

Efficient Multi-Domain Learning by Covariance Normalization

Yunsheng Li Nuno Vasconcelos University of California San Diego La Jolla, CA 92093 yul554@ucsd.edu, nvasconcelos@ucsd.edu

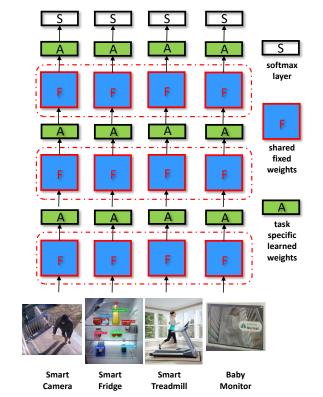
Abstract

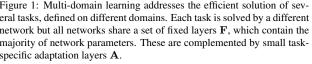
The problem of multi-domain learning of deep networks is considered. An adaptive layer is induced per target domain and a novel procedure, denoted covariance normalization (CovNorm), proposed to reduce its parameters. CovNorm is a data driven method of fairly simple implementation, requiring two principal component analyzes (PCA) and fine-tuning of a mini-adaptation layer. Nevertheless, it is shown, both theoretically and experimentally, to have several advantages over previous approaches, such as batch normalization or geometric matrix approximations. *Furthermore, CovNorm can be deployed both when target* datasets are available sequentially or simultaneously. Experiments show that, in both cases, it has performance comparable to a fully fine-tuned network, using as few as 0.13%of the corresponding parameters per target domain.

1. Introduction

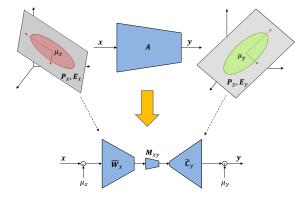
Convolutional nerual networks (CNNs) have enabled transformational advances in classification, object detection and segmentation, among other tasks. However they have non-trivial complexity. State of the art models contain millions of parameters and require implementation in expensive GPUs. This creates problems for applications with network but all networks share a set of fixed layers **F**, which contains computational constraints, such as mobile devices or consumer electronics. Figure 1 illustrates the problem in the specific adaptation layers **A**. context of a smart home equipped with an ecology of de- few models can be cached in the GPU, and moving mod As devices are added to the ecology, the GPU server in the reside on the GPU. A remaining small number of task sp

is to use a different CNN to solve each task. Since only a





vices such as a camera that monitors package delivery and els in and out of cache adds too much overhead to enable theft, a fridge that keeps track of its content, a treadmill that real-time task switching, there is a need for very efficient adjusts fitness routines to the facial expression of the user, parameter sharing across tasks. The individual networks or a baby monitor that keeps track of the state of a baby. should share most of their parameters, which would always house must switch between a larger number of classifica- cific parameters would be switched per task. This problem tion, detection, and segmentation tasks. Similar problems is known as *multi-domain learning* (MDL) and has been will be faced by mobile devices, robots, smart cars, etc. addressed with the architecture of Figure 1 [34, 38]. This Under the current deep learning paradigm, this task consists of set of *fixed* layers (denoted as '**F**') shared by all switching is difficult to perform. The predominant strategy tasks and a set of task specific *adaptation* layers (denoted



1 by three transformations: $\tilde{\mathbf{W}}_{m}$ which implements a

of the network ecology while minimizing the ratio of task and extensive parameter sharing between them. specific (A) to total parameters (both types F and A) per **Domain Adaptation:** In domain adaptation, the source *mation:* a projection into input PCA space, followed by to do the transfer at the image level, e.g. using GANs [11] t

tion (CovNorm). CovNorm is shown to outperform, with for MDL [2]. We show that these mechanisms underpendent of the show that these mechanisms underpendent of the show that the show the show that the show the s both theoretical and experimental arguments, purely geo-

malization [2]. It is also quite simple, requiring two PCAs and the finetuning of a very small mini-adaptation layer pe A layer and task. Experimental results show that it can outperform full network fine-tuning while reducing A layers to as little as 0.53% of the total parameters. When all tasks can be learned together, A layers can be further reduced to 0.51% of the full model size. This is achieved by combining the individual PCAs into a global PCA model, of parameters shared by all tasks, and only fine-tunning miniadaptation layers in a task specific manner.

. Related work

MDL is a transfer learning problem, namely the transfe nto the PCA space of the input x (principal component matrix \mathbf{P}_x and of a model trained on a *source* learning problem to an ecoleigenvalue matrix \mathbf{E}_x), $\tilde{\mathbf{W}}_y$, which reconstructs the PCA space of the ogy of *target* problems. This makes it related to different output y (matrices P_y and E_y), and a mini-adaptation layer M_{xy} . types of transfer learning problems, which differ mostly in as 'A') fine-tunned to each task. If the A layers are much terms of input, or *domain*, and range space, or *task*.

smaller than the **F** layers, many models can be cached simultaneously. Ideally, the F layers should be pre-trained, model trained on a source task to the solution of a target e.g. on ImageNet, and used by all tasks without additional task. The two tasks can be defined on the same or different training, enabling the use of special purpose chips to implement the majority of the computations. While A layers a CNN pre-trained on a large source dataset, such as Ima would still require a processing unit, the small amount of geNet, is usually fine-tunned [21] to a target task. While computation could enable the use of a CPU, making it cost-extremely effective and popular, full network fine-tunning effective to implement each network on the device itself. changes most network parameters, frequently all. MDL ad-In summary, MDL aims to maximize the performance dresses this problem by considering multiple target tasks

network. [34, 38] have shown that the architecture of Fig- and target tasks are the same, and a model trained on a ure 1 can match the performance of fully fine-tuning each source domain is transferred to a target domain. Domain network in the ecology, even when A layers contain as few adaptation can be supervised, in which case labeled data is as 10% of the total parameters. In this work, we show that available for the target domain, or unsupervised, where it A layers can be substantially further shrunk, using a data-is not. Various strategies have been used to address these driven low-rank approximation. As illustrated in Figure 2, problems. Some methods seek the network parameters that this is based on transformations that match the 2nd-order minimize some function of the distance between feature disstatistics of the A layer inputs and outputs. Given prin-tributions in the two domains [24, 4, 43]. Others introduce cipal component analyses (PCAs) of both input and out- an adversarial loss that maximizes the confusion between put, the layer is approximated by a *recoloring transfor*- the two domains [8, 45]. A few methods have also proposed a reconstruction into the output PCA space. By control-map source images into (labeled) target images, then used to ling the intermediate PCA dimensions, the method enables learn a target classifier [3, 41, 14]. All these methods exploit low-dimensional approximations of different input and output dimensions. To correct the mismatch (between PCA and target domains. This is unlike MDL, where source and components) of two PCAs learned independently, a small target tasks are different. Nevertheless, some mechanisms *mini-adaptation* layer is introduced between the two PCA proposed for domain adaptation can be used for MDL. For matrices, and fine-tunned on the target target. example, [5, 28] use a batch normalization layer to match Since the overall transformation generalizes *batch nor*- the statistics of source and target data, in terms of means malization, the method is denoted covariance normaliza- and standard deviation. This is similar to an early proposal

