

Multi-Domain Learning with Modulation Adapters

Ekaterina Iakovleva
Inria

Karteek Alahari
Inria

Jakob Verbeek
Meta AI

Abstract

Deep convolutional networks are ubiquitous in computer vision, due to their excellent performance across different tasks for various domains. Models are, however, often trained in isolation for each task, failing to exploit relatedness between tasks and domains to learn more compact models that generalise better in low-data regimes. Multi-domain learning aims to handle related tasks, such as image classification across multiple domains, simultaneously. Previous work on this problem explored the use of a pre-trained and fixed domain-agnostic base network, in combination with smaller learnable domain-specific adaptation modules. In this paper, we introduce Modulation Adapters, which update the convolutional filter weights of the model in a multiplicative manner for each task. Parameterising these adaptation weights in a factored manner allows us to scale the number of per-task parameters in a flexible manner, and to strike different parameter-accuracy trade-offs. We evaluate our approach on the Visual Decathlon challenge, composed of ten image classification tasks across different domains, and on the ImageNet-to-Sketch benchmark, which consists of six image classification tasks. Our approach yields excellent results, with accuracies that are comparable to or better than those of existing state-of-the-art approaches.

1. Introduction

Deep learning models are ubiquitous across many tasks and application domains [20]. Yet, one of their biggest limitations is the reliance on large amounts of labelled data to train these models, which often have millions if not billions of parameters. Knowledge transfer is one of the main solutions to this problem, where knowledge obtained by learning to solve one task is leveraged to learn another task in a more sample-efficient manner. A popular example of knowledge transfer is pre-training of a deep neural network on a large (un)labelled dataset, e.g., on ImageNet [33], and then fine-tuning it on a smaller target dataset.

Beyond pre-training and fine-tuning, there is a wide

range of learning paradigms that involve some form of knowledge transfer. In *multi-domain learning* (MDL) [30], which is the focus of this paper, a single model is trained to solve several related tasks on a number of different domains, e.g., digit recognition and object recognition. The idea is to use knowledge from related tasks to induce regularisation across the considered domains. The common view on MDL is to transfer knowledge across domains to maximise predictive accuracy, while sharing as many parameters as possible across tasks, and adding as few task-specific parameters as possible. A common approach is to learn a base network on one of the domains, and then adapt it with a small number of parameters for tasks on other domains. A wide range of solutions to such base network adaptations has been proposed in the literature, which can be divided into two major groups of approaches. The first group applies a task-specific binary mask to either the features or the weights in each layer of the model [3, 24, 26]. While this adds a small number of parameters per task, essentially one bit per weight or feature channel, its ability to adapt to new tasks is limited for the same reason. The second group leaves the weights of the base network unchanged, but adds a number of task-specific adaptation layers, for example, a 1×1 convolution parallel to each 3×3 convolutional layer of the base network [21, 30, 32]. Depending on the design of the adaptation modules, this approach adds more parameters to adapt the model to the task. However, by leaving the parameters of the base network layers unchanged, task adaptation can only be achieved by added model capacity in the adaptation layers themselves, without capitalising on adaptations of all of the backbone parameters.

To overcome the limitations of existing approaches, we propose a novel type of adaptation module, which we call *Modulation Adapters*. Our approach adapts the weights of the existing layers of the base network to the task by means of non-binary scaling. More specifically, to adapt the base network to a new domain d , we learn a set of domain-specific adapters α^d , each modulating the fixed convolutional filters by multiplying them with a scalar α_{mn}^d that is specific to the pair (m, n) of the output and input channels to which the filter applies. See Figure 1c for an illustration.

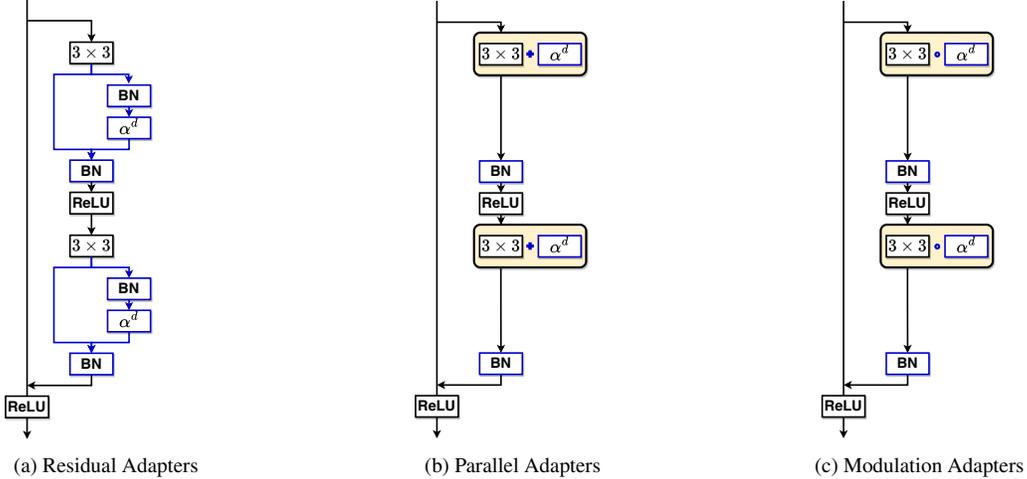


Figure 1. Different adaptation units embedded into a residual block. Blocks in blue, adapters α^d and batch-norm (BN) layers, contain domain-specific parameters. Blocks in black in the base network are domain-agnostic. Yellow blocks show that the 3×3 convolutions of the base network are modulated before being applied. Note that the BN blocks contribute a negligible number of parameters compared to the Modulation Adapter blocks.

The latter allows for nine-fold parameter reduction compared to 3×3 convolutions that constitute the core of common base networks, such as ResNets [13]. To further reduce the memory footprint, we propose a factorised representation of each adapter α^d as a product of two matrices with a smaller intermediate dimensionality. In contrast to other approaches, our adapters allow to modulate all weights in the base network in a non-binary manner, while offering a scalable number of parameters.

We perform several ablative studies to assess the impact of multiplicative adaptation compared to additive, and demonstrate the benefits of the proposed approach. We evaluate our Modulation Adapters on the Visual Decathlon Challenge and the ImageNet-to-Sketch benchmark. We obtain excellent results, with accuracies that are comparable or better than those of existing state-of-the-art approaches across a wide range of parameter budgets.

