76.3 ± 0.8
Textures	65.2 ± 0.8	64.1 ± 0.8	68.3 ± 0.8	58.3 ± 0.7	72.8 ± 0.7	71.5 ± 0.8	71.8 ± 0.7
QuickDraw	64.9 ± 0.9	59.7 ± 1.1	66.8 ± 0.9	72.5 ± 0.8	61.2 ± 0.9	81.3 ± 0.6	82.5 ± 0.6
Fungi	40.3 ± 1.1	33.5 ± 1.1	42.0 ± 1.2	47.4 ± 1.0	44.5 ± 1.1	63.1 ± 1.0	63.5 ± 1.0
VGGFlower	86.9 ± 0.7	79.9 ± 0.8	88.7 ± 0.7	86.0 ± 0.5	90.6 ± 0.6	82.8 ± 0.7	88.2 ± 0.6
TrafficSigns	46.5 ± 1.0	42.9 ± 1.3	52.4 ± 1.1	60.2 ± 0.9	57.5 ± 1.0	70.4 ± 0.8	69.4 ± 0.8
MSCOCO	39.9 ± 1.1	29.4 ± 1.1	41.7 ± 1.1	42.6 ± 1.1	51.9 ± 1.0	52.4 ± 1.1	52.2 ± 1.1
MNIST	-	-	-	92.7 ± 0.4	-	94.3 ± 0.4	94.8 ± 0.4
CIFAR10	-	-	-	61.5 ± 0.7	-	66.8 ± 0.9	67.3 ± 0.8
CIFAR100	-	-	-	50.1 ± 1.0	-	56.6 ± 1.0	56.9 ± 1.0

3.3 Training Strategy

We train representations produced by the URT layer by following the approach of Prototypical Networks [19], where the probability of a label y for a query example \mathbf{x} given the support set of a task is modeled as:

$$p(y = c | \mathbf{x}, S; \Theta) = \frac{\exp(-d(\phi(\mathbf{x}) - \mathbf{p}_c))}{\sum_{c'=1}^N \exp(-d(\phi(\mathbf{x}) - \mathbf{p}_{c'}))}, \quad (9)$$

where d is a distance metric and $\mathbf{p}_c = 1/|S_c| \sum_{\mathbf{x} \in S_c} \phi(\mathbf{x})$ corresponds to the centroid of class c , referred to as its prototype. We use (negative) cosine similarity as the distance. The full training algorithm is presented in Algorithm 1.

4 Experiments

In this section, we seek to answer three key experimental questions:

- Q1** How does URT compare with previous state-of-the-art on multi-domain few-shot classification?
- Q2** Do the URT attention heads generate interpretable and meaningful attention scores?
- Q3** Does the URT layer provide consistent benefits, even when pre-trained backbones are trained in different ways?

In addition, we investigate architectural choices made, such as our models for keys/queries and their regularization, and study their contribution to achieving strong performance with URT.

4.1 Datasets and Setup

We test our methods on the large-scale few-shot learning benchmark Meta-Dataset [20]. It consists of ten datasets with various data distributions across different domains, including natural images (Birds, Fungi, VGG Flower), hand-written characters (Omniglot, Quick Draw), and human created objects (Traffic Signs, Aircraft). Among the ten datasets, eight provide data that can be used during either training, validation and testing (with each class assigned to only one of those sets), while two datasets are solely used for testing. Following Requeima et al. [16], we also include MNIST [9], CIFAR10 and CIFAR100 [7] as additional unseen test datasets. Following Triantafillou et al. [20], few-shot tasks are sampled with varying number of classes N , varying number of shots K and class imbalance. The performance is reported as the average accuracy over 600 sampled tasks. More details of Meta-Dataset can be found in Triantafillou et al. [20].

The domain-specific backbones are pre-trained following the setup in [5]. Then, we freeze the backbone and train the URT layer for 10,000 episodes, with an initial learning rate of 0.01 and a cosine