metric methods for matrix approximation, such as the sin- Multitask learning: Multi-task learning [6, 49] adgular value decomposition (SVD) [35], fine-tuning of the dresses the solution of multiple tasks by the same model. It original A layers [34, 38], or adaptation based on batch nor-

a) b) c)

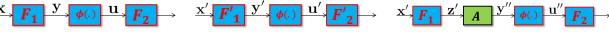
Figure 3: a) original network, b) after fine-tuning, and c) with adaptation layer **A**. In all cases, \mathbf{W}_i is a weight layer and $\phi(.)$ a non-linearity.

sion in object detection [9, 37], joint estimation of surface different domains. More powerful architectures were pro

Most multitask learning approaches emphasize the learn- adapters of multiple tasks into a large matrix, which is a ing domain agnostic lower-level network layers with task was shown to reduce adaptive parameter counts to approx those of all other tasks. Even when multi-task learning with joint optimization. is addressed with multiple tower networks, the emphasis tends to be on inter-tower connections, e.g. through cross- **3. MDL by covariance normalization** stitching [29, 17]. In MDL, such connections are not feasible, because different networks can join the ecology of In this section, we introduce the CovNorm procedure for Figure 1 asynchronously, as devices are turned on and off. MDL with deep networks.

Lifelong learning: Lifelong learning aims to learn mul-**3.1. Multi-domain learning** tiple tasks sequentially with a shared model. This can be upon its use, constraints are needed to force the model to non-linear layer $\phi(.)$ in between. Since the fixed layers are cantly underperform MDL with CovNorm. Methods that ing and converts the network into an ML solution for ity than MDL, since several towers can be needed to solve the weights are changed accordingly, into \mathbf{F}'_1 and \mathbf{F}'_2 . a single task [40], and there is no sharing of fixed layers While very effective, this procedure has two drawbacks, across tasks.

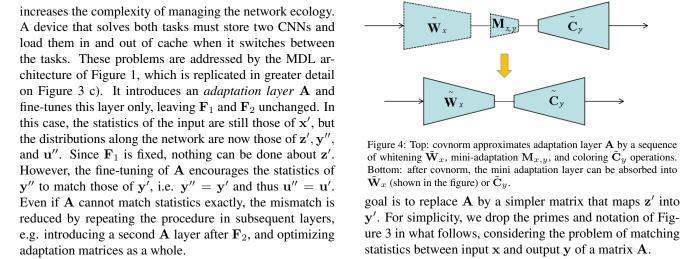
tectures for the adaptation layers of Figure 1. [2] used a BN optimal for S, i.e. the CNN forgets the source task, there is a



examples include classification and bounding box regres-grees of freedom to support transfer of large CNNs across normals and depth [7] or segmentation [29], joint represen-posed by [38], who used a 1×1 convolutional layer and tation in terms of attributes and facial landmarks [50, 33], [34], who proposed a ResNet-style residual layer, known as among others. Multitask learning is sometimes also used to a residual adaptation (RA) module. These methods were solve auxiliary tasks that strengthen performance of a task shown to perform surprisingly well in terms of recogniof interest, e.g. by accounting for context [10], or represent- tion accuracy, equaling or surpassing the performance of ing objects in terms of classes and attributes [15, 29, 30, 25]. full network fine tunning, but can still require a substan-Recently, there have been attempts to learn models that tial number of adaptation parameters, typically 10% of the solve many problems jointly [18, 19, 48]. network size. [35] addressed this problem by combining ing of the interrelationships between tasks. This is frequently accomplished by using a single network, combin- target dataset. Compressing adaptation layers in this way specific network heads and loss functions [50, 7, 10, 15, 37, mately half of [34]. However, all tasks have to be optimized 9], or some more sophisticated forms of network branch-simultaneously. We show that CovNorm enables a further ing [25]. The branching architecture is incompatible with ten-fold reduction in adaptation layer parameters, without MDL, where each task has its own input, different from this limitation, although some additional gains are possible

done by adapting the parameters of a network or adapting Figure 3 a) motivates the use of A layers in MDL. The the network architecture. Since training data is discarded figure depicts two fixed weight layers, F_1 and F_2 , and a remember what was previously learned. Methods that only pre-trained on a *source* dataset S, typically ImageNet, all change parameters either use the model output on previous weights are optimized for the source statistics. For standard tasks [23], previous parameters values [22], or previous net-losses, such as cross entropy, this is a maximum likelihood work activations [44] to regularize the learning of the target (ML) procedure that matches \mathbf{F}_1 and \mathbf{F}_2 to the statistics task. They are very effective at parameter sharing, since of activations \mathbf{x}, \mathbf{y} and \mathbf{u} in \mathcal{S} . However, when the CNN a single model solves all tasks. However, this model is is used on a different *target* domain, the statistics of these not optimal for any specific task, and can perform poorly variables change and F_1 , F_2 are no longer an ML solution. on all tasks, depending on the mismatch between source Hence, the network is sub-optimal and must be finetunned and target domains [36]. We show that they can signifion on a target dataset \mathcal{T} . This is denoted full network finetunadapt the network architecture usually add a tower per new with the outcome of Figure 3 b). In the target domain, the task [40, 1]. These methods have much larger complex-intermediate random variables become \mathbf{x}', \mathbf{y}' , and \mathbf{u}' and