2. Related work

Weight and feature masking. One way to approach multi-domain learning is to adapt the convolutional layers of a base network that has been pre-trained on a large dataset. The “piggy-back” approach [24] learns element-wise binary masks for the weights in the convolutional layers of the base network. Such parameterisation allows memory-efficient storage of domain-specific parameters: requiring to store 1 bit for each real-valued weight, i.e., a 32 fold reduction in storage for single-precision parameters. The idea of adaptation through binary masks is further generalised in WTPB [26], where the authors use three scalars to define an affine transformation of convolutional

filter weights based on the binary masks. In BA² [3], binary masks are applied across feature channels rather than network weights. As a result, after training the masked feature maps can be removed from the computational graph, reducing both the memory footprint and the computational complexity. While reduction in memory budget is an appealing property of the masking approach, the common downside is limited capability of model adaptation due to simplicity of the binary masks. In our work, we propose a method for non-binary multiplicative adaptation of the base network, which allows to span a large range of parameter budgets.

Adaptation modules. Domain-specific trainable modules to adapt a pre-trained base network have been explored in several methods. Residual Adapters [30] use an intermediate residual block with trainable 1×1 convolution after each pre-trained and fixed 3×3 convolution. Compared to the base network, such types of adapters lead to nine-fold reduction in the number of additional parameters learned for a new domain. In Parallel Adapters [31], a domain-specific 1×1 convolution is used as an additive bypass to each domain-agnostic 3×3 convolution. The Parallel Adapter configuration provides more flexibility than the sequential one, since it can be used on top of existing pre-trained networks as the adapters need not be present when pre-training. DAN [32] learns a linear transformation of the output filters in the base convolutions, which is equivalent to insertion of intermediate 1×1 convolution after each base convolution. Unlike Residual Adapters, this model does not use additional skip connections and batch normalisation layers during adaptation. CovNorm [21] learns approximation of adaptation layers which consists of whitening, mini-adaptation and colouring operations. For each

trained adapter, the whitening and colouring matrices are obtained from PCA of corresponding input and output activations, while the mini-adaptation layer is learned from scratch. Our method also falls into the category of models which contain adaptation modules. Different from the other approaches, we perform modulation of fixed convolutional filters by scaling their weights with a non-binary factor specific to each input-output channel combination.

Adapter modules have been investigated in the context of NLP [14] to specialize Transformer models [38] to specific tasks. In the context of generative image models, adapters have been used to adapt pre-trained GANs to new domains with few samples [2, 9]. Our model is related to weight modulation approaches for GANs used in StyleGAN2 [15] and in CISP [1].

Other approaches. Besides the two main lines of work described above, several other directions have been explored. SpotTune [11] uses a policy network that for each data sample from a new domain, and for each convolutional layer, predicts whether the weights should be fine-tuned or re-used from the pre-trained network. There is, however, no incentive to reduce the number of parameters. In the work of Guo *et al.* [10], the authors replace standard convolutions in the base network with depth-wise separable convolutions [6]. Parameter-heavy 1×1 convolutions are then shared across domains, while the per-feature 3×3 convolutions with less parameters per task are learned. Such parameterisation of convolutions significantly decreases the size of the base network, as well as the number of parameters learned for each domain. Approach of Senhaji *et al.* [35] divides the fixed pre-trained base network with Parallel Adapters learned for a new domain into three non-overlapping blocks, and place an exit module after each of them. Each exit module takes the output of the previous block as input, and produces the softmax output for the domain of interest. The authors learn all exit modules and choose the one with the highest performance.

3. Method

We consider the multi-domain learning problem where a model is required to solve D classification tasks on D data domains, while minimising the total number of parameters being learned. We begin with the popular framework, e.g., [24, 26, 30], of pre-training a base feature extractor on a large dataset, which is then fixed and shared across all domains. We then learn a separate adapter for each domain, which modifies the feature extractor to improve its performance on the given domain. Along with the adapters, we also learn domain-specific batch normalisation layers and a classification head which predicts the domain-specific labels. The objective function is a sum of domain-specific objective functions that correspond to individual tasks: in

our case the sum of corresponding cross-entropy losses.

We detail our approach and its properties in Section 3.1, and discuss the factored parametrisation in Section 3.2. Then, in Section 3.3, we provide a detailed comparison of our approach with respect to the most related prior works.

3.1. Modulation Adapters

We suppose the CNN base network has already been pre-trained. Let $\mathbf{f} \in \mathbb{R}^{M \times N \times K \times K}$ be the filter of a specific layer in the base network, where K is the spatial filter size, N is the number of input features, and M is the number of output features.

To adapt the convolutional layers of the base network to a new domain $d \in \{1, \dots, D\}$, we define a modulation adaptation unit ρ , with domain specific parameters $\alpha^d \in \mathbb{R}^{M \times N}$ that modulate the convolutional filter in a multiplicative manner. For an input \mathbf{x} with N feature maps, this unit produces an output \mathbf{y} with M feature maps:

$$\mathbf{y} = \rho(\mathbf{f}, \mathbf{x}, \alpha^d) = (\mathbf{f} \circ \alpha^d) * \mathbf{x}, \quad (1)$$

where $*$ is the convolutional operator, and \circ is an element-wise product with the adapter parameters $\alpha^d \in \mathbb{R}^{M \times N}$ broadcasted across the spatial dimensions of the convolutional weight tensor $\mathbf{f} \in \mathbb{R}^{M \times N \times K \times K}$. Adaptation module ρ applied to the pre-trained filter \mathbf{f} is equivalent to a new convolutional filter \mathbf{g}^d , with the weights obtained as:

$$[\mathbf{g}]_{mnkl}^d = \alpha_{mn}^d \cdot [\mathbf{f}]_{mnkl}. \quad (2)$$

Note that modulation of the filter weights for each domain can also be interpreted as applying scaling of the input features, which is specific for each output channel. Therefore, with modulation each output channel can mix the input channels in different ways, e.g., amplify some and ignore others. The input channels are, however, processed by the same spatial filter for each input-output channel combination, leading to reduction in number of trainable parameters relative to the size of the base network.

We learn a separate set of modulation weights for each convolutional layer in the base network, and denote our adaptation unit as a *Modulation Adapter* (MAD). Figure 1c shows how Modulation Adapter units are incorporated into a single residual block. While all convolutional weights are altered for a new domain, we only need to store the pre-trained filters \mathbf{f} as well as the Modulation Adapters, which are learned from scratch. The latter contain MN parameters for each convolution in the base model. This is a reduction in the number of parameters by a factor K^2 (the spatial filter size), relative to corresponding convolution in the base network (e.g., for 3×3 convolutions $K^2 = 3 \times 3 = 9$).

