¹ UC San Diego, ² Microsoft

Abstract

on learning a domain-agnostic model, of which the parameters are **static**. However, such a static model is difficult to handle conflicts across multiple domains, and suffers from a performance degradation in both source domains and target domain. In this paper, we present dynamic transfer to address domain conflicts, where the model parameters are adapted to samples. The key insight is that adapting model across domains is achieved via adapting model across samples. Thus, it breaks down source domain barriers and turns multi-source domains into a single-source domain. This also simplifies the alignment between source and target domains, as it only requires the target domain to be aligned with any part of the union of source domains. Furthermore, we find dynamic transfer can be simply modeled by trix. Experimental results show that, without using domain 'Static Transfer' implements domain adaptation with a labels, our dynamic transfer outperforms the state-of-the-static model f_{θ_c} , which has fixed parameters θ_c to average art method by more than 3% on the large multi-source do-domain conflict. (b) 'Dynamic Transfer' $(f_{\theta(x)})$ adapts the main adaptation datasets – DomainNet. Source code is at model parameters $\theta(x)$ according to samples, which gen-

Multi-source domain adaptation addresses the adaptation cult, since different domains can give rise to very different from multiple source domains to a target domain. It is chalimage distributions. When forcing a model to be domain lenging because a clear domain discrepancy exists not only agnostic, it essentially averages the domain conflict. Thus between source and target domains, but also among multiple source domains (see exemplar images in Figure 2). This dated by our preliminary study. As shown in Figure 2, comsuggests that successful adaptation requires significant *elas*ticity of the model to adapt. A nature way to achieve this model consistently degrades in each source domain. elasticity is to make model dynamic i.e. the mapping imIn this paper, we propose dynamic transfer to address plemented by the model should vary with the input sample. this issue. As shown in Figure 1(b), it contains a param-

Yunsheng Li¹, Lu Yuan², Yinpeng Chen², Pei Wang¹, Nuno Vasconcelos¹

Recent works of multi-source domain adaptation focus $x \qquad f_{\theta_c} \qquad y = f_{\theta_c}(x) \qquad x \qquad f_{\theta(x)} \qquad y = f_{\theta(x)}(x)$ → Model parameter flow → Data flow (a) Static Transfer (b) Dynamic Transfer aggregating residual matrices and a static convolution ma-

This hypothesis has not been explored by existing work, eter predictor that changes the model parameters on a per e.g. [20, 26], which instead aims to learn a domain agnostic sample basis, i.e. implements mapping $f_{\theta(x)}$. It has the model f_{θ_c} , of static parameters θ_c , that works well for all advantage of not requiring the definition of domains or the source $\{S_1, S_2, ..., S_N\}$ and target \mathcal{T} domains. We refer to collection of domain labels. In fact, it unifies the probthis approach as *static transfer*. As illustrated in Figure 1 lems of single-source and multi-source domain adaptation. (a), the model implements a fixed mapping across all do-

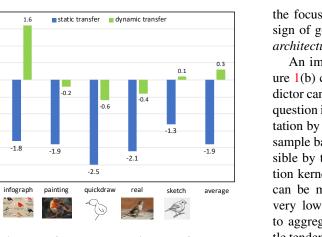


Figure 2: Static Transfer vs. Dynamic Transfer on the performance degradation of source domains compared enhance the domain adaptation performance. Experimento the oracle results. Both transfer models are tested across tal results show that the proposed dynamic residual transsource domains. A clear performance degradation (1.9% fer (DRT) can model domain variation in source domains in average) exists when using static transfer, indicating the (see Figure 2) and outperform its static counterpart (MCD conflicts across domains. The degradation is significantly [23] method) by a large margin (11.2% on DomainNet). reduced when using dynamic transfer, as it handles the domain variations well. (Best view in color)

source domain adaptation into a single-source domain problem. The only difference is the complexity of this domain.

2. Related Work

mains is achieved statistically by adapting model per sam-model from a source to a target domain. A common method ple, since each domain is viewed as a distribution of imis to minimize the distance between the two domains. While age samples. The dynamic transfer learns how to adapt the some methods [3, 24] minimize distance functions defined model's parameters and fit to the union of source domains. in terms of first and second order data statistics, others learn Thus the alignment between source domains and target domain is significantly simplified, as it is no longer necessary ing [25, 23, 9, 14]. Although these methods are effective to pull all source domains together with the target domain. for single-source domain distributions and relatively sim-In this case, as long as the target domain is aligned to any ple datasets (such as VisDA [21] or Office-31 [22]), they

erates a different model per sample and turns multi-source domain adaptation into single-source domain adaptation. mains. However, learning a domain agnostic model is diffi-