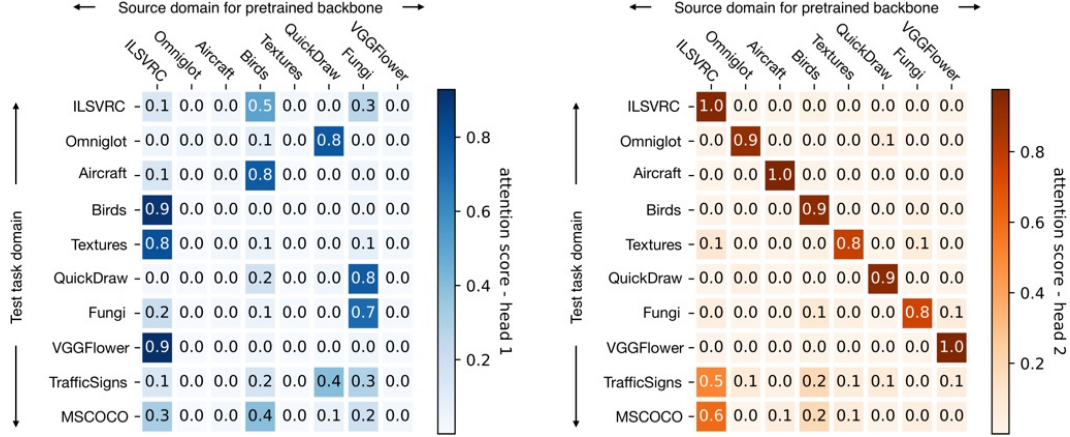


Figure 2: Average attention scores generated by URT with two heads. Rows correspond to the domain of the test tasks and the columns correspond to the pre-trained backbones $r_i(x)$ trained on the eight training domains.

learning rate scheduler. Following Chen et al. [3], the training episodes have 50% probability coming from the ImageNet data source. Since different pre-trained backbones may produce representations with different vector norms, we normalize the outputs of the backbones as in Dvornik et al. [5]. URT is trained with parameter weight decay of $1e-5$ and with a regularization factor $\lambda = 0.1$. The number of heads (H in Equation 7), is set to 2 and the dimension of the keys and queries (l in Equation 4) is set to 1024. We choose the hyper-parameters based on the performance of the validation set. Details of the hyper-parameter selection and how the performance is influenced by them are outlined in Section 4.5.

4.2 Comparison with Previous Approaches

Table 1 presents a comparison of URT with SUR, as well as other baselines based on transfer learning by fine-tuning [17] or meta-learning (Prototypical Networks [19], first-order MAML [6], ProtoMAML [20], CNAPs [16]). Among the listed datasets, eight (above the middle line) have some of their classes used for training and five (below the middle line) do not.

We observe that URT establishes a new state-of-the-art, by outperforming SUR on 3 datasets without compromising performance on others, which is challenging to achieve in the multi-domain setting. Of note, the average inference time for URT is 0.04 second per task, compared to 0.43 for SUR, on a single V100. Thus, getting rid of the optimization procedure for every episode with our meta-trained URT layer also significantly increases the latency, by more than $10\times$. Additionally, URT outperforms all other meta-learning methods, on all datasets.

4.3 Interpreting and Visualizing Attention by URT

To better understand how the URT model of Section 4.2 uses its two heads to build adapted representations, we visualize the attention scores produced on the test tasks of Meta-Dataset in Figure 2.

The blue (first head) and orange (second head) heatmaps summarize the values of the attention scores (Equation 5), averaged across several tasks for each test domain. Specifically, the element on row t and column i is the averaged attention scores α_i computed on test set domain t for the backbone from domain i . Note that the last two rows are the two unseen domain datasets. We found that for datasets from the seen domains, i.e. the first eight rows, one head (right, orange) consistently puts most of its weight on the backbone pre-trained on the same domain, while the other head (left, blue) learns relatively smoother weight distributions that blends other related domains. For unseen datasets, the right head puts half of its weight on ImageNet and the left head learned to blend the representations from four backbones.

4.4 URT using FiLM Modulated Backbones

As additional evidence of the benefit of URT on universal representations, we also present experiments based on a different set of backbone architectures. Following SUR [5], we consider the backbones from a parametric network family, obtained by training a base backbone on one dataset (ILSVRC) and then learning separate FiLM layers [12] for each other dataset, to modulate the backbone so it is adapted to the other domains. These backbones collectively have only 0.5% more parameters than a single backbone.

A comparison between SUR and URT using these backbones (referred to as SUR-pf and URT-pf) is presented in Table 2. Once again, URT improves the performance on three datasets without sacrificing performance on others. Additionally, URT-pf now achieves better performance than BOHB-E on VGGFlower.

Table 2: Performance comparison using parametric network family (pf) backbones.