3.2. Factorisation of Modulation Adapters

The number of parameters in our Modulation Adapters is given by the product of the number of input channels and

the number of output channels of a convolutional layer. For residual blocks the number of input and output channels is the same, and thus the number of adaptation parameters is quadratic in the number of channels. Therefore, the number of adaptation parameters is typically largest in deeper layers of the base network which tend to contain hundreds of feature channels. For example, the ResNet-26 architecture we use in our experiments has 256 channels in the deeper layers, resulting in 64k parameters in each adaptation module. An alternative, allowing for linear scaling of the number of parameters with the number of feature channels, would be to simply scale the input and output feature channels of the convolutional layer. This scaling, however, could be absorbed in the domain specific BatchNorm layers, which has been found to be too restrictive for effective adaptation in previous work [30].

Inspired by parameter reduction demonstrated by matrix decomposition in CovNorm [21], we adopt a similar strategy in our work. However, instead of learning an unconstrained full adapter and approximating it with a product of two (or more) matrices using some form of matrix decomposition, e.g. SVD as in ConvNorm, we propose to directly represent and learn MAD as a product of two matrices, each with a smaller intermediate dimension $I < \min(M, N)$:

$$\alpha^d = \beta^d \times \gamma^d, \quad (3)$$

where $\beta^d \in \mathbb{R}^{M \times I}$, $\gamma^d \in \mathbb{R}^{I \times N}$ and \times denotes matrix multiplication. Such representation limits the maximum rank of the MAD to I , and can be interpreted as a rank factorisation. In practice, we observe that for certain base networks the full Modulation Adapters are, indeed, learned to be sparse, which we show in Appendix B. This justifies restriction of the rank promoted by the proposed representation, while direct learning of the factors reduces the memory requirements during training.

We learn the factors β^d and γ^d from scratch, and compute their product prior to scaling the base filter bank \mathbf{f} . The domain-specific convolution \mathbf{g}^d is obtained as:

$$[\mathbf{g}]_{mnkl}^d = \left(\sum_{i=1}^I \beta_{mi}^d \cdot \gamma_{in}^d \right) \cdot [\mathbf{f}]_{mnkl}. \quad (4)$$

In this case MAD can still be interpreted as scaling the input channels in a specific way for each output channel, but the scaling coefficients are no longer completely independent. This adapter factorisation reduces the number of parameters being trained from MN in the full adapter to $I(M + N)$ in the factorised case. The intermediate dimension I is a hyper-parameter that can be used to trade-off the number of parameters with prediction accuracy.

3.3. Comparison to related approaches

Some earlier work described in Section 2 can also be understood in terms of domain specific filters \mathbf{g}^d obtained

with domain specific adapters α^d in order to modify the filters \mathbf{f} of the base network. This allows us to further clarify the differences and similarities between these approaches and our Modulation Adapters. In this discussion, both \mathbf{g}^d and \mathbf{f} have the same shapes as before, the form of α^d varies across approaches.

Masking methods apply element-wise scaling of convolutional weights:

$$[\mathbf{g}]_{mnkl}^d = \alpha_{mnkl}^d \cdot [\mathbf{f}]_{mnkl}, \quad (5)$$

where the scaling coefficients α_{mnkl}^d are either binary [24], or take two unique non-binary values [26]. Rather than masking individual weights, BA² [3] masks the entire features:

$$[\mathbf{g}]_{mnkl}^d = \alpha_m^d \cdot [\mathbf{f}]_{mnkl}. \quad (6)$$

Although these adapters are compact to store, their binary nature and factored scaling of features limit task adaptation. Non-binary weights in our Modulation Adapters provide more degrees of freedom for adaptation, while our factorisation approach allows to control the parameter budget.

Linear feature combination methods [32] learn adapters that linearly combine filter outputs:

$$[\mathbf{g}]_{mnkl}^d = \sum_{i=1}^M \alpha_{mi}^d \cdot [\mathbf{f}]_{inlk}, \quad (7)$$

which is equivalent to inserting a domain specific 1×1 convolution after the (fixed) convolution of the base network. Residual Adapters [30], depicted in Figure 1a, add a skip-connection on top of linear combination of features:

$$[\mathbf{g}]_{mnkl}^d = \sum_{i=1}^M (1 + \alpha_{mi}^d) \cdot [\mathbf{f}]_{inlk}. \quad (8)$$

Linear combination approaches have comparable parameter budget as our (non-factored) Modulation Adapters: M^2 for the former, and MN for the latter. These adapters are implemented as additional linear layers applied to the output of the corresponding unaltered convolutions. Our approach can not be viewed as a linear transformation applied to either inputs or outputs, since we perform element-wise multiplication of adapters with parameters of convolutions.

Parallel Residual Adapters [31], depicted in Figure 1b, adopt a domain specific 1×1 convolution in parallel to the fixed 3×3 convolutions of the base network. This is equivalent to adapting the central element of the pre-trained filters:

$$[\mathbf{g}]_{mnkl}^d = [\mathbf{f}]_{mnkl} + \begin{cases} \alpha_{mn}^d & \text{if } k = l = (K - 1)/2 + 1, \\ 0 & \text{otherwise.} \end{cases} \quad (9)$$

While the number of adapter parameters is MN as in our approach, this adapter only affects the central elements of the filters, which is restrictive and limits the space of possible adaptations.

4. Experimental Evaluation

We describe our experimental setup in Section 4.1. After that, we present ablation experiments in Section 4.2, followed by comparison to previous work in Section 4.3.

4.1. Experimental setup

Datasets. In our experiments we use the two most common multi-domain learning benchmarks: the Visual Decathlon Challenge [30] and ImageNet-to-Sketch [24], which contain image classification datasets from heterogeneous visual domains, and each dataset contains different classes to recognise.

The *Visual Decathlon Challenge* consists of ten image classification datasets: ImageNet [33], CIFAR-100 [18], Aircraft [23], Daimler pedestrian classification [27], Describable textures [7], German traffic signs [37], Omniglot [19], Street view house numbers [28], UCF101 Dynamic Images [4, 36], and VGG-Flowers [29].

The second benchmark, *ImageNet-to-Sketch*, consists of six image classification datasets: ImageNet [33], VGG-Flowers [29], Stanford Cars [17], Caltech-UCSD Birds [39], Skteches [8] and WikiArt [34].

For the base network architectures we follow previous work. For the Visual Decathlon challenge we use the ResNet26 architecture [21, 30, 32], and for ImageNet-to-Sketch we use the DenseNet-121 architecture [3, 25, 26]. More details on the network architectures, data processing and training details are described in Appendix A.