the literature, most works assume the static network of Fig- a variety of styles. [31] pioneered this problem by adapure 1(a) and focus on loss functions. The goal is to detively picking the best among a set of hypothesis learned fine losses that somehow "pull all the domains" together for different source domains. [1] derived an upper bound into a shared latent representation. The problem is that on the classification error achievable in the target domain, the domains are usually very different at the network in-based on the $\mathcal{H}\Delta\mathcal{H}$ divergence. Several methods have been put. Hence, the force introduced by the loss at the output, to proposed after the introduction of deep learning. Some of bring them together, is counter-balanced by an input force to these align domains pair-wise. [29] uses a discriminator to keep them apart. This usually leads to a difficult optimizaalign each source domain with the target domain, while [20] tion and compromises adaptation performance. The introduction of a dynamic network, as in Figure 1(b), enables a mains. These methods learn one classifier per domain and more elastic mapping. In this case, it is not necessary to pull use their weighted combination to predict the class of target

the focus of the domain adaptation problem from the design of good *loss functions* to the design of good *network architectures* for dynamic transfer.

dictor cannot generate all parameters for a large model. The question is whether it is possible to perform the model adaptation by only modifying a small subset of parameters on a sample basis. In this work, we show that this is indeed possible by the addition of dynamic residuals to the convolution kernels of a static network. Since the residual blocks can be much smaller than the static ones, this has both very low additional computational cost (less than 0.1%) to aggregate dynamic residuals with static kernel and little tendency to overfit. However, it is shown to significantly tion methods [26], it achieves a sizeable gain (3.9%) with a

The key insight is that adapting model according to dopart of the source domains, the model can be easily adapted are not competitive for the multi-source domain adaptation

When compared to the domain adaption literature, dy
Multi-Source Domain Adaptation considers the domain namic transfer introduces a significant paradigm shift. In adaptation problem when the source contains domains with

An immediate difficulty is that the architecture of Figure 1(b) can be very hard to train, since the parameter preenhance multi-source domain adaptation.

much simpler loss function and training algorithm.

problem, due to a more complex data distribution.

tween domains with a knowledge graph. Target sample pre-

dictions are based on both their features and relationship The model f_{θ} is denoted static or dynamic depending on to different domains. [13] proposes a meta-learning techwhether the model parameters θ vary with samples x. Static nique to search the best initial conditions for multi-source models have constant parameters $\theta = \theta_c$, while dynamic domain adaptation. [32] uses an auxiliary network to premodels have parameters $\theta = \theta(x)$ that depend on x. In dict the transferability of each source sample and use it as the case of deep networks, this implies that layer transfer a weight to learn a domain discriminator. All these works functions depend on the input x. Figure 1 illustrates the use a static transfer model. In this paper, we propose that static transfer and dynamic transfer model built for multithe model should instead be dynamic, i.e. a function that source domain adaptation

changes with samples, and show that this can significantly **Static Transfer.** Static transfer, shown on Figure 1(a), consists of learning of a single model f_{θ_0} that is applied to **Dynamic Networks** have architectures based on blocks

[16, 28, 30, 4] or channels [10, 2, 27] that change depending on the input sample. [16, 28] proposed an input dependent block path that decides whether a network block should be kept or dropped. [30, 4] widen the network by adding new parallel blocks and train an attention module to choose the best combination of features dynamically. [10, 2, 27] rely on feature based attention modules that reweigh features depending on the input example. [27] shows that, for object detection, objects from different domains are best detected with domain dependent features. In this paper, we propose a dynamic convolution residual branch, which adds an inputdomains *implicitly*, relying on the distribution of samples dependent residual matrix to a static kernel, to implement x. Dynamic transfer learns to adapt the parameters to fit

dynamic multi-source domain adaptation.

with any specific domains S_i and there are no rigid domain In this section, we introduce dynamic transfer for multi-boundaries. The model parameters $\theta(x)$ can be similar for source domain adaptation, in which the model is adaptive to examples from different domains and different for examples the domain implicitly, but adaptive to the input explicitly. It from the same domain. not only has better performance, but also turns multi-source The key insight is that adapting model per domain is domains into a single-source domain. achieved statistically by adapting model per sample, as each domain can be considered as a distribution of image

samples. The dynamic transfer learns to adapt model pa-Multi-source domain adaptation (MSDA) aims to trans-

