

RETHINKING HARD-PARAMETER SHARING IN MULTI-DOMAIN LEARNING

Lijun Zhang*, Qizheng Yang†, Xiao Liu‡, Hui Guan◊

University of Massachusetts Amherst, Amherst, MA, USA

* lijunzhang@cs.umass.edu; † yangqizheng@cs.umass.edu;

‡ xiaoliu1990@cs.umass.edu; ◊ huiguan@cs.umass.edu.

ABSTRACT

Hard parameter sharing in multi-domain learning (MDL) allows domains to share some of the model parameters to reduce storage cost while improving prediction accuracy. One common sharing practice is to share the bottom layers of a deep neural network among domains while using separate top layers for each domain. In this work, we revisit this common practice via an empirical study on image classification tasks from a diverse set of visual domains and make two surprising observations. (1) Using separate bottom-layer parameters could achieve significantly better performance than the common practice and this phenomenon holds with different experimental settings. (2) A multi-domain model with a small proportion of domain-specific parameters from bottom layers can achieve competitive performance with independent models trained on each domain separately. Our observations suggest that people adopt the new strategy of using separate bottom-layer parameters as a stronger baseline for model design in MDL.

Index Terms— Multi-domain Learning, Hard-parameter Sharing, Empirical Study

1. INTRODUCTION

Recent years have witnessed the rapid development of Deep Neural Network (DNN) and their superior performance in many areas of artificial intelligence (AI) and vision tasks. AI-powered applications thus increasingly adopt DNNs for solving tasks in single or multiple data streams, leading to more than one DNN model running simultaneously on resource-constrained embedded devices. Although the recent progress on efficient model design and model compression has made it easier to deploy a single model on device, supporting many models on device is still challenging due to the linearly increased bandwidth, energy, and storage costs.

An effective approach to address the problem is multi-domain learning (MDL), where several different domains are learned jointly. Compared to learning each domain separately, MDL allows some parameter sharing across domains to take advantage of the domains’ similarities for improved task performance and reduced storage cost. MDL is closely related to multi-task learning (MTL), where a set of tasks are learned

jointly. Typically, MTL refers to learning different downstream tasks (e.g., depth estimation and semantic segmentation) together while MDL aims to learn on multiple datasets (e.g., data collected from multiple sources with differing statistics) addressing tasks corresponding to each dataset simultaneously. Despite the subtle differences, both MDL and MTL face an open question on how to share parameters across domains or tasks. In this paper, we use the term “domain” and “task” interchangeably as each domain is associated with one task.

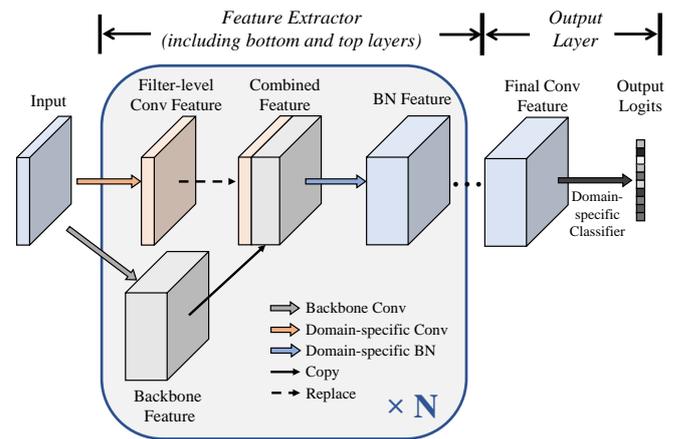


Fig. 1. The forward propagation of a multi-domain model given a domain’s input. The model consists of domain-specific convolution filters, BN layers, and a classifier. When a set of filters in a convolutional layer is domain-specific, the activation maps from these filters (called filter-level conv feature) will replace the corresponding activation maps from the backbone architecture (called backbone feature). The combined feature will be feed into the next layer.