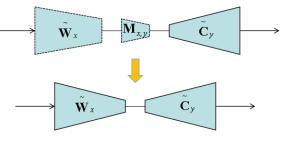
which follow from updating all weights. First, it can be **Multi-domain learning:** This work builds on previous computationally expensive, since modern CNNs have large attempts at MDL, which have investigated different archi-weight matrices. Second, because the weights \mathbf{F}'_i are not layer [16] of parameters tunned per task. While perform-need to store and implement two CNNs to solve both tasks. ing well on simple datasets, this does not have enough de-This is expensive in terms of storage and computation and



3.2. Adaptation layer size

lar to \mathbf{F}_1 . In this case, each domain has as many adaptation parameters as the original network, all networks have twice trix \mathbf{A} is given by the SVD, $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^T$. More prethe size, task switching is complex, and training complexity cisely, the minimum Frobenius norm approximation $\tilde{\mathbf{A}}$ is equivalent to full fine tunning of the original network. On $\arg\min_{\{\mathbf{B}|rank(\mathbf{B})=r\}} ||\mathbf{A} - \mathbf{B}||_F^2$, where $r < rank(\mathbf{A})$, is the other hand, if **A** is much smaller than $\tilde{\mathbf{F}}_1$, MDL is com- $\tilde{\mathbf{A}} = \mathbf{U}\tilde{\mathbf{S}}\mathbf{V}^T$ where $\tilde{\mathbf{S}}$ contains the *r* largest singular values putationally light and task-switching much more efficient. of **A**. This can be written as $\tilde{\mathbf{A}} = \mathbf{C}\mathbf{W}$, where $\mathbf{C} = \mathbf{U}\sqrt{\tilde{\mathbf{S}}}$ In summary, the goal is to introduce an adaptation layer **A** and $\mathbf{W} = \sqrt{\tilde{\mathbf{S}}\mathbf{V}^T}$. If $A \in \mathbb{R}^{d \times d}$, these matrices have a as small as possible, but still powerful enough to match the total of 2rd parameters. An even simpler solution is to destatistics of y' and y''. A simple solution is to make A a fine $\mathbf{C} \in \mathbb{R}^{d \times r}$ and $\mathbf{W} \in \mathbb{R}^{r \times d}$, replace A by their product batch normalization layer [16]. This was proposed in [2] in Figure 3 c), and fine-tune the two matrices instead of but, as discussed below, is not effective. To overcome this **A**. We denote this as the *fine-tunned approximation* (FTA) problem, [38] proposed a linear transformation A and [34] These approaches are limited by their purely geometric naadopted the residual structure of [13], i.e. an adaptation ture. Note that d is determined by the source model (output layer $\mathbf{T} = (\mathbf{I} + \mathbf{A})$. To maximize parameter savings, \mathbf{A} was dimension of \mathbf{F}_1) and fixed. On the other hand, the dimension

parameters, especially in upper network layers. Let \mathbf{F}_1 consimpler, it should be possible to use a smaller r than othvolve a bank of d filters of size $k \times k \times l$ with l feature maps. erwise. There is also no reason to believe that a single n Then, \mathbf{F}_1 has size dk^2l , \mathbf{y} is d dimensional, and \mathbf{A} a $d \times d$ or even a single ratio r/d, is suitable for all network laymatrix. Since in upper network layers k is usually small and ers. While r could be found by cross-validation, this be d > l, A can be only marginally smaller than \mathbf{F}_1 . [35] exploited redundancies across tasks to address this problem, throughout the CNN. We next introduce an alternative, data creating a matrix with the A layer parameters of multiple driven, procedure that bypasses these difficulties. tasks and computing a low-rank approximation of this matrix with an SVD. The compression achieved with this ap- **3.4. Covariance matching** proximation is limited, because the approximation is purely Assume that, as illustrated in Figure 2, x and y are Gaus geometric, not taking into account the statistics of \mathbf{z}' and sian random variables of means μ_x, μ_y and covariance y'. In this work, we propose a more efficient solution, mo- Σ_x, Σ_y , respectively, related by $\mathbf{y} = \mathbf{A}\mathbf{x}$. Let the covaritivated by the interpretation of A as converting the statistics ances have eigendecomposition of \mathbf{z}' into those of \mathbf{y}' . It is assumed that the fine-tuning of \mathbf{A} produces an output variable y'' whose statistics match those of y'. This could leverage adaptation layers in other layers of the network, but that is not important for the discussion where $\mathbf{P}_x, \mathbf{P}_y$ contain eigenvectors as columns and $\mathbf{E}_x, \mathbf{E}_y$



statistics between input \mathbf{x} and output \mathbf{y} of a matrix \mathbf{A} .

3.3. Geometric approximations

Obviously, MDL has limited interest if A has size simi-Geometrically, the closest low rank approximation of a maimplemented with a 1×1 convolutional layer in both cases. sion r should depend on the target dataset \mathcal{T} . Intuitively, This can, however, still require a non-trivial number of \mathcal{T} is much smaller than \mathcal{S} , or if the target task is much

$$\mathbf{\Sigma}_x = \mathbf{P}_x \mathbf{E}_x \mathbf{P}_x^T$$
 $\mathbf{\Sigma}_y = \mathbf{P}_y \mathbf{E}_y \mathbf{P}_y^T$

that follows. The only assumption is that y'' = y'. The are diagonal eigenvalue matrices. We refer to the triplet

 $\mathcal{P}_x = (\mathbf{P}_x, \mathbf{E}_x, \mu_x)$ as the PCA of **x**. Then, it is well known that the statistics of \mathbf{x} and \mathbf{x} are related by that the statistics of \mathbf{x} and \mathbf{y} are related by

and, combining (1) and (2), $\mathbf{P}_{y}\mathbf{E}_{y}\mathbf{P}_{y}^{T} = \mathbf{A}\mathbf{P}_{x}\mathbf{E}_{x}\mathbf{P}_{x}^{T}\mathbf{A}^{T}$. 2 Store the layer input and output PCAs $\mathcal{P}_{x}, \mathcal{P}_{x}$ This holds when $\mathbf{P}_y \sqrt{\mathbf{E}}_y = \mathbf{A} \mathbf{P}_x \sqrt{\mathbf{E}}_x$ or, equivalently, select the k_x, k_y non-zero eigenvalues and

$$\mathbf{A} = \mathbf{P}_y \sqrt{\mathbf{E}_y} \sqrt{\mathbf{E}_x^{-1}} \mathbf{P}_x^T. \qquad (3)$$
$$= \mathbf{C}_y \mathbf{W}_x \qquad (4)$$