 $f_{\theta(x)} = f_0 + \Delta f_{\theta(x)},$

all examples from source and target domains. The model might, for instance, map images into a latent space where all the distributions are aligned. Since the big variation among the input samples, this is a difficult problem and the model usually has sub-optimal performance on all domains. **Dynamic Transfer.** In this case the model parameters are a function of the input example x directly, i.e. the model has the form $f_{\theta(x)}$ where $x \in \mathcal{S}_1 \cup \cdots \cup \mathcal{S}_N \cup \mathcal{T}$. This is illustrated in Figure 1(b), where there exists a model per sample. Compared to the static transfer, dynamic transfer varies the model according to sample explicitly and chooses

domains. The target domain is not required to be aligned

the model to the distribution formed by the union of source

3.1. Multi-Source Domain Adaptation

fer a model learned on a source data distribution drawn

This simplifies the alignment between source and target dofrom several domains $S = \{S_1, ..., S_N\}$ to a target domains, as it is not necessary to pull all source domains and main \mathcal{T} . While the following ideas can be applied to vartarget domain together. As long as the target domain is ious tasks, we consider a classification model f_{θ} , of paaligned with any part of the union of source domains, the rameters θ , which maps images $x \in \mathcal{X}$ to class premodel can be easily adapted to the target samples. dictions $y \in \mathcal{Y} = \{1, \dots, C\}$, where C is the number Dynamic transfer has two advantages over static transof classes and \mathcal{X} is some image space. The goal is to fer. First, it turns multi-source domains into a single-source adapt the parameters θ of a model learned from a dataset domain, voiding the need for domain labels. Second, it sim- $\mathcal{D}^{\mathcal{S}} = \{(x_i^{\mathcal{S}}, y_i)\}_{i=1}^{N_{\mathcal{S}}}$ of examples from the source distribuplifies learning, since domain labels can be arbitrary. In tion $S(\mathbf{v}_i)$ is the one-hot encoding of the label of example practice, any "domain" can contain a mixture of unlabeled x^{S}) to a dataset $\mathcal{D}^{T} = \{x^{T}_{i}\}_{i=1}^{N}$ of unlabeled examples sub-domains and some of these can be shared by multiple from the target distribution. Note that, in the most gen"domains". Due to this, explicit assignment of data to doeral formulation of the problem, the domain of origin of mains can be difficult, and models learned over single doeach source example, (x_i^S, y_i) is unknown. This is ignored main can loose access to shared sub-domain data.

parameters of modern networks, it is impossible to simply predict all parameter values at inference time. The key is to restrict the model's dependence on input x to a small num-

Attention Branch Subspace Routing ber of parameters. To guarantee this, we propose a model composed by a static network and dynamic residual blocks

where f_0 represents the static component and $\Delta f_{\theta(x)}$ the Figure 3: Subspace routing of DRT: dynamic coefficients dynamic residual that depends on the input sample x. As are generated by a dynamic branch given the input x. Each usual, the residual is implemented by adding residual blocks dynamic coefficient $\pi_i(x)$ is then multiplied by a matrix Φ_i , to the various network layers. Since the static component f_0 and the K matrices are aggregated as the residual kernelis shared by all samples, static transfer is a special case of $\Delta W_0(x)$. For channel attention, softmax is replaced by the proposed approach, where $\Delta f_{\theta(x)} = 0$. This approach sigmoid and the resulting coefficients in $\Lambda(x)$ are multiis denoted as dynamic residual transfer (DRT). plied to corresponding channels of W_0 .

To implement DRT in convolution neural networks (CNNs), we represent a $k \times k$ convolution kernel as a **Combination:** the two mechanisms are combined into $C_{out} \times C_{in}k^2$ weight matrix, where C_{in} and C_{out} are the number of input and output channels. We ignore bias terms $\Delta W(x) = \Lambda(x)W_0 + \sum \pi_i(x)\Phi_i. \tag{5}$ in this discussion for the sake of brevity. DRT is implemented by applying Equation 1 to each convolution kernel Similar to squeeze-and-excitation block [11], the dynamic in a CNN, i.e. defining the network convolutions as

$$W(x) = W_0 + \Delta W(x)$$
, (2) weight attention branch that includes average pooling and

two fully connected layers (See Figure 3). A sigmoid where W_0 is a static convolution kernel matrix, and is used to normalize $\Lambda(x)$ and a softmax to normalize $\Delta W(x)$ a dynamic residual matrix. We next discuss sev- $\{\pi_i(x)\}$. As explained by [4, 30], the extra FLOPs caused eral possibilities for the latter. by dynamic coefficient generation and residual aggregation **Channel Attention:** in this case, the residual only rescales of dynamic transfer is negligible (less than 0.1% in our imthe output channels of W_0 . This is implemented as plementation) compared to the static model.