Inspired by the effectiveness of the hard parameter sharing strategy in MTL for improving task performance, researchers believe the same strategy is also a promising solution for MDL. Hard parameter sharing is generally applied by sharing the bottom layers among all tasks, while keeping several top layers and an output layer task-specific [1]. It is commonly used in designing multi-task DNN models in the literature [2, 3] and

gains popularity in MDL [4, 5].

However, it is unclear whether the effectiveness of sharing parameters in bottom layers in MTL can transfer to MDL or not. As shown in Figure 1, a DNN model usually relies on a stack of layers to transform inputs to features and then an output layer to produce predictions based on the features. Researchers [6, 7] believe that the first several layers (bottom layers) serve as low-level feature extractors such as edge detectors and corner detectors, which should be able to be shared across multiple tasks. On the contrary, the top layers are more sensitive to the input data (i.e., the training task), implying that different tasks should possess distinct top layers to generate diverse high-level features. This belief is intuitively reasonable but lacks solid experimental verification for tasks in diverse domains. Whether the conventional hard parameter sharing strategy could lead to the satisfying performance in MDL remains an open question.

In this work, we show that the above sharing strategy is not suitable in MDL through an empirical study on fine-grained image classification tasks and a popular multi-domain learning benchmark Decathlon [8]. Specifically, the common hard parameter sharing strategy, which shares low-level (bottom-layer) parameters while keeping high-level (top-layer) ones domain-specific, is compared to its counterpart, which uses separate bottom-layer parameters for each domain and shares top-layer ones. For the description purpose, we refer to the common hard parameter sharing strategy as *top-specific* while the counterpart as *bottom-specific*. We refer to a model that is trained on diverse domains as a *multi-domain model*. Based on an extensive evaluation on four representative convolutional neural networks (CNN) architectures, we make two major observations.

- Multi-domain models constructed using the bottom-specific strategy could achieve significantly better performance than the ones constructed using the top-specific strategy (i.e., the common practice). Controlled experiments show that this phenomenon can be reproduced with the different number of domains trained together on different backbone architectures using different quantities of domain-specific parameters.
- Multi-domain models with few domain-specific parameters from bottom layers can achieve the same, if not better, performance, as independent models trained separately on each domain in terms of the validation accuracy, which introduces the bottom-specific strategy as a strong baseline for model design in MDL.

These observations advocate for people to rethink the design of hard parameter sharing strategies in MDL. Particularly, because top layers of a modern CNN architecture are usually wider, they tend to have higher redundancy and representation capability that is not fully exploited when trained on a single task. Several prior studies [9, 10] empirically demonstrate that bottom layers are less redundant than top layers in existing architectures. These studies raise a potential expla-

nation to our observation: the bottom-specific strategy could achieve better performance than the top-specific strategy because it makes full use of the capacity in top layers while alleviating task interference by increasing the representation power of bottom layers. We further validate the explanation via a set of controlled experiments based on network pruning. Our experimental evidence and discussions suggest using the bottom-specific strategy as a stronger baseline for model design in MDL.

2. RELATED WORK

Multi-Task Learning. Recent works in multi-task learning (MTL) create multi-task models based on popular DNN architectures called *backbone architectures*. They fall into either hard parameter sharing or soft parameter sharing [1]. Compared with soft parameter sharing where each task still keeps its own model and parameters, hard parameter sharing allows multiple tasks to share some of the model parameters and enjoys the benefits of reduced storage cost and inference latency. This paper thus focuses on hard parameter sharing.

The most widely-used hard parameter strategy is proposed by Caruana [11], which shares the bottom layers of a model across tasks. For instance, Multi-linear Relationship Networks [2] share the first five convolutional layers of AlexNet and use task-specific fully-connected layers for different tasks. Meta Multi-Task Learning [3] could also be materialized as a hard parameter sharing structure. Hard parameter sharing may suffer from task interference because tasks compete for the same set of parameters in the shared bottom layers. However, it remains the most popular paradigm due to its effectiveness in reducing the risk of overfitting and storage cost.

Multi-Domain Learning. Multi-domain learning (MDL) aims at utilizing a single network to perform target tasks in a diverse set of domains. This paper focuses on MDL and studies how to design a compact model that jointly learns representations for all the domains with a small number of domain-specific parameters.