Evaluation metrics. We report average classification accuracy across all domains, as well as domain-specific accuracies. Each model is trained 10 times with random initialisations, and the average across these runs is reported. We also use the score function introduced by Rebuffi *et al.* [30]. This score S is computed as:

$$S = \sum_{d=1}^D a_d \max\{0, E_d^{\max} - E_d\}^{b_d}, \quad (10)$$

where D is the number of domains in the benchmark, $a_d = 1000 (E_d^{\max})^{-b_d}$ is the weight which ensures a perfect score of 1000 on a single domain, E_d^{\max} and E_d are the baseline and the model’s test errors, respectively, and $b_d = 2$ is the exponent that rewards more significant reductions of the classification error. The baseline error E_d^{\max} is twice the error of a per-domain fully finetuned model. The score of the finetuning baseline is therefore 250 per dataset. We report the total number of trainable parameters relative to the size of the base feature extractor network, not counting the linear classification head.

4.2. Ablation studies

We conduct ablation studies on the Visual Decathlon Challenge in order to assess several variants of our model,

as well as the effect of scaling, selecting the parameter budget, and the complementarity with additive adapters.

Scaling central element vs. entire filter. We experiment with a version of Modulation Adapter in which we only scale the central element of the 3×3 filters, rather than the entire filter. This is similar to additive Parallel Adapters (PA) [31], which use 1×1 filters to adapt the central elements of the 3×3 filters of the base network. With this central-element-only scaling approach, we obtain a mean accuracy of 72.5% and a decathlon score of 2252, see Table 1. This is significantly worse than the performance of our full model (78.7% mean accuracy, and score 3828), showing the benefit of scaling the entire 3×3 filter for each input-output channel combination. Interestingly, scaling only the central element is also worse than the result of Parallel Adapters, which reported a mean accuracy of 78.1% and a score of 3412 (see Table 2). Although scaling the central element is equivalent after training, our hypothesis is that optimising the scaling adapter with stochastic gradient descend is more difficult than optimising the additive adapter. Using an optimiser with adaptive learning rates, such as Adam [16], may alleviate this.

Additive adaptation of all filter elements. In our second ablation, we consider a variant of our approach in which we adapt all filter elements, as in Eq. (2), but in an additive rather than multiplicative manner. This is similar to PA, but now the adapter updates all the filter elements rather than only the central one. This variant leads to a substantially worse mean accuracy of 61.1% and a decathlon score of 916, see line “Additive 3×3 adapt.” in Table 1. Our interpretation of this deterioration is that strong additive adaptation of all filter weights leads to spatially uniform filters, that fail to detect spatial patterns.

Decomposition of Modulation Adapters. In Figure 2, we consider the effect of decomposing our Modulation Adapters as a product of two smaller matrices of weights, as in Eq. (3). We vary the intermediate dimension I from 2 to 92, and plot the mean accuracy against the parameter budget. We note that the mean accuracy quickly increases when allowing more parameters, reaches a maximum at $I = 36$, and then slightly declines. The mean accuracy of the full, non-decomposed model, and the model with $I = 36$ are very close at 78.7% and 78.8% respectively, see Table 1. The decomposed model with $I = 36$, however, yields a saving of more than 60% in the parameter budget.

While on average there is a smooth trend between the accuracy and the parameter budget, the optimal parameter budget differs per dataset. In Table 1 (line “Modulation Adapter (best)”), we show what the optimal accuracy would be if we were to set this budget for each dataset separately. We notice that this only leads to minor improvements of 0.1% and 0.2% w.r.t. the $I = 36$ and the full models, which

	Params	ImNet	Airc	C100	DPed	DTD	GTSR	Flwr	Ogt	SVHN	UCF	Mean	Score
Modulation adapter (full)	2	60.8	66.8	79.2	97.9	56.9	99.2	86.4	90.0	97.0	52.5	78.7 ± 0.3	3828 ± 178
Scaling central element	2	60.8	39.6	74.7	97.7	53.1	98.0	76.5	85.8	94.8	44.0	72.5 ± 0.3	2252 ± 120
Additive 3×3 adapt.	2	60.8	30.5	53.6	91.5	25.3	95.1	44.1	88.8	94.5	27.2	61.1 ± 0.2	916 ± 11
Modulation adapter ($I = 36$)	1.39	60.8	65.3	81.7	98.2	59.4	99.2	84.8	89.5	96.8	52.1	78.8 ± 0.1	3798 ± 134
Modulation adapter (best)	1.52	60.8	65.3	81.7	98.3	59.4	99.3	85.1	89.6	96.8	52.5	78.9 ± 0.1	3878 ± 68
Modulation adapter (LOO)	1.38	60.8	64.9	81.7	98.2	59.1	99.2	84.4	89.5	96.8	52.1	78.7 ± 0.3	3802 ± 123

Table 1. Ablation experiments. Classification accuracy is reported per dataset, along with the mean, and the decathlon score. For the latter two we also report their standard deviations across the ten repetitions of the experiments.

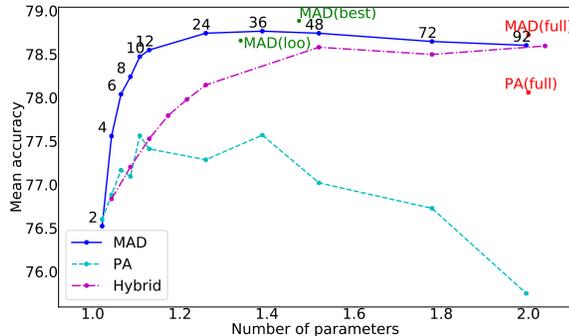


Figure 2. Mean accuracy vs. parameter budget for different adapter-based methods: decomposed version of our Modulation Adapters (MAD), decomposed version of Parallel Adapters (PA, our implementation), hybrid adapters that combine MAD and PA.

is within the standard deviation of the reported results. We further illustrate the stability of the trade-off between accuracy and parameter budget across datasets with a leave-one-out (LOO) experiment. In this case, for each dataset we select the parameter budget by taking the optimal budget on all the other datasets. We then average the results across all the datasets. As in previous cases, we only find minor deviations in the mean accuracy, reaching 78.7%, which is the same as the performance of the full non-decomposed model.