 $\Delta \boldsymbol{W}(\boldsymbol{x}) = \boldsymbol{\Lambda}(\boldsymbol{x}) \boldsymbol{W}_0,$

Subspace Routing: the dynamic residual is a linear combi-

 $\Delta W(x) = \sum \pi_i(x) \Phi_i$

trix in the corresponding weight subspaces. By choosing

these projections in an input dependent manner, the network

based attention mechanism.

nation of K static matrices Φ_i

As usual for domain adaptation problems, the DRT netwhere $\Lambda(x)$ is a diagonal $C_{out} \times C_{out}$ matrix, whose entries work is learned with a combination of two losses, are functions of x. This can be seen as a dynamic feature-

coefficients $\Lambda(x)$ and $\{\pi_i(x)\}$ are implemented by a light-

 $\mathcal{L} = \mathcal{L}_{ce} + \lambda \mathcal{L}_d,$ where λ is a hyperparameter that controls the trade-off be-

tween the loss components. The first loss

 $\mathcal{L}_{ce} = \frac{1}{N} \sum_{i} \mathbf{y}_{i}^{T} \log f_{\boldsymbol{\theta}(\boldsymbol{x}^{\mathcal{S}})}(\boldsymbol{x}_{i}^{\mathcal{S}}), \tag{7}$

 $\mathcal{L}_d = \mathcal{H}\left(f_{\boldsymbol{\theta}(\mathcal{D}^{\mathcal{S}})}(\mathcal{D}^{\mathcal{S}}), f_{\boldsymbol{\theta}(\mathcal{D}^{\mathcal{T}})}(\mathcal{D}^{\mathcal{T}})\right), \tag{8}$

whose weights depend on x. The matrices Φ_i can be seen as a basis for CNN weight space, although they are not necis the cross entropy loss over the source data $\mathcal{D}^{\mathcal{S}}$. The secessarily linearly independent. And the dynamic coefficients ond is a domain alignment loss that encourages the mini- $\pi_i(\boldsymbol{x})$ can be seen as the projections of the residual mamization of the distance between source and target domains

chooses different feature subspaces to route different x. To reduce the number of parameters and computation, where $\mathcal{D}^{\mathcal{T}}$ is the target data and \mathcal{H} a measure of discrepancy

inf,pnt,qdr clp,pnt,qdr clp,inf,qdr clp,inf,pnt clp,inf,pnt clp,inf,pnt $rel, skt \rightarrow clp \quad rel, skt \rightarrow inf \quad rel, skt \rightarrow pnt \quad rel, skt \rightarrow qdr \quad qdr, skt \rightarrow rel \quad qdr, rel \rightarrow skt$ 54.3 ± 0.64 22.1 ± 0.70 45.7 ± 0.63 7.6 ± 0.49 58.4 ± 0.65 43.5 ± 0.57 38.5 ± 0.61 Channel Attention 67.8 ± 0.46 30.9 ± 0.85 57.1 ± 0.36 6.9 ± 1.12 66.7 ± 0.42 57.4 ± 0.33 47.8 ± 0.59 Subspace Routing 69.7 ± 0.24 31.0 ± 0.56 59.5 ± 0.43 9.9 ± 1.03 68.4 ± 0.28 59.4 ± 0.21 49.7 ± 0.46 Combination 69.1 ± 0.35 **31.6** ±0.61 58.2 ± 0.25 **11.9** ±0.96 67.8 ± 0.36 58.8 ± 0.44 49.6 ± 0.50

> Table 1: Comparison of different implementations for dynamic residual transfer: Channel Attention (Equation 3), Subspace Routing (Equation 4) and Combination (Equation 5).

entire source dataset \mathcal{D}^S , i.e. there is no need for domain plementations of the dynamic transfer, (b) the number of labels and not even a difference between the single domain basis used for subspace routing (Equation 4), and (c) differand multiple domains adaptation problems. For the domain ent alignment losses \mathcal{L}_d . The default model uses subspace alignment losses commonly used in multi-source domain routing with K = 4 and is trained with the MCD [23] loss. adaptation, Equation 8 also does require the evaluation of **DRT Implementations:** Table 1 shows that all implemenpairwise distances between all source domains and target tations of DRT have significantly better adaptation perfordomain, which is not necessary in dynamic transfer. mance than the static model. The average gains are of 9.3%

learning rate is decayed by 0.1 every 5 epochs.

ated.

In this section, adaptation performance of DRT is evalutention suggests that it is not enough to re-scale the features of the static model. Routing the input x through different

subspaces appears to be more effective, although the dif-4.1. Datasets and Experimental Settings: ferences are not staggering. While combining the two ap-

Following [20], we consider two datasets, Digit-five and proaches has no additional overall benefit, the combination DomainNet [20], which contain images from several dowards beneficial for specific transfer problems. When 'infomains but shared classes. Each domain is alternatively used graph' and 'quickdraw' were used as target domains, the as the target domain and the remaining ones as the source combination model outperformed subspace routing. Since domain. All experiments are repeated with 5 times and these are the hardest transfer problems, this suggests that the enhanced dynamics of the combined implementation can be mean and variance are reported. **Digit-five:** Digit-five contains digit images from 5 domains: beneficial as the domain gap increases. It is because the en-