There are two types of approaches to developing multi-domain models. The first type of approaches designs various adapter modules (e.g., Batch Normalization Adapter [12] and Residual Adapter [8]) and plugs them in the backbone architecture. The entire backbone architecture keeps domain-agnostic and is shared across domains while the adapters are domain-specific. A recent study [13] has shown that the choice of adapters and the locations they are plugged in depend on the set of domains. It leverages neural architecture search to figure out what adapter to use and where to add adapters for a given set of domains.

The second type of approaches follows the common practice of hard parameter sharing in MTL. Some researchers [4, 5] proposed to share bottom layers and design sophisticated top layers for each domain. However, it remains an open question whether the effectiveness of sharing parameters in bottom lay-

ers in MTL can transfer to MDL. In other words, it is unclear whether the common practice adopted for tasks from a single domain could lead to similarly satisfied performance for tasks from different domains. In this work, we aim to answer the open question and provide insights on designing a better hard parameter sharing strategy for MDL.

3. METHODOLOGY

Our goal is to study the performance of different hard parameter sharing strategies via controlled experiments, in which the number of tasks, the backbone architectures, the quantity of domain-specific parameters are taken into account. We first describe three design considerations for the studied sharing strategies in this section and then report experimental settings and results in the following two sections.

Domain-specific parameters in the filter granularity.

Our experiments focus on image classification tasks, where each task has its own dataset. CNN models naturally become the backbone architectures due to their superior performance on vision tasks. To create multi-domain models, we determine domain-specific parameters in the granularity of filters instead of layers. Specifically, we will compare the performance of multi-domain models created using the common practice (i.e., the top-specific strategy) and its counterparts (i.e., the bottom-specific strategy) given the same targeted amount of domain-specific weights. The filter-level granularity allows us to precisely control the percentage of domain-specific parameters and ensures that each strategy in comparison can have the same percentage in controlled experiments. Our experiments show that both the amount of domain-specific parameters and where these parameters come from (e.g., bottom layers or top layers) have a significant impact on the performance of a multi-domain model.

Separate classifier for each domain. It is common that different domains expect a different size of outputs or even have diverse prediction goals. In our experiments, the backbone architectures use a separate classifier (i.e., the output layer) for each domain to fit the needs of different output dimensions. Although it is still feasible to allow multiple image classification tasks to share the same classifier, we observe serious performance degradation due to the aggressive sharing. Thus, following the practice in prior work [14], we adopt a separate classifier for each domain.

Separate batch normalization layers for each domain.

We use separate Batch Normalization (BN) layers for each domain in our multi-domain models. It is motivated by a prior study [15], which shows that re-learning a set of scales and biases is sufficient to achieve comparable performance as re-learning the entire set of parameters when a pre-trained model is transferred to another task. The scales and biases correspond to parameters in BN layers in typical CNN architectures. In our experiments, we adopt the idea of making BN parameters domain-specific.

Figure 1 illustrates the forward propagation of a multi-domain model given a specific domain. All domains to be learned together share the same backbone architecture. Each domain has its own BN layers, the output layer that produces logits, and a subset of convolutional filters from the backbone architecture. The remaining convolutional filters are shared by all domains. When a set of filters in a convolutional layer is domain-specific, we use the activation maps produced by these filters to replace the corresponding activation maps from the backbone architecture before the activation maps are fed into the next layer.

4. EXPERIMENTAL SETTINGS

To comprehensively and fairly compare the performance of different hard parameter sharing strategies, we need to consider several potentially influential factors including the backbone architectures, the set of domains and domain-specific parameters, parameter initialization, and hyper-parameters. We next explain each factor in detail.

Backbone architectures. Our experiments use four popular CNNs, MobileNetV2, ResNet50, MNasNet, and SqueezeNet. When creating multi-domain models based on each backbone architecture, we add a separate set of BN layers and a separate output layer for each domain. We also select a subset of filters as domain-specific parameters and share the rest among all domains to be learned jointly.

