
Network Transplanting

Quanshi Zhang^{††}, Yu Yang[†], Ying Nian Wu[†], and Song-Chun Zhu[†]

[†]University of California, Los Angeles ^{††}Shanghai Jiao Tong University.

Abstract

This paper focuses on a novel problem, *i.e.* transplanting a category-and-task-specific neural network to a generic, distributed network without strong supervision. Like playing LEGO blocks, incrementally constructing a generic network by asynchronously merging specific neural networks is a crucial bottleneck for deep learning. Suppose that the pre-trained specific network contains a module f to extract features of the target category, and the generic network has a module g for a target task, which is trained using other categories except for the target category. Instead of using numerous training samples to teach the generic network a new category, we aim to learn a small adapter module to connect f and g to accomplish the task on a target category in a weakly-supervised manner. The core challenge is to efficiently learn feature projections between the two connected modules. We propose a new distillation algorithm, which exhibited superior performance. Our method without training samples even significantly outperformed the baseline with 100 training samples.

1 Introduction

Problem: In this paper, we focus on a new problem, *i.e.* how to merge several specific convolutional networks that are pre-trained for different categories and different tasks to a generic, distributed neural network. In recent years, convolutional neural networks (CNNs) [14, 13, 8, 15, 11] have achieved superior performance in many visual tasks, such as object classification and detection. However, current studies mainly either learn a neural network for each specific task, or train a neural network for multiple tasks and multiple categories at the same time. In contrast, asynchronously transplanting pre-trained specific networks to a generic network allows the generic network to incrementally grow new modules for new tasks and new categories.

As shown in Fig. 1, the target generic network has a distributed structure, and we call it a *transplant network*. We divide the transplant network into category modules, task modules, and adapters. Each category module extracts general features for a category, and each task module is learned for a certain task and is shared by different category modules. Category modules and task modules are transplanted from specific networks. We also design a small adapter module to connect each pair of category and task modules.

We believe that learning via network transplanting will exhibit significant efficiency in representation, learning, and application.

1. **Generic representations:** A theoretical solution to network transplanting makes an important step towards building a universal, distributed network for all categories and tasks, which is one of the ultimate objectives for high-level artificial intelligence.
2. **Interpretability & compactness:** The transplant network has a compact structure, where each category/task module is functionally meaningful and is oriented to different tasks/categories.
3. **For learning:** Unlike fine-tuning, network transplanting used very limited training samples. Network transplanting does not require people to prepare training samples of all tasks and all

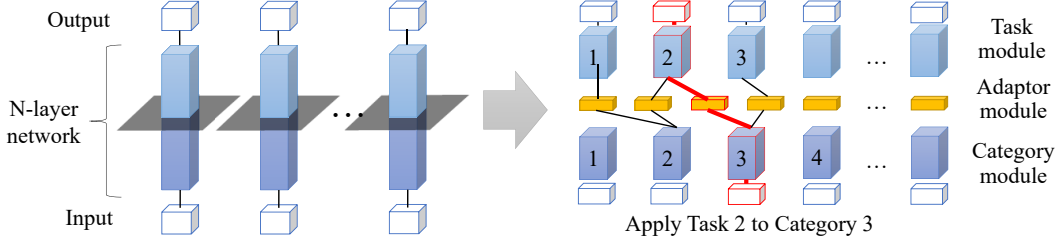


Figure 1: Problem of building a generic, distributed transplant network. We aim to explore a theoretical solution to asynchronously merging specific networks into a generic transplant network in a weakly-supervised (or even unsupervised) manner. We divide the transplant network into three types of modules, *i.e.* the task modules, category modules, and adapters. A category module, *e.g.* a cat module, provides cat features to different task modules. A task module, *e.g.* a segmentation module, services for various categories. We incrementally distill specific networks to grow new task/category modules for the transplant network on-the-fly, which enables asynchronous and weakly-supervised learning for different tasks. In addition, the transplant network has an interpretable, distributed structure. People can conduct a certain task by manually linking its corresponding modules (red).

categories at the same time for learning. This is quite practical when the transplant network deals with a large number of categories and tasks.