MNIST [12] (mt), Synthetic [6] (sy), MNIST-M [6] (mm), hanced dynamics make the model more elastic. Therefore, SVHN [19] (sv) and USPS [6] (up). These domains conit is more likely to adapt models to target domain with larger tribute 25,000 images for training and 9000 for validation, gap. On the other hand, for the problems of smaller domain with the exception of USPS which uses 29752 and 1860. respectively. Since these datasets are relatively small, LeNet source domain, as is the case for the remaining target do-[12] is used as the backbone model. A dynamic residual mains. More experiments on datasets with more domains is added on each convolutional layer. The model is trained will likely be needed to resolve this question. In any case, from scratch with initial learning rate 0.002 and SGD opti-

mizer. The learning rate is decayed by 0.1 every 100 epochs performance. and decreased to 2e - 5 in 300 epochs. **Number of Residual Basis.** The impact of the number of

images of 345 classes from 6 domains of different image For different values of $K \in \{2, 4, 6, 8\}$, DRT achieves styles: clipart (clp), infograph (inf), painting (pnt), quick- {48.8, 49.7, 49.5, 49.3}, all of which improve the adaptadraw (gdr), real (rel) and sketch (skt). Results are obtained tion performance of static transfer (38.5%) by a large marwith ImageNet [5] pretrained ResNet-101 [8]. The dynamic gin (more than 10%). Best performance is achieved with residual is only added on the 3×3 kernel of each bottleneck K = 4, although the results are not highly sensitive to this block. The networks are trained with SGD for 15 epochs parameter. with initial learning rate of 0.001 and batch size as 64. The **Alignment Loss Function.** Three different domain align-

> ment losses with different forms of \mathcal{H} (see Equation 8) were compared: ADDA [25], MCD [23] and M^3SDA [20].

for channel attention, 10.8% for combined and 11.2% for

subspace routing. The weaker performance of channel at-

inf,pnt,qdr clp,pnt,qdr clp,inf,qdr clp,inf,pnt clp,inf,pnt clp,inf,pnt $| \text{rel,skt} \rightarrow \text{clp} \quad \text{rel,skt} \rightarrow \text{inf} \quad \text{rel,skt} \rightarrow \text{pnt} \quad \text{rel,skt} \rightarrow \text{qdr} \quad \text{qdr,skt} \rightarrow \text{rel} \quad \text{qdr,rel} \rightarrow \text{skt}$ 52.1 ± 0.51 23.4 ± 0.28 47.7 ± 0.96 **13.0** ±0.72 60.7 ± 0.23 46.5 ± 0.56 40.6 ± 0.56 Source Only + **DRT** | 63.1 ± 0.62 **25.9** ±0.84 **48.4** ±1.02 6.4 ±0.98 **66.4** ±0.54 **46.8** ±0.44 **42.8** ±0.74 47.5 ± 0.76 11.4 ± 0.67 36.7 ± 0.53 14.7 ± 0.50 49.1 ± 0.82 33.5 ± 0.49 32.2 ± 0.63 ADDA+DRT 63.6 \pm 0.52 27.6 \pm 0.43 52.3 \pm 0.68 8.2 \pm 1.44 67.9 \pm 0.42 49.6 \pm 0.33 44.9 \pm 0.64 54.3 ± 0.64 22.1 ± 0.70 45.7 ± 0.63 7.6 ± 0.49 58.4 ± 0.65 43.5 ± 0.57 38.5 ± 0.61 69.7 ± 0.24 31.0±0.56 59.5±0.43 9.9±1.03 68.4±0.28 59.4±0.21 49.7±0.46 M^3 SDA- β [20] 58.6+0.53 26.0+0.89 52.3+0.55 6.3+0.58 62.7+0.51 49.5+0.76 42.6+0.64 $M^3SDA-\beta+DRT$ | 67.4±0.52 31.3±0.83 56.5±0.67 13.6±0.34 66.9±0.42 56.8±0.49 48.8±0.55

Table 2 shows that dynamic residual transfer (DRT) outparing performance in individual adaptation problems, DRT performs static transfer for all loss functions, by a large has the best performance on four of the five problems conwill give more benefits to the dynamic model than the static (MNIST-M as target domain). one. However, the gains over the 'source only' case, where **Evaluation on DomainNet Dataset.** For DomainNet [20],

problems, except when 'quickdraw' is the target domain. Table 4 shows that DRT improves on the state-of-the-

In this case, DRT is only effective with the $M^3SDA-\beta$ [20] art method- CMSS by more than 3% (49.7% vs 46.5%). loss. It is because when 'quickdraw' is the target domain, When DRT is combined with a naive self-training method the domain discrepancy is much larger and makes it harder (DRT+ST), it achieves gains of 3.9% over LtC-MSDA [26], for DRT to adapt model to this domain. Thus, $M^3SDA-\beta$ a methods that generates pseudo-labels for the target sam-[20] which proposed a more powerful alignment loss, can ples (during self-training, pseudo-labels for target samples shift 'quickdraw' closer to source domains and works better with DRT. However, the strong alignment loss will cause with source samples). Compared to the adaptation meth-'over-alignment' for the domains e.g. 'clipart' that have ods that use no domain labels (single-source), DRT immuch smaller gap. The 'over-alignment' reduces the adaptability of the dynamic model, which causes performance (49.7% vs. 43%).