Datasets. We conduct extensive experiments on five fine-grained image classification tasks. For simplicity, we call it the *FGC benchmark*. We also verify our observations on Decathlon [8], which contains ten well-known datasets from diverse visual domains.

Domain-specific parameters. We experiment with different quantities of domain-specific parameters and sharing strategies. Specifically, we pick different quantities of domain-specific filters such that the number of weights in these filters accounts for 0% to 100% of the total number of convolution parameters in the backbone architecture. We use a step size of 10%. Given the same percentage of domain-specific parameters, we compare three hard parameter sharing strategies:

- *top-specific.* This strategy shares filters in bottom layers and makes filters in top layers domain-specific. It is the common practice used in hard parameter sharing [1, 14].
- *bottom-specific.* It shares filters in top layers while making filters in bottom layers domain-specific. This is a direct counterpart of the top-specific strategy.
- *random.* It randomly selects a subset of filters from all convolutional layers as domain-specific parameters and shares the rest.

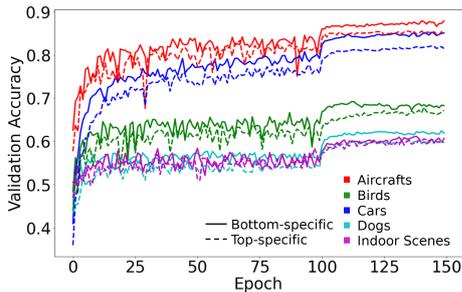
Initialization. Our goal is to eliminate randomness caused by initialization during controlled experiments to ensure the fairness of the comparison and the reproducibility of the experimental results. All multi-domain models are initialized with their corresponding backbone weights pre-trained on Im-

ageNet, while the domain-specific classifiers are pre-trained on each domain separately used to initialize the classifiers in multi-domain models.

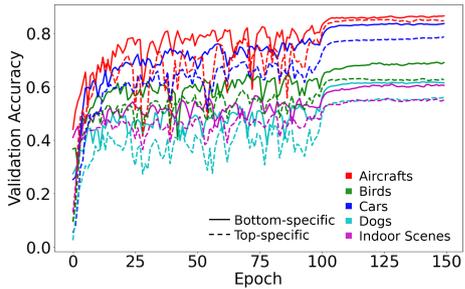
5. RESULTS AND ANALYSIS

5.1. Comparisons between sharing strategies on FGC

Our first surprising discovery is that when constructing multi-domain models, using separate bottom-layer parameters could achieve much better performance than using separate top-layer parameters, which contradicts the common belief in hard parameter sharing in MDL.



(a) MobileNetV2



(b) ResNet50

Fig. 2. The accuracy curves of multi-domain models created with the bottom-specific or top-specific strategy on the five domains in FGC (present by five colors). The percentage of domain-specific parameters for each domain is 20% of the total amount of convolution parameters in the backbone model, (a) MobileNetV2 and (b) ResNet50.

In the experiment, we construct three multi-domain models from each backbone architecture using the three sharing strategies (top-specific, bottom-specific, and random) and train these models on FGC. We strictly control the percentage of domain-specific parameters to be the same for the multi-domain models created using the three strategies.

Table 1 reports the validation accuracy of the 12 multi-domain models (4 backbone architectures \times 3 sharing strategies) on the five domains. The number of domain-specific parameters is 20% of the total amount of convolution parameters in the backbone model. The rows “random” report the mean accuracy of two multi-domain models constructed using

the random strategy. The results indicate that the bottom-specific strategy consistently achieves better performance than the top-specific one and outperforms the random strategy in most cases.

Different backbone architectures. We observed the same phenomenon using different backbone architectures. Figure 2 shows the accuracy curves of multi-domain models created with different strategies. The accuracy curves of the bottom-specific strategy are always above the ones with the top-specific strategy on all five domains with different backbone architectures, indicating the better performance of the bottom-specific strategy. The curves of multi-domain models built on MNasNet and SqueezeNet show a similar pattern.

Different number of domains. The same observation also holds with the different number of domains. Table 1 shows the performance when all five domains are trained together and Table 2 presents the results when fewer domains are used to jointly train multi-domain models built with MobileNetV2. The results suggest that the bottom-specific strategy consistently yields a better prediction performance than the top-specific strategy no matter the number of domains to learn jointly.