Decomposition of Parallel Adapters. Similar to the decomposition used in our Modulation Adapters, we implemented a decomposition version of Parallel Adapters (PA) [31], where each 1×1 adapter is given as a product of two matrices with smaller intermediate dimension I . As shown in Figure 2, unlike our approach, PA does not behave well with decomposition relative to the full version of PA. Moreover, decomposed PA performs consistently worse than our Modulation Adapters across all parameter budgets. For reference, we also included the “full” non-decomposed version of MAD and PA in the graph.¹

¹Result for PA (full) taken from original paper: 78.1% mean accuracy and score 3412. We implemented the decomposed version of PA (as well as our models) based on the public code-base of PA (full) from [31]. In our experiments, we obtained mean accuracy 77.7% and score 3206 for

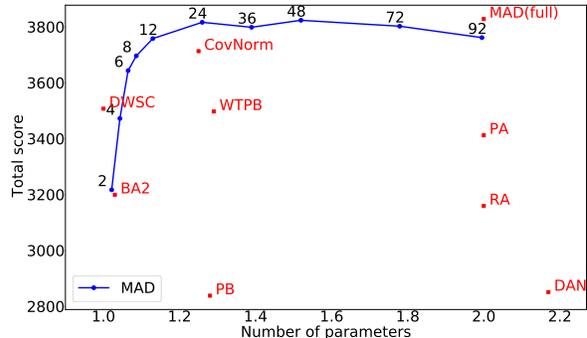


Figure 3. Total score vs. parameter budget for Modulation Adapters (MAD) and the other approaches.

Hybrid adapters. To study the complementarity of multiplicative and additive adaptation, we experiment with a hybrid approach that mixes our Modulation Adapters with Parallel Adapters [31]. We use decomposed versions of MAD and PA, since our goal is to keep the total number of parameters below the size of the full version of a single model. The available parameter budget is evenly split between the two types of adapters. The results in Figure 2 show that these adapters are not complementary: the hybrid approach is consistently worse than our Modulation Adapters, in particular for small parameter budgets, while providing a consistent improvement over the decomposed version of Parallel Adapters.

4.3. Comparison to the state of the art

Visual Decathlon Challenge. In Figure 3, we compare the total score of our full and decomposed Modulation Adapters to state-of-the-art methods across a range of parameter budgets. While other methods provide a single operating point, our decomposed Modulation Adapters enable adaptation to new domains across a wide range of parameter budgets, and yield better or comparable score. In particular, CovNorm [21] has a score 3713, and for a similar parameter budget our Modulation Adapters ($I = 24$) obtain an improved score 3816. Parallel Adapters [31] re-

PA (full). If we take best results across these ten runs, we obtain mean accuracy 78.2% and score 3460, consistent with the original paper.

Number of images	Params	ImNet	Airc	C100	DPed	DTD	GTSR	Flwr	Oglt	SVHN	UCF	Mean	Score	Rank
		1.3m	7k	50k	30k	4k	40k	2k	26k	70k	9k			
Feature [30]	1	59.7	23.3	63.1	80.3	45.4	68.2	73.7	58.8	43.5	26.8	54.3	544	14.7
BN Adapt. [5]	~1	59.9	43.1	78.6	92.1	51.6	95.8	74.1	84.8	94.1	43.5	71.8	1363	13.4
BA2 [3]	1.03	56.9	49.9	78.1	95.5	55.1	<u>99.4</u>	86.1	88.7	96.9	50.2	75.7	3199	8.9
Finetune [30]	10	59.9	60.3	82.1	92.8	55.5	97.5	81.4	87.7	96.6	51.2	76.5	2500	9.4
PB [24]	1.28	57.7	65.3	79.9	97.0	57.5	97.3	79.1	87.6	97.2	47.5	76.6	2838	9.1
DAN [32]	2.17	57.7	64.1	80.1	91.3	56.5	98.5	86.1	<u>89.7</u>	96.8	49.4	77.0	2851	8.3
RA [30]	2	60.3	61.9	81.2	93.9	57.1	99.3	81.7	89.6	96.6	50.1	77.2	3159	7.6
WTPB (full) [26]	1.29	<u>60.8</u>	52.8	<u>82.0</u>	96.2	58.7	99.2	88.2	89.2	96.8	48.6	77.2	3497	5.8
DWSC [10]	~1	64.0	61.1	81.2	97.0	55.5	99.3	85.7	89.1	96.2	49.3	77.8	3507	7.4
SpotTune [11]	11	60.3	63.9	80.5	96.5	57.1	99.5	85.2	88.8	96.7	<u>52.3</u>	78.1	3612	6.4
PA [31]	2	60.3	64.2	81.9	94.7	58.8	<u>99.4</u>	84.7	89.2	96.5	50.9	78.1	3412	6.5
CovNorm [21]	1.25	60.4	69.4	81.3	98.8	59.9	99.1	83.4	87.7	96.6	48.9	78.6	3713	6.6
MAD (full)	2	<u>60.8</u>	<u>66.8</u>	79.2	97.9	56.9	99.2	<u>86.4</u>	90.0	<u>97.0</u>	52.5	<u>78.7</u>	3828	4.1
MAD ($I = 24$)	1.26	<u>60.8</u>	64.9	81.5	<u>98.2</u>	<u>59.1</u>	99.2	85.1	89.5	96.8	<u>52.3</u>	<u>78.7</u>	<u>3816</u>	4
MAD ($I = 36$)	1.39	<u>60.8</u>	65.3	81.7	<u>98.2</u>	<u>59.4</u>	99.2	84.8	89.5	96.8	52.1	78.8	3798	3.9

Table 2. Comparison to the state-of-the-art methods on the Visual Decathlon Challenge. Best results per metric (mean and per dataset classification accuracy, as well as score) are in bold, and the second best are underlined.

port a score 3412, while at the same parameter budget our non-decomposed Modulation Adapters obtain a score 3828. With 3507, DWSC [10] achieves a lower score than the top-performing methods, but uses a smaller parameter budget. Unlike the other methods, however, DWSC does not use the ResNet-26 base network, but an architecture based on depthwise separable convolutions, and is therefore not directly comparable. The number of parameters is still given relative to the ResNet-26 model.

In Table 2, we additionally report the accuracy per dataset as well as the decathlon score. Here, we also include three more baselines. The “Feature” baseline does not adapt the base network, and only learns the linear classification head per dataset. This baseline obtains poor mean accuracy (54.3%) and score (544), due to insufficient adaptation. The “BN Adapt.” baseline only learns dataset-specific BatchNorm layers, but leaves all the other parameters of the base network unchanged. This adds very few parameters, but leads to a remarkable improvement with respect to the feature baseline, by scaling the output of the convolutional layers in a dataset-specific manner. After training, this approach is equivalent to only scaling the output channels of the filters. Our approach is similar, but scales filter weights in both input and output dependent manner, leading to more adaptation and far superior results. The “Finetune” baseline finetunes all the parameters of the base network for each dataset, and does not offer any savings in the parameter budget. It is outperformed by other methods that limit the number of free parameters, allowing for better generalisation to datasets with few samples. We also include SpotTune, which improves over the Finetune baseline in prediction accuracy, but does not offer reduction in the number of parameters.