	SUR-pf [5]	URT-pf	VS.
ILSVRC	56.4 \pm 1.2	55.5 \pm 1.1	=
Omniglot	88.5 \pm 0.8	90.2 \pm 0.6	+
Aircraft	79.5 \pm 0.8	79.8 \pm 0.7	=
Birds	76.4 \pm 0.9	77.5 \pm 0.8	=
Textures	73.1 \pm 0.7	73.5 \pm 0.7	=
Quick Draw	75.7 \pm 0.7	75.8 \pm 0.7	=
Fungi	48.2 \pm 0.9	48.1 \pm 0.9	=
VGG Flower	90.6 \pm 0.5	91.9 \pm 0.5	+
Traffic Signs	65.1 \pm 0.8	67.5 \pm 0.8	+
MSCOCO	52.1 \pm 1.0	52.1 \pm 1.0	=
MNIST	93.2 \pm 0.4	93.9 \pm 0.4	=
CIFAR10	66.4 \pm 0.8	66.1 \pm 0.8	=
CIFAR100	57.1 \pm 1.0	57.3 \pm 1.0	=

4.5 Hyper-Parameter and Ablation Studies

We analyze the importance of the various components of URT’s attention mechanism structure and training strategy in Table 3. First we analyze the importance of using the support set to model queries and/or keys. To this end, we consider setting the matrices \mathbf{W}^q / \mathbf{W}^k of the query / key linear transformation to 0, which only leaves the bias term. We found that the support set representation is most crucial for building the keys (row w/o \mathbf{W}^k in the table) and has minor benefits for queries (row w/o \mathbf{W}^q) in the table. This observation is possibly related to the success of attention-based models with learnable constant queries [11, 10]. We also found that adding a regularizer $\Omega(\Theta)$ as in Equation 8 is important for some datasets, specifically VGG Flower and Birds.

Table 3: Meta-Dataset performance variation on ablations of elements of the URT layer.

	ILSVRC	Omniglot	Aircraft	Birds	Textures	Draw	Fungi	Flower	Signs	MSCOCO
w/o \mathbf{W}^q	+0.2	-0.2	-0.6	-0.1	-0.3	-0.2	0.0	-0.2	-0.8	-0.1
w/o \mathbf{W}^k	-14.2	-2.8	-10.7	-18.1	-7.6	-9.3	-22.4	-3.6	-0.26	-10.9
w/o $r(S_c)$	-14.2	-2.8	-10.7	-18.1	-7.6	-9.2	-22.4	-3.6	-0.26	-10.9
w/o $\Omega(\Theta)$	0.0	-0.9	-0.4	-3.3	-1.2	-0.2	+0.3	-9.0	-2.0	0.0

An important hyper-parameter in URT is the number of heads H . We chose this hyper-parameter based on the performance on validation set of tasks in Meta-Dataset. In Table 4, we show the validation performance of URT for varying number of heads. As suggested by Triantafillou et al. [20], we considered looking at the rank of the performance achieved by each choice of H for each validation domains, and taking the average across domains as a validation metric. However, since the performances when using two to four heads are similar and yield the same average rank, we instead simply consider the average accuracy as the selection criteria.

Table 4: Validation performance on Meta-Dataset using different number of heads

	1	2	3	4	5	6	7	8
Average Accuracy	74.605	77.145	76.943	76.984	76.602	75.906	75.454	74.473
Average Rank	2.875	1.000	1.000	1.000	2.250	2.250	2.25	2.50

In general, we observe a large jump in performance when using multiple heads instead of just one. However, since the number of heads controls the capacity, predictably we also observe that having too many heads leads to overfitting.

5 Conclusion

We proposed the URT layer to effectively integrate representations from multiple domains and demonstrated improved performance in multi-domain few-shot classification. Notably, our URT approach was able to set a new state-of-the-art on Meta-Dataset, while also being $10\times$ more efficient at inference compared to the next-best approach (SUR). This work suggests that combining meta-learning with pre-trained universal representations is a promising direction for new few-shot learning methods. Specifically, we hope that future work can investigate the design of richer forms of universal representations that go beyond simply pre-training a single backbone for each domain, and developing meta-learners adapted to those settings.

Broader Impact

Our URT model may present an interesting element of solution for applications that present difficulties in the collection and sharing of data. This could include settings where each user of an application has limited private data, and as such desires that a classification task be executed directly and solely on their devices. Any deployment of the proposed model however should be preceded by an analysis of the potential biases captured by the dataset sources used for training and the correction of any such undesirable biases captured by the pre-trained backbones and model.