Different number of domain-specific parameters. The superiority of the bottom-specific strategy still holds with different quantities of domain-specific parameters. Figure 3 shows the validation accuracy of the multi-domain models whose domain-specific parameters account for 0% to 100% of the total amount of convolution parameters of their backbone architecture MobileNetV2. Note that picking 0% domain-specific parameters results in a multi-domain model with only separate BN layers and domain-specific classifiers, while possessing 100% domain-specific parameters is equivalent to building up a completely independent model for each domain. The results show that the performance of the bottom-specific strategy is always higher than the top-specific one, as indicated by the significant gap between the red and green curves. Besides, randomly selecting domain-specific filters also consistently produces higher validation accuracy than the top-specific strategy but is worse than the bottom-specific one in most cases.

5.2. Comparisons between sharing strategies on Decathlon

The same phenomenon can be observed when constructing multi-domain models on Decathlon with the bottom-specific and the top-specific sharing strategy.

Figure 4 shows the accuracy curves of multi-domain models built on MobileNetV2 with different sharing strategies under the same amount of domain-specific parameters (20% of the total number of convolution parameters in the backbone model). It can be seen that the bottom-specific strategy can produce better accuracy performance than the top-specific one for all or most of the ten domains, which is consistent with the

Table 1. The validation accuracy of the 12 multi-domain models (4 backbone architectures \times 3 sharing strategies) on FGC with the same amount of domain-specific parameters for each domain (20% of the total number of convolution parameters in the backbone model). “independent” reports the results of independent models trained on each domain separately.

Architectures	# Params (M)	Sharing Strategy	Aircraft	Birds	Cars	Dogs	Indoor Scenes
MobileNetV2	10.67	top-specific	0.8536	0.6714	0.8145	0.6006	0.5985
	10.61	random	0.8617	0.6699	0.8259	0.6021	0.6013
	10.63	bottom-specific	0.8782	0.6920	0.8506	0.6186	0.6119
	17.52	independent	0.8749	0.6920	0.8496	0.6202	0.6074
ResNet50	52.84	top-specific	0.8464	0.6243	0.7887	0.5625	0.5530
	52.74	random	0.8650	0.6740	0.8218	0.6211	0.6037
	52.97	bottom-specific	0.8657	0.6925	0.8363	0.6128	0.6067
	127.78	independent	0.8680	0.6799	0.8419	0.6248	0.6037
MNasNet	12.31	top-specific	0.8422	0.6369	0.7747	0.5880	0.6104
	12.34	random	0.8482	0.6634	0.7946	0.6080	0.6082
	12.24	bottom-specific	0.8596	0.6740	0.8081	0.6212	0.6164
	21.91	independent	0.8668	0.6660	0.8174	0.6177	0.6037
SqueezeNet	3.91	top-specific	0.7555	0.5666	0.6423	0.5034	0.5000
	3.90	random	0.7702	0.5817	0.6743	0.5293	0.5119
	3.90	bottom-specific	0.7705	0.5733	0.6764	0.5218	0.5131
	6.24	independent	0.7819	0.6018	0.7016	0.5427	0.5321

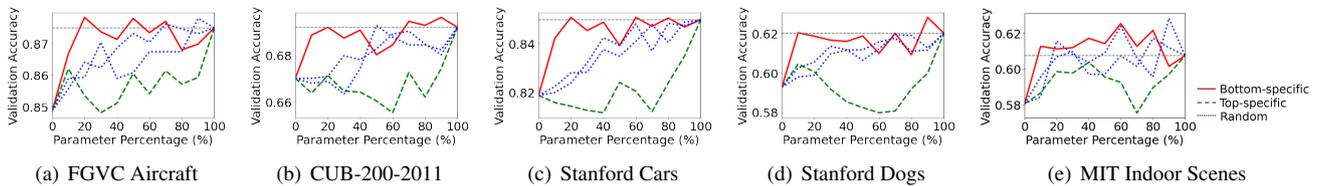


Fig. 3. The validation accuracy trends with an increasing percentage of domain-specific parameters (0% to 100% of the total convolution parameters in MobileNetV2) for the five domains in FGC. In each line chart, the four lines correspond to applying different hard parameter sharing strategies.