Our Modulation Adapters (MAD) not only obtain high mean accuracy and score, but also offer the best or comparable accuracy per domain. This is unlike the other methods whose performance is often less consistent across datasets. To quantify this, we rank the methods by accuracy (best first) on each domain, and then average the rank of each method across domains. On the Visual Decathlon Challenge, MAD ($I=36$), MAD ($I=24$) and MAD (full) have average ranks 3.9, 4.0 and 4.1, respectively. These are the best three results on this benchmark. CovNorm, which is our next competitor in terms of the score and accuracy, has the average rank 6.6, which is the seventh-best result. Furthermore, MAD is significantly easier to train than CovNorm. The latter consists of a complex multi-step algorithm: (i) learning Residual Adapters [30], (ii) computing PCA for input and output activations, (iii) learning additional mini-adaptation layers, and (iv) joint finetuning of all components of the final adapters. Training our model is comparable to the first step of CovNorm: we learn a single multiplicative adapter per domain in a single training run.

ImageNet-to-Sketch benchmark. In Table 3, we compare our Modulation Adapters with the state of the art on the ImageNet-to-Sketch benchmark. The full (non-decomposed) model outperforms all the other methods both in terms of the average accuracy and the total score, and is the only one to outperform finetuning the base network (“Finetune”). This result further confirms the effectiveness of the proposed multiplicative adaptation strategy. While our parameter budget is larger than those of our competitors (due to abundance of 1×1 convolutions in the DenseNet-121 base network), it is still below the budget of the fully finetuned models, which allows us to reduce the performance gap between the latter and learning only domain-

Number of images	Params	ImNet 1.3m	CUBS 6k	Stanford Cars 8k	Flowers 2k	WikiArt 42k	Sketch 16k	Mean	Score	Rank
Feature [24]	1	74.4	73.5	56.8	83.4	54.9	53.1	66.0	324	6.8
PB [24]	1.21	74.4	81.4	90.1	95.5	73.9	79.1	82.4	1209	5.8
BA2 [3]	1.17	74.4	82.4	92.9	96.0	71.5	79.9	82.9	1434	4
WTPB (full) [26]	1.21	74.4	81.7	91.6	96.9	75.7	79.8	83.4	1534	3.5
Finetune [24]	6	74.4	81.9	91.4	<u>96.5</u>	<u>76.4</u>	<u>80.5</u>	83.5	1500	3
MAD (full)	4.61	74.4	83.9	<u>91.9</u>	96.9	76.9	81.0	84.2	1668	1.2
MAD ($I=36$)	2.33	74.4	<u>83.1</u>	<u>90.6</u>	96.4	75.3	80.2	83.3	1446	3.7
MAD (full) + WTPB (full)	1.26	74.4	<u>82.7</u>	91.4	96.9	76.2	80.1	<u>83.6</u>	<u>1569</u>	<u>2.7</u>

Table 3. Comparison to the state-of-the-art methods on the ImageNet-to-Sketch benchmark. Best results per metric (mean and per dataset classification accuracy, as well as score) are in bold, and the second best are underlined.

specific classification heads (“Feature”). As for individual tasks, Modulation Adapters demonstrate the best performance on four datasets (CUBS, Flowers, WikiArt and Sketch), and the second-best on Stanford Cars. Results on WikiArt and Sketch are of special interest, since these datasets are quite different from ImageNet which was used for pre-training. While the latter consists of natural images, the former two contain paintings and sketches, which makes the multi-domain task on them more challenging. On ImageNet-to-Sketch, MAD (full) obtains the best average rank 1.2, while next competitors “Finetune” and WTPB obtain 3 and 3.5.

We found that factorisation of our Modulation Adapters is less effective on the ImageNet-to-Sketch dataset: using $I = 36$ intermediate dimensions for decomposition (reducing the parameter budget roughly by a factor two), we obtained an average accuracy of 83.3. The different behaviour of decomposition across the two benchmarks may be related to the relatively large number of parameters in 1×1 convolutions in the DenseNet-121 architecture, whereas in the ResNet-26 architecture most parameters reside in 3×3 convolutions. We also found that the learned adapters for the ImageNet-to-Sketch benchmark do not show the sparse structure that we observe in the adapters learned for the Visual Decathlon Challenge, see Appendix B, making rank restriction less effective. Despite the less consistent performance, MAD ($I=36$) still outperforms other competitors on Sketch and CUBS, improving over full finetuning of the base network on the latter.

As an alternative strategy to reduce the parameter budget, we explored the possibility of combining our method with other approaches. In particular, we trained a hybrid model where WTPB [26] is applied to 1×1 convolutions, while Modulation Adapters are applied to 3×3 convolutions (“MAD (full) + WTPB (full)”). The hybrid model improves over WTPB on three datasets, bringing it closer to individual finetuning on WikiArt and Sketch, and even outperforming it on CUBS. In terms of the average accu-

racy, the total score and the average rank, the hybrid model is slightly better than individual finetuning, second only to our Modulation Adapters. This shows that even more flexible models could be created from existing approaches depending on the desired parameter budget and performance requirements.

5. Conclusion

In this paper, we address the problem of learning across multiple domains while reducing the total parameter budget. In contrast to previous works which rely on adaptation modules that are limited in their capacity, or only adapt a part of the base network parameters, we introduce Modulation Adapters, a novel type of adapters based on flexible weight modulation. We design these adapters to update the weights of the pre-trained convolutional filters by scaling their input and output channels, and provide an efficient parameterisation through decomposition for further parameter reduction. Ablative analyses as well as our model’s excellent performance on the popular Visual Decathlon Challenge and ImageNet-to-Sketch benchmarks validate the benefits of our multiplicative adaptation method. At the same time other means of parameter reduction for the full Modulation Adapter could be considered to cover the setups where the learned adapters are not sparse.