where $\mathbf{W}_x = \sqrt{\mathbf{E}_x^{-1}} \mathbf{P}_x^T$ is the "whitening matrix" of x and $\mathbf{C}_y \mathbf{M}_{x,y} \mathbf{W}_x \mu_x + \mu_y$ can be implemented with a $\mathbf{L}_{y} = \mathbf{P}_{y} \sqrt{\mathbf{E}_{y}}$ the "coloring matrix" of y. It follows that vector of biases. (2) holds if $\mathbf{v} = \mathbf{A}\mathbf{x}$ is implemented with a sequence of 4 fine-tune $\mathbf{M}_{T,u}$ with \mathbf{W}_{T} and \mathbf{C}_{u} on \mathcal{T} and absorb two operations. First, **x** is mapped into a variable **w** of zero $\mathbf{M}_{x,y}$ into the larger of $\tilde{\mathbf{W}}_x$ and $\tilde{\mathbf{C}}_y$. mean and identity covariance, by defining

$$\mathbf{w} = \mathbf{W}_x(\mathbf{x} - \mu_x).$$

Second, w is mapped into y with

$$\mathbf{y} = \mathbf{C}_y \mathbf{w} + \mu_y.$$

combination of a whitening of x followed by a colorization and replacing the two matrices by their product reduces the with the statistics of \mathbf{y} .

3.5. Covariance normalization

The interpretation of the adaptation layer as a recoloring operation (whitening + coloring) sheds light on the number f parameters effectively needed for the adaptation, since **3.6. The importance of covariance normalization** the PCAs $\mathcal{P}_x, \mathcal{P}_y$ capture the *effective* dimensions of x and y. Let k_x (k_y) be the number of eigenvalues significantly larger than zero in \mathbf{E}_x (\mathbf{E}_y). Then, the whitening and coloring matrices can be approximated by

 $ilde{\mathbf{W}}_x = \sqrt{ ilde{\mathbf{E}}_x^{-1}} ilde{\mathbf{P}}_x^T \qquad ilde{\mathbf{C}}_y = ilde{\mathbf{P}}_y \sqrt{ ilde{\mathbf{E}}_y}$

where $\tilde{\mathbf{E}}_x \in \mathbb{R}^{k_x \times k_x}$ ($\tilde{\mathbf{E}}_y \in \mathbb{R}^{k_y \times k_y}$) contains the non-zero eigenvalues of Σ_x (Σ_y), and $\tilde{\mathbf{P}}_x \in \mathbb{R}^{d \times k_x}$ ($\tilde{\mathbf{P}}_y \in \mathbb{R}^{d \times k_y}$) which is the batch normalization equation. Hence, Covthe corresponding eigenvectors. Hence, A is well approximated by a pair of matrices $(\tilde{\mathbf{W}}_x, \tilde{\mathbf{C}}_y)$ totaling $d(k_x + k_y)$ ever, important differences. First, there is no batch. The parameters.

On the other hand, the PCAs are only defined up to a permutation, which assigns an ordering to eigenval-the goal is not to facilitate the learning of \mathbf{F}_2 , but produce ues/eigenvectors. When the input and output PCAs are a feature vector \mathbf{y} with statistics matched to \mathbf{F}_2 . This turns computed independently, the principal components may not out to make a significant difference. Since, in regular batch be aligned. This can be fixed by introducing a permutation matrix between \mathbf{C}_{u} and \mathbf{W}_{x} in (4). The assumption that all distributions are Gaussian also only holds approximately in real networks. To account for all this, we augment the recoloring operation with a mini-adaptation layer $M_{x,y}$ of size $k_x \times k_y$. This leads to the covariance normalization (Cov-

$$\tilde{\mathbf{y}} = \tilde{\mathbf{C}}_y \mathbf{M}_{x,y} \tilde{\mathbf{W}}_x (\mathbf{x} - \mu_x) + \mu_y,$$

Data: source S and target T

- $\mu_y = \mathbf{A}\mu_x$ $\Sigma_y = \mathbf{A}\Sigma_x \mathbf{A}^T$ (2) 1 Insert an adaptation layer **A** on a CNN trained on *S* and fine-tune A on \mathcal{T} .
 - corresponding eigenvectors from each PCA, and (3) compute $\tilde{\mathbf{C}}_{u}, \tilde{\mathbf{W}}_{r}$ with (7).
 - (4) 3 add mini-adaptation layer $M_{x,y}$ and replace A b (8). Note that, as usual, the constant

where $\mathbf{M}_{x,y}$ is learned by fine-tuning on the target dataset (5) \mathcal{T} . Beyond improving recognition performance, this has the advantage of further parameters savings. The direct implementation of (8) increases the parameter count to (6) $d(k_x + k_y) + k_x k_y$. However, after fine-tuning, \mathbf{M}_{xy} can be absorbed into one of the two other matrices, as shown in In summary, for Gaussian x, the effect of A is simply the Figure 4. When $k_x > k_y$, $\mathbf{M}_{x,y} \tilde{\mathbf{W}}_x$ has dimension $k_y \times$ total parameter count to $2dk_y$. In this case, we say that $\mathbf{M}_{x,y}$ is absorbed into \mathbf{W}_{x} . Conversely, if $k_{x} < k_{y}$, \mathbf{M}_{y} can be absorbed into \mathbf{C}_{y} . Hence, the total parameter count is $2d \min(k_x, k_y)$. CovNorm is summarized in Algorithm 1.

The benefits of covariance matching can be seen by com parison to previously proposed MDL methods. Assume first, that \mathbf{x} and \mathbf{y} consist of *independent* features. In this case, \mathbf{P}_x , \mathbf{P}_y are identity matrices and (5)-(6) reduce to

$$y_{i} = \sqrt{e_{y,i}} \frac{x_{i} - \mu_{x,i}}{\sqrt{e_{x,i}}} + \mu_{y,i},$$
(9)

normalizing distribution x is now the distribution of the feature responses of layer \mathbf{F}_1 on the target dataset \mathcal{T} . Second, normalization, \mathbf{F}_2 is allowed to change, it can absorb an initial mismatch with the independence assumption. This i not the case for MDL, where \mathbf{F}_2 is *fixed*. Hence, (9) usually

Next, consider the geometric solution. Since CovNorm reduces to the product of two tall matrices, e.g. $\mathbf{K} =$ $\mathbf{C}_{u}\mathbf{M}_{x,u}$ and $\mathbf{L} = \mathbf{W}_{x}$ of size $d \times k_{x}$, it should be possible sible to replace it with the fine-tuned approximation based (8) on two matrices of this size. Here, there are two difficulties.