Table 2. The validation accuracy of multi-domain models trained on different number of domains. The amount of domain-specific parameters for each domain accounts for 20% of the total number of convolution parameters in the backbone architecture MobileNetV2.

Strategy	Aircraft	Birds	Cars	Dogs
top-specific	0.8719	-	0.8390	-
random	0.8762	-	0.8387	-
bottom-specific	0.8767	-	0.8460	-
top-specific	0.8548	0.6655	0.8336	-
random	0.8687	0.6716	0.8357	-
bottom-specific	0.8677	0.6786	0.8474	-
top-specific	0.8629	0.6723	0.8256	0.5948
random	0.8683	0.6902	0.8318	0.6118
bottom-specific	0.8713	0.6917	0.8423	0.6186

observations on FGC.

5.3. Comparison with independent models

Our second discovery is that multi-domain models with a relatively small proportion of parameters selected from bottom layers for each domain can achieve competitive performance with independent models trained on each domain separately.

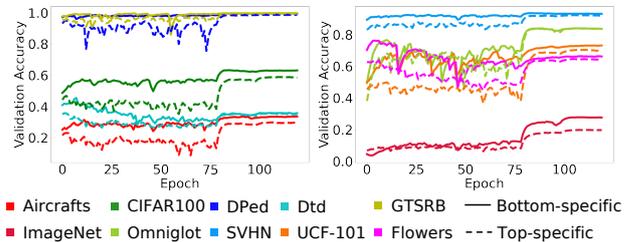


Fig. 4. The accuracy curves of multi-domain models created with the bottom-specific or top-specific strategy on the ten domains in Decathlon (present by ten colors). The percentage of domain-specific parameters for each domain is 20% of the total amount of convolution parameters in the backbone model MobileNetV2.

In Figure 3, the performance of independent models correspond to the point where the percentage of domain-specific parameters is 100%. Overall, multi-domain models constructed with the bottom-specific strategy can achieve competitive validation accuracy as independent models when the percentage of domain-specific parameters is over 20% for all the five domains. This observation is consistent with the well-recognized benefits of MDL in reducing overfitting and improving prediction accuracy.

6. DISCUSSIONS

We summarize the two main observations from our experiments as follows.

- Multi-domain models with domain-specific parameters from bottom layers could achieve better performance than the ones with the same amount of domain-specific parameters from top layers.
- Multi-domain models with a relatively small quantity of domain-specific parameters from bottom layers achieve the same, if not better, performance as their independent counterparts.

Why the bottom-specific sharing strategy outperforms the top-specific one? A potential explanation is that top layers of a modern CNN architecture are usually much wider than bottom layers and thus have a higher representation capability. Prior studies [9, 10] have shown that bottom layers have less redundancy than top layers in existing architectures. When tackling multiple domains together, top layers may have sufficient capacity to learn diverse features while the bottom layers are easily distracted by different domains during training. The bottom-specific strategy could achieve better performance than the top-specific strategy because it makes full use of the capacity in top layers while alleviating task interference by increasing the representation power of bottom layers.

Why owning a small proportion of domain-specific parameters from bottom layers is sufficient to surpass independent models? The reason comes from two aspects. Firstly, DNNs are well-known to be over-parameterized [16, 17]. It is reasonable to assume that a single DNN, especially its top layers, can largely accommodate the representation requirements of multiple domains by taking advantage of the redundant parameters that are not fully exploited on a single domain. A small portion of domain-specific parameters from bottom layers increases the representation power of bottom layers, alleviates domain interference, and thus improves model performance. Secondly, the commonalities among tasks serve as implicit data augmentation [1] and avoid overfitting. This is consistent with the common belief about the benefits of MDL.