Acknowledgement

We would like to thank Tianyi Zhou for paper review and suggestions. The computation support for this project is provided by Compute Canada and Google Cloud. This project was supported by the Canada CIFAR AI Chairs program.

References

- [1] Hakan Bilen and Andrea Vedaldi. Universal representations: The missing link between faces, text, planktons, and cat breeds. *arXiv preprint arXiv:1701.07275*, 2017.
- [2] Wei-Yu Chen, Yen-Cheng Liu, Zsolt Liu, Yu-Chiang Frank Wang, and Jia-Bin Huang. A closer look at few-shot classification. In *International Conference on Learning Representations (ICLR)*, 2019.
- [3] Yinbo Chen, Xiaolong Wang, Zhuang Liu, Huijuan Xu, and Trevor Darrell. A new meta-baseline for few-shot learning. *arXiv preprint arXiv:2003.04390*, 2020.
- [4] Guneet Singh Dhillon, Pratik Chaudhari, Avinash Ravichandran, and Stefano Soatto. A baseline for few-shot image classification. In *International Conference on Learning Representations (ICLR)*, 2020.
- [5] Nikita Dvornik, Cordelia Schmid, and Julien Mairal. Selecting relevant features from a universal representation for few-shot classification. *arXiv preprint arXiv:2003.09338*, 2020.
- [6] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *The International Conference on Machine Learning (ICML)*, 2017.
- [7] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009.
- [8] Brenden M Lake, Ruslan Salakhutdinov, and Joshua B Tenenbaum. Human-level concept learning through probabilistic program induction. *Science*, 350(6266):1332–1338, 2015.
- [9] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [10] Zhouhan Lin, Minwei Feng, Cicero Nogueira dos Santos, Mo Yu, Bing Xiang, Bowen Zhou, and Yoshua Bengio. A structured self-attentive sentence embedding. *arXiv preprint arXiv:1703.03130*, 2017.

- [11] Yang Liu, Chengjie Sun, Lei Lin, and Xiaolong Wang. Learning natural language inference using bidirectional lstm model and inner-attention. *arXiv preprint arXiv:1605.09090*, 2016.
- [12] Ethan Perez, Florian Strub, Harm de Vries, Vincent Dumoulin, and Aaron C. Courville. Film: Visual reasoning with a general conditioning layer. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2018.
- [13] Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. In *International Conference on Learning Representations (ICLR)*, 2016.
- [14] Sylvestre-Alvise Rebuffi, Hakan Bilen, and Andrea Vedaldi. Efficient parametrization of multi-domain deep neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [15] 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 *International Conference on Learning Representations (ICLR)*, 2018.
- [16] James Requeima, Jonathan Gordon, John Bronskill, Sebastian Nowozin, and Richard E Turner. Fast and flexible multi-task classification using conditional neural adaptive processes. In *The Conference on Neural Information Processing Systems (NeurIPS)*, pages 7957–7968, 2019.
- [17] Tonmoy Saikia, Thomas Brox, and Cordelia Schmid. Optimized generic feature learning for few-shot classification across domains. *arXiv preprint arXiv:2001.07926*, 2020.
- [18] Adam Santoro, Sergey Bartunov, Matthew Botvinick, Daan Wierstra, and Timothy Lillicrap. Meta-learning with memory-augmented neural networks. In *The International Conference on Machine Learning (ICML)*, 2016.
- [19] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In *The Conference on Neural Information Processing Systems (NeurIPS)*, 2017.
- [20] Eleni Triantafillou, Tyler Zhu, Vincent Dumoulin, Pascal Lamblin, Utku Evci, Kelvin Xu, Ross Goroshin, Carles Gelada, Kevin Swersky, Pierre-Antoine Manzagol, et al. Meta-dataset: A dataset of datasets for learning to learn from few examples. In *International Conference on Learning Representations (ICLR)*, 2020.
- [21] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *The Conference on Neural Information Processing Systems (NeurIPS)*, 2017.
- [22] Oriol Vinyals, Charles Blundell, Tim Lillicrap, Daan Wierstra, et al. Matching networks for one shot learning. In *The Conference on Neural Information Processing Systems (NeurIPS)*, 2016.
- [23] Han-Jia Ye, Hexiang Hu, De-Chuan Zhan, and Fei Sha. Few-shot learning via embedding adaptation with set-to-set functions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.