DomainNet: DomainNet [20] is a dataset with 0.6 million basis K used in Equation 4 for subspace routing is ablated. 4.3. Comparisons to the state-of-the-art DRT was compared to the results in the literature for Digit Five and DomainNet dataset. In these experiments,

> DRT is implemented with subspace routing (4 basis), using the MCD loss [23], and $\lambda = 50$ in Equation 6. Evaluation on Digit Five Dataset: Table 3 shows a com-

margin (12.7%, 11.2%, 6.2% respectively). Its improved sidered. The only exception occurs when SVHN is the tarperformance is in part, due to the fact that DRT takes a get domain, where DRT achieves the second best performuch larger advantage of the domain alignment losses. It mance of all methods. Beyond this, the smallest gains occonfirms our claim that DRT simplifies the domain alignment by unifying all source domains into a single domain. since MNIST is easier to transfer to and somewhat satu-Thus the target samples are more likely to be aligned with rated. In general, the gains of DRT increase with domain the union of source domains and the same alignment loss discrepancy, reaching 5.7% for the hardest transfer problem

Table 2: Static transfer vs. Dynamic transfer evaluated on DomainNet with different domain alignment loss functions

no alignment loss \mathcal{L}_d is used in Equation 6, is only 2%. It a ResNet-101 [8] was used as backbone and DRT was commeans alignment loss is very critical for dynamic transfer. pared to 11 baselines. Among these, ADDA [25], DANN Without alignment loss, even though the model can adapt [7] and MCD [23] were developed for traditional unsuperto the entire source domain very well, it can hardly adapt to vised domain adaptation (UDA), where a single-source dotarget samples due to a large domain gap. main is assumed. The remaining are multi-source domain These conclusions also apply to the individual transfer adaptation methods that require domain labels.

> better than that of the MCD method. On average, over all pairs of source and target domains, it outperforms MCD by
>
> 4.5. Visualization more than 8%. These results show that, for problems with
>
> To obtain further insight about dynamic residual transfer

degradation compared to that given by simpler alignment Regarding individual adaptation problems, DRT achieves the best performance for all target domains other than 'quickdraw'. This can be explained by the large domain gap between 'quickdraw' and the other domains, and the fact that the MCD loss does not fare well in this problem. Better results would likely be possible by using the M³SDA loss, as was the case in Table 2.

Models

 $sy \rightarrow mt$ $sy \rightarrow mm$ $sy \rightarrow up$ $sy \rightarrow sv$ $sv \rightarrow sy$ 63.37 ± 0.74 90.50 ± 0.83 88.71 ± 0.89 63.54 ± 0.93 82.44 + 0.65 77.71 + 0.81DANN [7] 71.30+0.56 97.60+0.75 92.33+0.85 63.48+0.79 85.34+0.84 82.01+0.76 ADDA [25] 71.57 ± 0.52 97.89 ± 0.84 92.83 ± 0.74 75.48 ± 0.48 86.45 ± 0.62 84.84 ± 0.64 MCD [23] 72.50 ± 0.67 96.21 ± 0.81 95.33 ± 0.74 78.89 ± 0.78 87.47 ± 0.65 86.10 ± 0.73 DCTN [29] 70.53 ± 1.24 96.23 ± 0.82 92.81 ± 0.27 77.61 ± 0.41 86.77 ± 0.78 84.79 ± 0.72 $M^3SDA-\beta$ [20] 72.82 ± 1.13 98.43 ± 0.68 96.14 ± 0.81 81.32 ± 0.86 89.58 ± 0.56 87.65 ± 0.75 CMSS [321 75.3 \pm 0.57 99.0 \pm 0.08 97.7 \pm 0.13 **88.4** \pm 0.54 93.7 \pm 0.21 90.8 \pm 0.31 **DRT** 81.03±0.34 99.31±0.05 98.40±0.12 86.67±0.38 93.89±0.34 91.86±0.25 Table 3: Comparison between **dynamic residual transfer (DRT)** with the state-of-the-art models on Digit-five dataset. The

mm,up,sv mt,up,sv mt,mm,sv mt,mm,up mt,mm,up

source domains and target domain are shown at the top of each column.