References

- [1] I. Anokhin, K. Demochkin, T. Khakhulin, G. Sterkin, V. Lempitsky, and D. Korzhenkov. Image generators with conditionally-independent pixel synthesis. In *CVPR*, 2021. 3
- [2] N. Atsuhiko and H. Tatsuya. Image generation from small datasets via batch statistics adaptation. In *ICCV*, 2019. 3
- [3] R. Berriel, S. Lathuillere, M. Nabi, T. Klein, T. Oliveira-Santos, N. Sebe, and E. Ricci. Budget-aware adapters for multi-domain learning. In *ICCV*, 2019. 1, 2, 4, 5, 7, 8, 10
- [4] H. Bilen and A. Vedaldi. Weakly supervised deep detection networks. In *CVPR*, 2016. 5

- [5] H. Bilen and A. Vedaldi. Universal representations: The missing link between faces, text, planktons, and cat breeds. *arXiv preprint*, 2017. 7
- [6] F. Chollet. Xception: Deep learning with depthwise separable convolutions. In *CVPR*, 2017. 3
- [7] M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, and A. Vedaldi. Describing textures in the wild. In *CVPR*, 2014. 5
- [8] M. Eitz, J. Hays, and M. Alexa. How do humans sketch objects? *ACM Trans. Graphics*, 31(4):1–10, 2012. 5
- [9] R. Esther, C. Wen-Sheng, K. Abhishek, and H. Jia-Bi. Few-shot adaptation of generative adversarial networks. *arXiv preprint*, 2020. 3
- [10] Y. Guo, Y. Li, L. Wang, and T. Rosing. Depthwise convolution is all you need for learning multiple visual domains. In *AAAI*, 2019. 3, 7
- [11] Y. Guo, H. Shi, A. Kumar, K. Grauman, T. Rosing, and R. Feris. SpotTune: Transfer learning through adaptive fine-tuning. In *ICCV*, 2019. 3, 7, 10
- [12] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 10
- [13] K. He, X. Zhang, S. Ren, and J. Sun. Identity mappings in deep residual networks. In *ECCV*, 2016. 2
- [14] N. Houlsby, A. Giurgiu, S. Jastrzebski, B. Morrone, Q. De Laroussilhe, A. Gesmundo, M. Attariyan, and S. Gelly. Parameter-efficient transfer learning for NLP. In *ICML*, 2019. 3
- [15] T. Karras, S. Laine, M. Aittala, J. Hellsten, J. Lehtinen, and T. Aila. Analyzing and improving the image quality of stylegan. In *CVPR*, 2020. 3
- [16] D. Kingma and J. Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 5
- [17] J. Krause, M. Stark, J. Deng, and L. Fei-Fei. 3d object representations for fine-grained categorization. In *ICCV Workshops*, 2013. 5
- [18] A. Krizhevsky. Learning multiple layers of features from tiny images. Master’s thesis, University of Toronto, 2009. 5
- [19] B. M. Lake, R. Salakhutdinov, and J. B. Tenenbaum. Human-level concept learning through probabilistic program induction. *Science*, 350(6266):1332–1338, 2015. 5
- [20] Y. LeCun, Y. Bengio, and G. Hinton. Deep learning. *Nature*, 52:436–444, 2015. 1
- [21] Y. Li and N. Vasconcelos. Efficient multi-domain learning by covariance normalization. In *CVPR*, 2019. 1, 2, 4, 5, 6, 7, 10
- [22] I. Loshchilov and F. Hutter. Decoupled weight decay regularization. In *ICLR*, 2019. 10
- [23] S. Maji, J. Kannala, E. Rahtu, M. Blaschko, and A. Vedaldi. Fine-grained visual classification of aircraft. *arXiv preprint*, arXiv:1306.5151, 2013. 5
- [24] A. Mallya, D. Davis, and S. Lazebnik. Piggyback: Adapting a single network to multiple tasks by learning to mask weights. In *ECCV*, 2018. 1, 2, 3, 4, 5, 7, 8, 10
- [25] A. Mallya and S. Lazebnik. Packnet: Adding multiple tasks to a single network by iterative pruning. In *CVPR*, 2018. 5
- [26] M. Mancini, E. Ricci, B. Caputo, and S. Rota Bulò. Adding new tasks to a single network with weight transformations using binary masks. In *ECCV Workshops*, 2018. 1, 2, 3, 4, 5, 7, 8, 10
- [27] S. Munder and D. M. Gavrilu. An experimental study on pedestrian classification. *PAMI*, 28(11):1863–1868, 2006. 5
- [28] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng. Reading digits in natural images with unsupervised feature learning. In *NeurIPS Workshop on Deep Learning and Unsupervised Feature Learning*, 2011. 5
- [29] M.-E. Nilsback and A. Zisserman. Automated flower classification over a large number of classes. In *ICVGIP*, 2008. 5
- [30] S.-A. Rebuffi, H. Bilen, and A. Vedaldi. Learning multiple visual domains with residual adapters. In *NeurIPS*, 2017. 1, 2, 3, 4, 5, 7, 10
- [31] S.-A. Rebuffi, H. Bilen, and A. Vedaldi. Efficient parametrization of multi-domain deep neural networks. In *CVPR*, 2018. 2, 4, 5, 6, 7
- [32] A. Rosenfeld and J.K. Tsotsos. Incremental learning through deep adaptation. *PAMI*, 42(3):651–663, 2018. 1, 2, 4, 5, 7, 10
- [33] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, and M. Bernstein. Imagenet large scale visual recognition challenge. *IJCV*, 115(3):211–252, 2015. 1, 5
- [34] B. Saleh and A. Elgammal. Large-scale classification of fine-art paintings: Learning the right metric on the right feature. *International Journal for Digital Art History*, 2016. 5
- [35] A. Senhaji, J. Raitoharju, M. Gabbouj, and A. Iosifidis. Not all domains are equally complex: Adaptive multi-domain learning. In *ICPR*, 2021. 3
- [36] K. Soomro, A. R. Zamir, and M. Shah. A dataset of 101 human action classes from videos in the wild. *Center for Research in Computer Vision*, 2(11), 2012. 5
- [37] J. Stallkamp, M. Schlipsing, J. Salmen, and C. Igel. Man vs. computer: Benchmarking machine learning algorithms for traffic sign recognition. *Neural Networks*, 32:323–332, 2012. 5
- [38] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 3
- [39] P. Welinder, S. Branson, T. Mita, C. Wah, F. Schroff, S. Belongie, and P. Perona. Caltech-ucsd birds 200. Technical report, California Institute of Technology, 2010. 5

In this supplementary material, we describe the experimental setup in Appendix A. In Appendix B, we provide visualization of our Modulation Adapters for some of the convolutions layers in the base network on the two benchmarks.