First, k_{τ} is not known in the absence of the PCA decompositions. Second, in our experience, even when k_x is set to $2d\min(k_x, k_y)$ task specific parameters (per layer) per the value used by PCA, the fine-tuned approximation does dataset. not work. As shown in the experimental section, when the matrices are initialized with Gaussian weights, performance can decrease significantly. This is an interesting observation because A is itself initialized with Gaussian weights. It appears that a good initialization is more critical for the low-rank matrices.

Finally, CovNorm can be compared to the SVD, A = The independent model is needed if, for example, the de $\sqrt{\mathbf{E}_{u}}\sqrt{\mathbf{E}_{x}^{-1}}$ and $\mathbf{U}=\mathbf{P}_{u}$. The problem is that the singular value matrix S conflates the variances of the input and **4. Experiments** output PCAs. The fact that $s_i = e_{y,i}/e_{x,i}$ has two impor-In this section, we present results for both the indepentant consequences. First, it is impossible to recover the dimensions k_x and k_y by inspection of the singular values. **Dataset:** [34] proposed the decathlon dataset for ever Second, the low-rank criteria of selecting the largest sin-uation of MDL. However, this is a collection of relatively gular values is *not* equivalent to CovNorm. For example, the principal components of x with largest eigenvalues $e_{x,i}$ have the smallest singular values s_i . Hence, it is impossible to tell if singular vectors \mathbf{v}_i of small singular values are the most important (PCA components of large variance for \mathbf{x}) or the least important (noise). Conversely, the largest singular values can simply signal the least important input dimensions. CovNorm eliminates this problem by explicitly selecting the important input and output dimensions.

3.7. Joint training

[35] considered a variant of MDL where the different is the same as assuming that a joint dataset $\mathcal{T} = \bigcup_i \mathcal{T}_i$ is available. For CovNorm, the only difference with respect to the single dataset setting is that the PCAs $\mathcal{P}_x, \mathcal{P}_y$ are now those of the joint data \mathcal{T} . These can be derived from the PCAs $\mathcal{P}_{x,i}, \mathcal{P}_{y,i}$ of the individual target datasets \mathcal{T}_i with

$$\mu_{\mathcal{T}} = \frac{1}{N} \sum_{i} N_{i} \mu_{i}$$
$$\boldsymbol{\Sigma}_{\mathcal{T}} = \sum_{i} \frac{N_{i}}{N} (\mathbf{P}_{i} \mathbf{E}_{i} \mathbf{P}_{i}^{T} + \mu_{i} \mu_{i}^{T})) - \mu_{\mathcal{T}} \mu_{\mathcal{T}}^{T} (10)$$

where N_i is the cardinality of \mathcal{T}_i . Hence, CovNorm can be nplemented by finetuning A to each T_i , storing the PCAs $\mathcal{P}_{x,i}, \mathcal{P}_{y,i}$, using (10) to reconstruct the covariance of \mathcal{T} , and computing the global PCA. When tasks are available equentially, this can be done recursively, combining the PCA of all previous data with the PCA of the new data. In summary, CovNorm can be extended to any number of tasks, with constant storage requirements (a single PCA), and no loss of optimality. This makes it possible to define In all experiments, fine-tuning used initial learning rate two CovNorm *modes*.

• *joint*: a global PCA is learned from \mathcal{T} and $\tilde{\mathbf{C}}_{u}, \tilde{\mathbf{W}}_{x}$ shared across tasks. Only a mini-adaptation layer is fine-tuned per \mathcal{T}_i . This requires $\min(k_x, k_y)$ taskspecific parameters (per layer) per dataset. All \mathcal{T}_i must be available simultaneously.

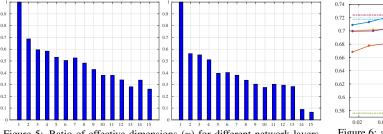
 \mathbf{USV}^T . From (3), this holds whenever $\mathbf{V} = \mathbf{P}_x$, $\mathbf{S} = \mathbf{V}$ vices of Figure 1 are produced by different manufacturers.

dent and joint CovNorm modes.

small datasets. While sufficient to train small networks, we found it hard to use with larger CNNs. Instead, we used collection of seven popular vision datasets. SUN 397 [47] contains 397 classes of scene images and more than a million images. **MITIndoor** [46] is an indoor scene dataset with 67 classes and 80 samples per class. **FGVC-Aircraft Benchmark** [26] is a fine-grained classification dataset of 10,000 images of 100 types of airplanes. Flowers102 [32 is a fine-grained dataset with 102 flower categories and 40 to 258 images per class. CIFAR100 [20] contains 60,000 tiny images, from 100 classes. Caltech256 [12] contains 30,607 images of 256 object categories, with at least 80 tasks of Figure 1 are all optimized simultaneously. This samples per class. **SVHN** [31] is a digit recognition dataset with 10 classes and more than 70,000 samples. In all cases images are resized to 224×224 and the training and testing splits defined by the dataset are used, if available. Other wise, 75% is used for training and 25% for testing.

Implementation: In all experiments, fixed F layers were extracted from a source VGG16 [42] model trained on ImageNet. This has convolution layers of dimensions ranging from 64 to 4096. In a set of preliminary experi-(0) ments, we compared the MDL performance of the architec ture of Figure 1 with these F layers and adaptation layers implemented with 1) a convolutional layer A of kernel size 1×1 [38], 2) the residual adapters $\mathbf{T} = \mathbf{B_2}(\mathbf{I} + \mathbf{AB_1})$ of [34], where \mathbf{B}_1 and \mathbf{B}_2 are batch normalization layer and A as in 1), and 3) the parallel adapters of [35]. Since residual adapters produced the best results, we adopted this structure in all our experiments. However, CovNorm can be used with any of the other structures, or any other matrix A. Note that \mathbf{B}_1 could be absorbed into \mathbf{A} after fine-tuning bu we have not done so, for consistency with [34].

of 0.001, reduced by 10 when the loss stops decreasing. A • *independent:* A layers of network *i* are adapted to ter fine-tuning the residual layer, features were extracted at target dataset \mathcal{T}_i . A PCA is computed for \mathcal{T}_i and the input and output of A and the PCAs $\mathcal{P}_x, \mathcal{P}_y$ computed



Left: MITIndoor. Right: CIFAR100.

only the i^* first eigenvalues/eigenvectors were kept. This dure was used on \mathcal{P}_x or \mathcal{P}_y). Unless otherwise noted, we used t = 0.99, i.e. 99% of the variance was retained.