clp rel,skt \rightarrow in 23.4 \pm 0.28	, 1	$rel,skt \rightarrow qdr$	$qdr,skt \rightarrow rel$		
$51 23.4 \pm 0.28$	1==1006		1 , . ,	$qdr,rel \rightarrow skt$	Avg
	47.7 ± 0.96	13.0 ± 0.72	60.7 ± 0.23	46.5 ± 0.56	40.6±0.56
76 11.4 ± 0.67	36.7 ± 0.53	14.7 ± 0.50	49.1 ± 0.82	33.5 ± 0.49	32.2 ± 0.63
54 22.1 ± 0.70	45.7 ± 0.63	7.6 ± 0.49	58.4 ± 0.65	43.5 ± 0.57	38.5 ± 0.61
$12 25.8 \pm 0.43$	50.4 ± 0.51	7.7 ± 0.68	62.0 ± 0.66	51.7 ± 0.19	43.0 ± 0.46
73 23.5 ± 0.59	48.8±0.63	7.2 ± 0.46	53.5±0.56	47.3 ± 0.47	38.2±0.57
26.0 ± 0.89	52.3 ± 0.55	6.3 ± 0.58	62.7 ± 0.51	49.5 ± 0.76	42.6 ± 0.64
79 26.2 ± 0.41	51.9 ± 0.20	19.1 ±0.31	57.0 ± 1.04	50.3 ± 0.67	44.3 ± 0.24
$22 21.4 \pm 0.07$	50.5 ± 0.08	15.5 ± 0.22	64.6 ± 0.16	50.4 ± 0.12	44.2 ± 0.07
5 28.7 ± 0.7	56.1 ± 0.5	16.3 ± 0.5	66.1 ± 0.6	53.8 ± 0.6	47.4 ± 0.6
8 28.0±0.20	53.6 ± 0.39	16.0 ± 0.12	63.4 ± 0.21	53.8 ± 0.35	46.5 ± 0.24
$\frac{1}{24}$ 31.0±0.56	59.5±0.43	9.9 ± 1.03	68.4 ± 0.28	59.4 ± 0.21	49.7±0.46
. - 51.0±0.50	61.0±0.32	12.3 ± 0.38	71.4 \pm 0.23	60.7 ±0.31	51.3 ±0.32
7			24 31.0 \pm 0.56 59.5 \pm 0.43 9.9 \pm 1.03 21 31.6 \pm 0.44 61.0 \pm 0.32 12.3 \pm 0.38		

Table 4: Comparison between dynamic residual transfer (DRT) with the state-of-the-art models on DomainNet. ('DRT+ST' represents the combination between dynamic residual transfer and self-training for domain adaptation)

> adaptation was performed from each of the other five domains (sources). The average and best performance among not require domain labels. Its universal nature makes it irthese adaptations is shown in Table 5, for each target domain. DRT again significantly outperforms the previous domain adaptation methods. For example, when 'clipart' is the figure out how to adapt the network to all settings. There is target domain, its average adaptation performance is 10.5\% no need to even define "source domains."

hundreds of classes, dynamic residual transfer can lead to (DRT), we visualized the dynamic coefficients of Equation very large adaptation gains even in the traditional domain 4 with t-SNE [18]. For each sample, we created a vector adaptation setting. This confirms the claim that even these $\Pi = \{\pi^l(x)\}, i \in \{1, 2, ..., K\}, l \in \{1, 2, ..., L\}$ by concateproblems tend to have many sub-domains. When this is the α nating the dynamic coefficients from the α network layers. case, the ability of dynamic residual transfer to adapt the

The vectors Π from different target domains are visualized model on a per-example basis can be a significant asset. in Figure 4. A more detailed visualization is given in Fig-Finally, comparing the results of Tables 4 and 5 shows ure 5, by splitting Π into Π_{low} and Π_{high} , which include that DRT trained on multi-source domains performs 8.2% the coefficients from lower and higher network layers. For better (49.7% vs. 41.5%) than the average of the best single-brevity, Figure 5, only visualizes the model trained with source domain transfers. This improvement is about 2% 'clipart' and 'real' as target domain and the other domains

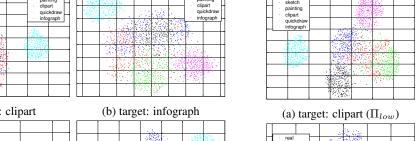
 $rel, skt \rightarrow clp$ $rel, skt \rightarrow inf$ $rel, skt \rightarrow pnt$ $rel, skt \rightarrow qdr$ $qdr, skt \rightarrow rel$ $qdr, rel \rightarrow skt$ ADDA [25] 28.2/39.5 9.3/14.5 20.1/29.1 **8.4/14.9** 31.1/41.9 21.7/30.7 19.8/28.4 MCD [23] 31.4/42.6 13.1/19.6 24.9/42.6 2.2/3.8 35.7/50.5 23.9/33.8 21.9/32.2 DRT 41.9/56.2 19.6/26.6 35.3/53.4 8.0/12.2 44.5/55.5 35.0/44.8 30.7/41.5