7. CONCLUSIONS

In this work, we revisit the common practice in hard parameter sharing for multi-domain learning (MDL) and conduct an empirical study on datasets from different visual domains to compare the performance of different sharing strategies. Experiments show that the common sharing strategy is outperformed by its direct counterpart—that is, selecting domain-specific parameters from bottom layers rather than top layers. The counterpart can also achieve competitive performance compared with independent models. We further provide explanations for the observations.

8. REFERENCES

- [1] Sebastian Ruder, “An overview of multi-task learning in deep neural networks,” *arXiv preprint arXiv:1706.05098*, 2017.
- [2] Mingsheng Long, Zhangjie Cao, Jianmin Wang, and S Yu Philip, “Learning multiple tasks with multilinear relationship networks,” in *NIPS*, 2017, pp. 1594–1603.
- [3] Sebastian Ruder, Joachim Bingel, Isabelle Augenstein, and Anders Søgaard, “Latent multi-task architecture learning,” in *AAAI*, 2019, vol. 33, pp. 4822–4829.
- [4] Syed Shakib Sarwar, Aayush Ankit, and Kaushik Roy, “Incremental learning in deep convolutional neural networks using partial network sharing,” *IEEE Access*, vol. 8, pp. 4615–4628, 2019.
- [5] Shuting He, Hao Luo, Weihua Chen, Miao Zhang, Yuqi Zhang, Fan Wang, Hao Li, and Wei Jiang, “Multi-domain learning and identity mining for vehicle re-identification,” in *CVPR Workshops*, 2020, pp. 582–583.
- [6] Aravindh Mahendran and Andrea Vedaldi, “Visualizing deep convolutional neural networks using natural pre-images,” *IJCV*, vol. 120, no. 3, pp. 233–255, 2016.
- [7] Zhuwei Qin, Fuxun Yu, Chenchen Liu, and Xiang Chen, “How convolutional neural network see the world: A survey of convolutional neural network visualization methods,” *arXiv preprint arXiv:1804.11191*, 2018.
- [8] Sylvestre-Alvise Rebuffi, Hakan Bilen, and Andrea Vedaldi, “Learning multiple visual domains with residual adapters,” *arXiv preprint arXiv:1705.08045*, 2017.
- [9] Kairen Liu, Rana Ali Amjad, and Bernhard C Geiger, “Understanding individual neuron importance using information theory,” *arXiv preprint arXiv:1804.06679*, 2018.
- [10] Shujian Yu, Kristoffer Wickstrøm, Robert Jenssen, and Jose C Principe, “Understanding convolutional neural networks with information theory: An initial exploration,” *IEEE Transactions on Neural Networks and Learning Systems*, pp. 435–442, 2020.
- [11] Rich Caruana, “Multitask learning,” *Machine Learning*, vol. 28, no. 1, pp. 41–75, 1997.
- [12] Hakan Bilen and Andrea Vedaldi, “Universal representations: The missing link between faces, text, planktons, and cat breeds,” *arXiv preprint arXiv:1701.07275*, 2017.
- [13] Hanbin Zhao, Hao Zeng, Xin Qin, Yongjian Fu, Hui Wang, Bourahla Omar, and Xi Li, “What and where: Learn to plug adapters via nas for multidomain learning,” *IEEE Transactions on Neural Networks and Learning Systems*, 2021.
- [14] Trevor Standley, Amir Zamir, Dawn Chen, Leonidas Guibas, Jitendra Malik, and Silvio Savarese, “Which tasks should be learned together in multi-task learning?,” in *ICML*, 2020, pp. 9120–9132.
- [15] Pramod Kaushik Mudrakarta, Mark Sandler, Andrey Zhmoginov, and Andrew Howard, “K for the price of 1: Parameter-efficient multi-task and transfer learning,” *arXiv preprint arXiv:1810.10703*, 2018.
- [16] Jonathan Frankle and Michael Carbin, “The lottery ticket hypothesis: Finding sparse, trainable neural networks,” *arXiv preprint arXiv:1803.03635*, 2018.
- [17] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song, “A convergence theory for deep learning via over-parameterization,” in *ICML*, 2019, pp. 242–252.