Benfits of CovNorm: We start with some independent MDL experiments that provide insight on the benefits of 0.3, and smallest for the top network layers.

was first approximated by the SVD and the matrices C, W derperformed RA. This is likely due to overfitting. finetuned on \mathcal{T} , and a mix of PCA and FTA (PCA+FTA), **CovNorm vs SVD:** Figure 7 provides empirical evi-% of parameters. Here, 100% parameters corresponds the PCA only accounts for the subspaces populated by data,

Several observations are possible. First, all geomet- a layer-dependent number of singular values.

0.65

Figure 5: Ratio of effective dimensions (η) for different network layers. Figure 6: accuracy vs. % of parameters used for adaptation. Left: MITIndoor. Right: CIFAR100.

and used in Algorithm 1. Principal components were se- (SVD+FTA) is as large as 2%. This is partly due to the use lected by the explained variance criterion. Once the eigen-of a constant low rank r throughout the network. This canvalues e_i were computed and sorted by decreasing magni-not match the effective, data-dependent, dimensions, which tude, i.e. $e_1 \ge e_2 \ge \ldots \ge e_d$, the variance explained by vary across layers (see Figure 5). CovNorm eliminates this the first *i* eigenvalues is $r_i = \frac{\sum_{k=1}^{i} e_i}{\sum_{k=1}^{i} d_{k-1}}$. Given a threshold *t*, problem. We experimented with heuristics for choosing the smallest index i^* such that $r_{i^*} > t$ was determined, and variable ranks but, as discussed below (Figure 7), could not achieve good performance. Among the geometric apset the dimensions k_x, k_y (depending on whether the procemance drops in most of datasets. It is interesting that, while A is fine-tuned with random initialization, the process is not effective for the low-rank matrices of FTA. In several datasets, FTA could not match SVD+FTA.

CovNorm over previous MDL procedures. While we only Even more surprising were the weaker results obtained report results for MITIndoor and CIFAR100, they are typi- when the random initialization was replaced by the two cal of all target datasets. Figure 5 shows the ratio $\eta = k_u/k_x$ PCAs (PCA+FTA). Note the large difference between of effective output to input dimensions, as a function of PCA+FTA and CovNorm (up to 4%), which differ by the adaptation layer. It shows that the input of A typically con-mini-adaptation layer $\mathbf{M}_{x,y}$. This is explained by the aligntains more information than the output. Note that η is rarely ment problem of Section 3.5. Interestingly, while minione, is almost always less than 0.6, frequently smaller than adaptation layers are critical to overcome this problem, they are as easy to fine-tune as A. In fact, the addition of these We next compared CovNorm to batch normalization layers (CovNorm) often outperformed the full matrix A (BN) [2], and geometric approximations based on the finetunned approximation (FTA) of Section 3.3. We also tested parameters, CovNorm matched the performance of RA, Fia mix of the geometric approaches (SVD+FTA), where **A** nally, as previously reported by [34], FNFT frequently un-

where the mini-adaptation layer $\mathbf{M}_{r,u}$ of CovNorm was removed and $\tilde{\mathbf{C}}_y, \tilde{\mathbf{W}}_x$ fine-tuned on \mathcal{T} , to minimize the PCA produced by CovNorm and the SVD. The figure shows a alignment problem. All geometric approximations were implemented with low-rank parameter values $r = d/2^i$, where put and output distributions of an adaptation layer A and d is the dimension of x or y and $i \in \{2, ..., 6\}$. For the corresponding plot for its singular values. Note how the CovNorm, the explained variance threshold was varied in PCA energy is packed into a much smaller number of coeffi-[0.8, 0.995]. Figure 6 shows recognition accuracies vs. the cients than the singular value energy. This happens because adaptation layers of [34]: a network with residual adapters restricting the low-rank approximation to these subspaces. whose matrix \mathbf{A} is fine-tunned on \mathcal{T} . This is denoted RA Conversely, the geometric approximation must approximate and shown as an upper-bound. A second upper-bound is the matrix behavior even outside of these subspaces. Note shown for full network fine tuning (FNFT). This requires that the SVD is not only less efficient in identifying the im- $10 \times$ more parameters than RA. BN, which requires close to portant dimensions, but also makes it difficult to determine zero parameters, is shown as a lower bound. how many singular values to keep. This prevents the use of

ric approximations underperform CovNorm. For comparable sizes, the accuracy drop of the best geometric method the recognition accuracy and % of adaptation layer param-

Table 1: Classification accuracy and % of adaptation parameters (with respect to VGG size) per target

	FGVC	MITIndoor	Flowers	Caltech256	SVHN	SUN39
FNFT	85.73%	71.77%	95.67%	83.73%	96.41%	57.29%
				100%	1	
			Ind	ependent learnii	ng	
BN [2]	43.6%	57.6%	83.07%	73.66%	91.1%	47.04%
		1		0%	1	1
LwF[23]	66.25%	73.43%	89.12%	80.02%	44.13%	52.85%
				0%		
RA [34]	88.92%	72.4%	96.43%	84.17%	96.13%	57.38%
				10%		
SVD+FTA	89.07%	71.66%	95.67%	84.46%	96.04%	57.12%
				5%		
FTA	87.31%	70.26%	95.43%	83.82%	95.96%	56.43%
				5%		
CovNorm	88.98%	72.51%	96.76%	84.75%	96.23%	57.97%
	0.34%	0.62%	0.35%	0.46%	0.13%	0.71%
				Joint learning		
SVD [35]	88.98%	71.7%	96.37%	83.63%	96%	56.58%
	5%					
CovNorm	88.99%	73.0%	96.69%	84.77%	96.22%	58.2
				0.51%		
		NI-4 Aim	C100	DD-1	DTD	CTCD
	11	nNet Airc	C100	DPed	DTD	GTSR