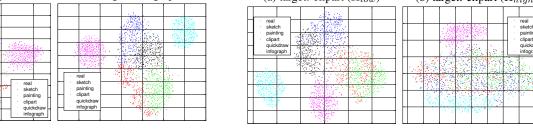
'sketch'. (Best view in color)

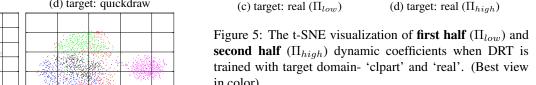
identifiable clusters, which confirms our claim that adapt-

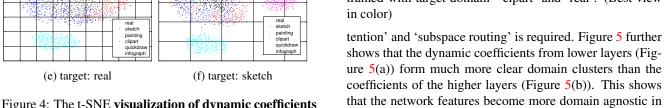
Table 5: Single source domain adaptation performance on DomainNet. Each column, shows the average/best classification accuracy for transfer from all source to the specified target domain.

inf,pnt,qdr clp,pnt,qdr clp,inf,qdr clp,inf,pnt clp,inf,pnt clp,inf,pnt









 $\Pi = \{\pi_i^l(x)\}$ when DRT is trained with target domainthe higher layers, confirming the effectiveness of DRT to

'clipart', 'infograph', 'painting', 'quickdraw', 'real' and reduce domain discrepancies.

model to samples. Secondly, the distance among clusters source domain adaptation, in which the model parameters reflects domain shifts that explain how adaptation perforare not static but adaptive to input samples. Dynamic transmance varies with target domain. For example, the fact that fer mitigates conflicts across multiple domains and unithe dynamic coefficients of 'quickdraw' are always quite fies multiple source domains into a single source domain, different from others, explains why adaptation performance which simplifies the alignment between source and target is weaker when this is the target domain. Thus for 'quick-domains. Experimental results show that dynamic trans-

ing model across domains can be achieved through adapting

In this paper, we introduce dynamic transfer for multi-

4.4. Single-Source to Single-Target Adaptation by many approaches e.g. [20, 15], that assume a source

3.3. Dynamic Residual Transfer all domains together. The model adaptation given by the samples. [15] uses mutual learning techniques to align feathe matrices can be further simplified into 1×1 convolution between feature distributions of the source and target dobetter than that given by MCD (8.2% vs. 6.3%). This show similar trend. draw', either a more powerful alignment loss is needed to fer achieves a better adaptation performance compared to dataset $\mathcal{D}^{\mathcal{S}} = \{(\boldsymbol{x}_{i}^{\mathcal{S}}, \mathbf{v}_{i}, z_{i})\}_{i=1}^{N_{\mathcal{S}}}$ contains domain labels kernels and applied to the narrowest layer of the bottleneck mains. \mathcal{H} can be any distance function previously proposed They are representative of previously proposed losses for reparison to 6 baselines on Digit Five. DRT outperforms all

The performance of DRT on the traditional domain adapdynamic transfer can be generalized to target domain easily ture distributions among pairs of source and target domains. shows that considering a variety of source domains imshift the samples close enough to the source domains to enwhen the target domain is shifted to the space formed by Other methods focus on the joint alignment of the feature $z_i \in \{1, \dots, N\}$ and aligning pairs of domains. We refer to The main difficulty of dynamic transfer is the model architecture in ResNet [8]. In this case, only C_{in} rows of for domain adaptation, e.g. the MMD [17] or adversarial An ablation study was performed on DomainNet to eval-ducing single-source domain shift at the domain level and other methods, beating the state of the art (CMSS) by more tation problem (single-source domain) was also evaluated proves domain adaptation performance, especially when ded into the dynamic coefficients $\{\pi_i^l(x)\}$. This can be able the dynamic model adapted to this domain or a more tation. We hope this paper can give a new understanding entire source domain. In this way, dynamic transfer shifts distributions of all domains. [26] models interactions bethis a domain supervised multi-source domain adaptation. $f_{\theta(x)}$ can be difficult to learn. Given the large number of learning [25]. Note that the two losses above operate on the uate the three key components of DRT: (a) the three imthan 1%, despite a much simpler implementation. Comdynamic residual transfer is used. A main advantage of observed by samples from same domains tend to group in complex dynamic model e.g. combination of 'channel at-about multi-source domain adaptation.