A. Network architecture and training

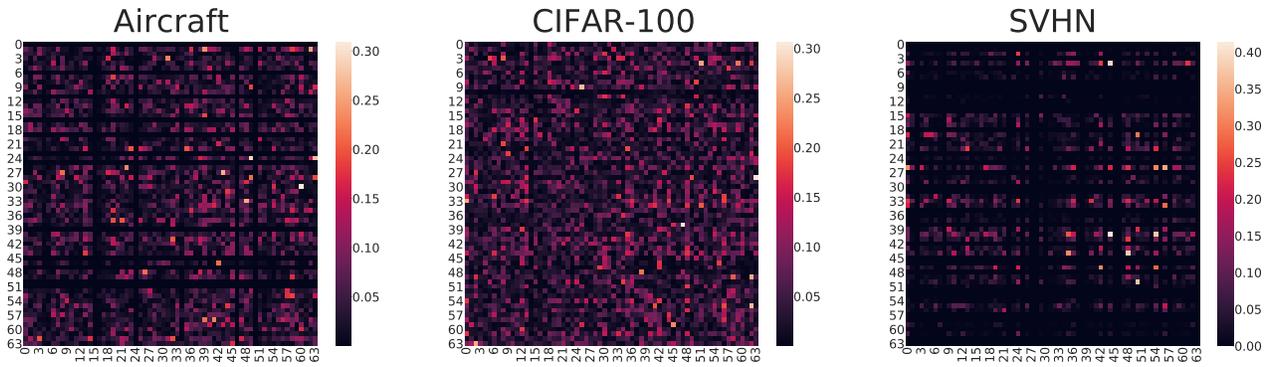
Visual Decathlon Challenge. For the base network, we use ResNet-26 [21, 30, 32] and the weights of its feature extractor pre-trained on ImageNet from [30]. All convolutional layers in the feature extractor consist of 3×3 filters, which are fixed after pre-training and accompanied by domain-specific Modulation Adapters. The latter are initialized with ones and learned for each domain separately. Batch normalization layers, which follow each of these modulated convolutions, are also finetuned for each domain separately, and are initialized with values from the pre-trained base network. Finally, domain-specific fully connected layers used as classification heads are learned from scratch, individually for each task. The model is trained with SGD with momentum 0.9 for 140 epochs. Initial learning rate is set to 0.1, and decreased by a factor ten after epochs 80 and 110. We set the weight decay to 0.0035 for VGG-Flowers, 0.0025 for Aircraft and Describable textures, 0.0015 for UCF101 Dynamic Images, 0.001 for Daimler pedestrian classification and Omniglot, and 0.0005 for CIFAR-100, German traffic signs and Street View House Numbers. All images are resized isotropically to have a shorter size of 72 pixels. We use official data splits from [30], as well as their data processing for all datasets except for Daimler pedestrian classification. For the latter, we use data processing from [11].

ImageNet-to-Sketch. For this benchmark, we use DenseNet-121 [12] as the base network, which is commonly used in previous work [3, 11, 24, 26]. Similarly to the Visual Decathlon Challenge, we apply Modulation Adapters to all convolutions in the base network, including those with 1×1 filters. Despite the absence of parameter reduction in the latter case, we empirically find that modulation of such convolutions leads to better performance compared to finetuning them directly. Domain-specific batch normalization layers and classification heads are initialized and trained the same way as for ResNet-26, while Modulation Adapters α are initialized with 0.15. The model is trained with AdamW [22] for 60 epochs. The initial learning rate is set to 0.001, and decreased by a factor ten after epochs 20 and 40. We set the weight decay to 0.0035 for Flowers, 0.001 for CUBS, Stanford Cars and Sketch, and 0.0005 for WikiArt. We use original data splits and data processing from [24]. Preprocessed images have resolution 224×224 .

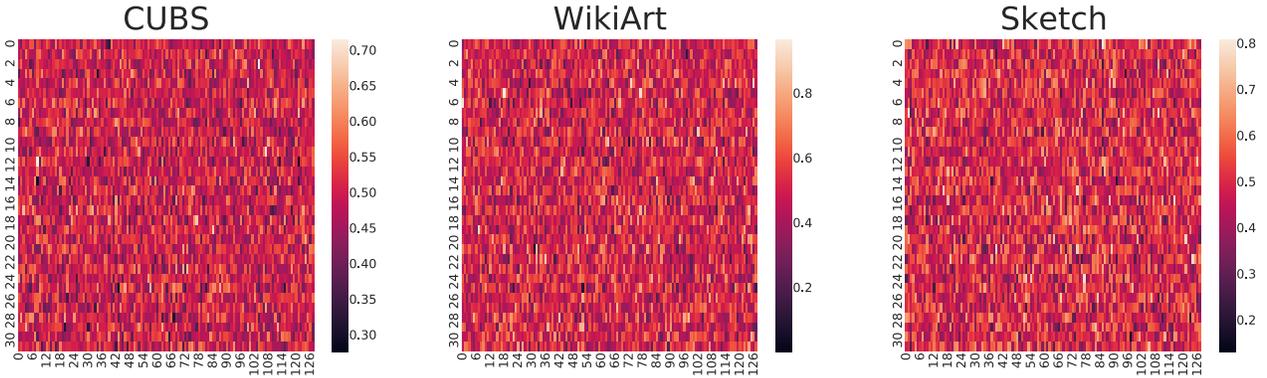
B. Visualization of Modulation Adapters

Figure 4(a) shows the absolute weight values of Modulation Adapters (full) which adapt the same 3×3 convolutional layer at depth five in ResNet-26 on Aircraft, CIFAR-100 and SVHN datasets from the Visual Decathlon Challenge. Figure 4(b) does the same for 3×3 convolutional layer at depth nine in DenseNet-121 on CUBS, WikiArt and Sketch datasets from the ImageNet-to-Sketch benchmark.

Non-trivial adapters are learned, i.e., the matrices do not contain rows or columns with uniform values. In addition to that, we can observe some feature selection happening in ResNet-26: dark spots indicate zeroing of corresponding input feature maps. This is not the case for DenseNet-121 used for ImageNet-to-Sketch benchmark. The low sparsity of Modulation Adapters on the DenseNet-121 might explain why the decomposed version is less effective on this benchmark.



(a) ResNet-26 modulation adapters learned on Visual Decathlon Challenge datasets.



(b) DenseNet-121 modulation adapters learned on ImageNet-to-Sketch benchmark datasets.

Figure 4. Visualization of the absolute weight values of Modulation Adapters. The vertical axis represents the output channels, while the horizontal one represents the input channels.