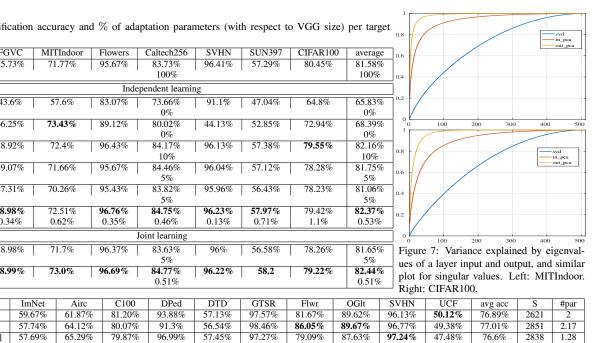
vNorm 60.37% 69.37% 81.34% 98.75% 59.95% 99.14% 83.44% 87.69% 96.55% 48.92% 78.55% 3713 1

Table 2: Visual Decathlon results

methods. All abbreviations are as above. Beyond MDL, Fourth, for joint training, CovNorm is substantially supewe compare to learning without forgetting (LwF) [23] a rior to the SVD [35], with higher recognition rates in all lifelong method to learn a model that shares all parame-datasets, gains of up to 1.62% (SUN397), and close to $10\times$ ters among datasets. The table is split into independent and less parameters. Finally, comparing independent and joint joint MDL. For joint learning, CovNorm is implemented CovNorm, the latter has slightly higher recognition for a

Several observations can be made. First, CovNorm are roughly equivalent. adapts the number of parameters to the task, according to **Results on Visual Decathlon** Table 2 presents result. its complexity and how different it is from the source (Ima-on the Decathlon challenge [34], composed of ten differgeNet). For the simplest datasets, such as the 10-digit class ent datasets of small images (72×72) . Models are trained SVHN, adaptation can require as few as 0.13% task-specific with a combination of training and validation set and results parameters. Datasets that are more diverse but ImageNet-obtained online. For fair comparison, we use the learning like, such as Caltech256, require around 0.46% parameters. protocol of [34]. CovNorm achieves state of the art perfor-Finally, larger adaptation layers are required by datasets that mance in terms of classification accuracy, parameter size, are either complex or quite different from ImageNet, e.g. and decathlon score S. scene (MITIndoor, SUN397) recognition tasks. Even here, **5.** Conclusion adaptation requires less than 1% parameters. On average, CovNorm requires 0.53% additional parameters per dataset. CovNorm is an MDL technique of very simple imple-

residual adapters significantly outperform BN and LwF. As ically reduces the number of adaptation parameters without shown by [34], RA outperforms FNFT. BN is uniformly loss of recognition performance. It was used to show that weak, LwF performs very well on MITIndoor and Cal-large CNNs can be "recycled" across problems as diverse tech256, but poorly on most other datasets. Third, Cov- as digit, object, scene, or fine-grained classes, with no loss, Norm outperforms even RA, achieving higher recognition by simply tuning 0.5% of their parameters. accuracy with $20 \times$ less parameters. It also outperforms **6.** Acknowledgment SVD+FTA and FTA by $\approx 0.6\%$ and $\approx 1.3\%$, respectively, while reducing parameter sizes by a factor of ≈ 10 . On a This work was partially funded by NSF awards IISper-dataset basis, CovNorm outperforms RA on all datasets 1546305 and IIS-1637941, a GRO grant from Samsung, and other than CIFAR100, and SVD+FTA and FTA on all of NVIDIA GPU donations.



- eters vs. VGG model size (100% parameters), for various them. In all datasets, the parameter savings are significant. with (10) and compared to the SVD approach of [35]. slightly higher parameter count. Hence, the two approaches

Second, for independent learning, all methods based on mentation. When compared to previous methods, it drama

References

- [1] R. Aljundi, P. Chakravarty, and T. Tuytelaars. Expert gate: Lifelong learning with a network of experts. In CVPR, pages 7120-7129, 2017.
- [2] H. Bilen and A. Vedaldi. Universal representations: The missing link between faces, text, planktons, and cat breeds. arXiv preprint arXiv:1701.07275, 2017.
- ishnan. Unsupervised pixel-level domain adaptation with generative adversarial networks. In *The IEEE Conference* [21] Y. LeCun, Y. Bengio, and G. Hinton. Deep learning. *nature*, on Computer Vision and Pattern Recognition (CVPR), volume 1, page 7, 2017.
-] K. Bousmalis, G. Trigeorgis, N. Silberman, D. Krishnan, and D. Erhan. Domain separation networks. In Advances in Neural Information Processing Systems, pages 343–351, 2016.
- pages 5077-5085 2017
- 95–133. Springer, 1998.
- Conference on Computer Vision, pages 2650–2658, 2015.
- [8] Y. Ganin and V. Lempitsky. Unsupervised domain adaptation volume 1, page 6, 2017. Learning, 2014.
- [9] R. Girshick. Fast r-cnn. arXiv preprint arXiv:1504.08083,
- [10] G. Gkioxari, R. Girshick, and J. Malik. Contextual action recognition with r* cnn. In *Proceedings of the IEEE inter*national conference on computer vision, pages 1080–1088,
- [11] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In Advances in neural information processing systems, pages 2672–2680, 2014.
- [12] G. Griffin, A. Holub, and P. Perona. Caltech-256 object category dataset. 2007.
- [13] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In Proceedings of the IEEE conrence on computer vision and pattern recognition, pages 770–778, 2016.
- [14] J. Hoffman, E. Tzeng, T. Park, J.-Y. Zhu, P. Isola, K. Saenko, A. A. Efros, and T. Darrell. Cycada: Cycle-consistent adver-
- [15] J. Huang, R. S. Feris, Q. Chen, and S. Yan. Cross-domain image retrieval with a dual attribute-aware ranking network. *puter vision*, pages 1062–1070, 2015.
- deep network training by reducing internal covariate shift. arXiv preprint arXiv:1502.03167, 2015.

- [18] A. Kendall, Y. Gal, and R. Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semar tics. arXiv preprint arXiv:1705.07115, 3, 2017.
- [19] I. Kokkinos. Ubernet: Training a universal convolutional neural network for low-, mid-, and high-level vision using diverse datasets and limited memory. In CVPR, volume 2,
- [3] K. Bousmalis, N. Silberman, D. Dohan, D. Erhan, and D. Krfeatures from tiny images. 2009.

 - [22] S.-W. Lee, J.-H. Kim, J. Jun, J.-W. Ha, and B.-T. Zhang. Overcoming catastrophic forgetting by incremental moment matching. In Advances in Neural Information Processing Systems, pages 4655-4665, 2017.
- [5] F. M. Carlucci, L. Porzi, B. Caputo, E. Ricci, and S. R. Bulò. [23] Z. Li and D. Hoiem. Learning without forgetting. *IEEE* Autodial: Automatic domain alignment layers. In ICCV, Transactions on Pattern Analysis and Machine Intelligence,
- [6] R. Caruana. Multitask learning. In Learning to learn, pages [24] M. Long, Y. Cao, J. Wang, and M. I. Jordan. Learnin transferable features with deep adaptation networks. Inter-[7] D. Eigen and R. Fergus. Predicting depth, surface normals *national Conference in Machine Learning*, 2015.
- and semantic labels with a common multi-scale convolu-[25] Y. Lu, A. Kumar, S. Zhai, Y. Cheng, T. Javidi, and R. S. tional architecture. In *Proceedings of the IEEE International* Feris. Fully-adaptive feature sharing in multi-task network with applications in person attribute classification. In CVF
- by backpropagation. International Conference in Machine [26] S. Maji, E. Rahtu, J. Kannala, M. Blaschko, and A. Vedaldi. Fine-grained visual classification of aircraft. arXiv preprint arXiv:1306.5151, 2013.
 - [27] A. Mallya, D. Davis, and S. Lazebnik. Piggyback: Adapting a single network to multiple tasks by learning to mask weights. In Proceedings of the European Conference of Computer Vision (ECCV), pages 67–82, 2018.
 - [28] M. Mancini, L. Porzi, S. R. Bulò, B. Caputo, and E. Ricci. Boosting domain adaptation by discovering latent domains. arXiv preprint arXiv:1805.01386. 2018.
 - [29] I. Misra, A. Shrivastava, A. Gupta, and M. Hebert. Cross stitch Networks for Multi-task Learning. In CVPR, 2016.
 - [30] P. Morgado and N. Vasconcelos. Semantically consistent regularization for zero-shot recognition. In CVPR, volume 9,
 - [31] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng. Reading digits in natural images with unsupervised feature learning. In NIPS workshop on deep learning and u pervised feature learning, volume 2011, page 5, 2011
- sarial domain adaptation. *arXiv preprint arXiv:1711.03213*, [32] M.-E. Nilsback and A. Zisserman. Automated flower classification over a large number of classes. In Computer Vision, Graphics & Image Processing, 2008. ICVGIP'08. Sixth Indian Conference on, pages 722–729. IEEE, 2008.
- In Proceedings of the IEEE international conference on com-[33] R. Ranjan, V. M. Patel, and R. Chellappa. Hyperface: A deep multi-task learning framework for face detection, landmark [16] S. Joffe and C. Szegedy. Batch normalization: Accelerating localization, pose estimation, and gender recognition. *IEEE* Transactions on Pattern Analysis and Machine Intelligence,
- [17] B. Jou and S.-F. Chang. Deep cross residual learning for multitask visual recognition. In Proceedings of the 2016 ACM on visual domains with residual adapters. In Advances in Neuro Multimedia Conference, pages 998–1007. ACM, 2016. Information Processing Systems, pages 506–516, 2017.

- [35] S.-A. Rebuffi, H. Bilen, and A. Vedaldi. Efficient parametrization of multi-domain deep neural networks. arXiv preprint arXiv:1803.10082, 2018.
- [36] S.-A. Rebuffi, A. Kolesnikov, G. Sperl, and C. H. Lampert. icarl: Incremental classifier and representation learning. In *Proc. CVPR.* 2017.
- [37] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: towards real-time object detection with region proposal ne works. IEEE transactions on pattern analysis and machine intelligence, 39(6):1137–1149, 2017.
- [38] A. Rosenfeld and J. K. Tsotsos. Incremental learning throu deep adaptation. arXiv preprint arXiv:1705.04228, 2017
- [39] A. Rosenfeld and J. K. Tsotsos. Incremental learning throug deep adaptation. IEEE transactions on pattern analysis and machine intelligence, 2018.
- [40] A. A. Rusu, N. C. Rabinowitz, G. Desjardins, H. Soyer, J. Kirkpatrick, K. Kavukcuoglu, R. Pascanu, and R. Hadsell. Progressive neural networks. arXiv preprint arXiv:1606.04671, 2016.
- [41] A. Shrivastava, T. Pfister, O. Tuzel, J. Susskind, W. Wang, and R. Webb. Learning from simulated and unsupervised images through adversarial training. In CVPR, volume 2,
- [42] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint* arXiv:1409.1556, 2014.
- [43] B. Sun and K. Saenko. Deep coral: Correlation alignment for deep domain adaptation. In European Conference on Computer Vision, pages 443–450. Springer, 2016.
- [44] A. R. Triki, R. Aljundi, M. B. Blaschko, and T. Tuytelaars. Encoder based lifelong learning. IEEE Conference Computer Vision and Pattern Recognition, 2017.
- [45] E. Tzeng, J. Hoffman, K. Saenko, and T. Darrell. Adversarial discriminative domain adaptation. In Computer Vision and Pattern Recognition (CVPR), volume 1, page 4, 2017.
- [46] M. Valenti, B. Bethke, D. Dale, A. Frank, J. McGrew, S. Ahrens, J. P. How, and J. Vian. The mit indoor multivehicle flight testbed. In Robotics and Automation, 2 IEEE International Conference on, pages 2758–2759. IEEE
- [47] J. Xiao, J. Hays, K. A. Ehinger, A. Oliva, and A. Torralba. Sun database: Large-scale scene recognition from abbey to zoo. In Computer vision and pattern recognition 2010 IEEE conference on, pages 3485–3492. IEEE, 2010
- [48] A. R. Zamir, A. Sax, W. Shen, L. Guibas, J. Malik, and S. Savarese. Taskonomy: Disentangling task transfer learning. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 3712-3722, 2018
- [49] Y. Zhang and Q. Yang. A survey on multi-task learning. arXiv preprint arXiv:1707.08114, 2017.
- [50] Z. Zhang, P. Luo, C. C. Loy, and X. Tang. Facial landmark detection by deep multi-task learning. In European Con ence on Computer Vision, pages 94-108. Springer, 2014