4. **For application:** Just like building LEGO blocks, for a specific application, people can manually activate its corresponding category module and task module to accomplish this application.

In addition, we need to consider the following four issues to learn a transplant network.

1. **Generality:** We expect a generic solution that can be applied to different categories and tasks.
2. **Asynchronous learning:** We need to learn the transplant network by transplanting specific networks one-by-one in an asynchronous manner. It is usually difficult to collect all pre-trained specific networks at the same time. Thus, we incrementally grow modules of the transplant network given each specific network on the fly.
3. **Weakly-supervised/unsupervised learning:** Network transplanting needs to be conducted with a limited number of training samples or even without any training samples, which ensures broad applicability of our method.
4. **Local transplantation vs. global fine-tuning of category and task modules:** During network transplanting, we exclusively learn the adapter without fine-tuning the task and category modules. It is because each category/task module is potentially responsible for multiple tasks/categories, and we do not want the insertion of a new module to damage the generality of existing modules.

When we have constructed the transplant network with lots of category modules and task modules, we may fine-tune each task module using features of different categories or modify each category module using gradients from multiple task modules to ensure their generality.

Task: As shown in Fig. 2, we are given 1) a specific network with a category module f that is pre-trained for a certain task and 2) a transplant network with a task module g_S that are learned to accomplish the same task for other categories. We may (and may not) be also given a few training samples of the target category for the target task. The task is to learn the adapter between f and g_S that projects the output feature space of f to the input feature space of g_S . We fix the task and category modules during the learning process to avoid damaging their generality.

The core technical challenge of this study is that because we can only modify the lower adapter to fit a fixed upper task module g_S , without permission to push the fixed g_S towards the lower adapter. There is vast *forgotten space* for the input feature of the upper g_S , *i.e.* input features in the forgotten space cannot pass their information to the final network output. Learning correct projection from the output of f to the input of g_S is difficult considering the forgotten space. However, the comfortable regime of exclusively learning the task loss is to refine the upper g_S based on feature responses of the lower adapter, instead of discovering potentially effective input feature spaces of the upper g_S . Meanwhile, it is difficult to directly distill representations from the specific network to the adapter. It

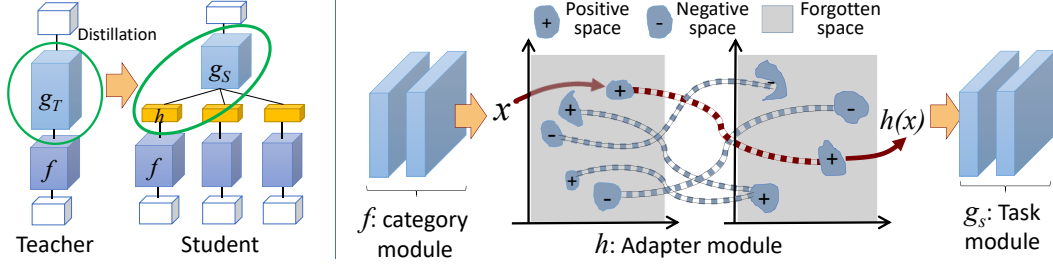


Figure 2: Overview. (left) Our method distills knowledge representations from a pre-trained teacher CNN ($f \circ g_T$) for a specific category to learn a generic student network ($f \circ h \circ g_S$). Let the teacher network have learned the category module f , and let the student network have modeled a task module g_S for other categories. Thus, our task is to transplant g_S to f by learning an adapter h , so that the task module g_S can deal with the new category f . (right) During the transplantation, the adapter h needs to learn the potential projection between f 's output feature space and g_S 's input feature space. Well-learned network modules usually have huge *forgotten space*, and features $h(x)$ in the forgotten space cannot pass through the task module g_S . Incorrect space projections may paralyze the network.

is because except for the final network output, high-layer features in the task module of the transplant network and those in the specific network are not aligned in semantic meanings.

Therefore, in this study, we propose a strategy, namely *back distillation*, to learn the adapter when we fix parameters of the category and task modules. The back-distillation algorithm forces the combination of the adapter and the task module in the transplant network to approximate representations of upper layers in the pre-trained specific network. This algorithm requires the specific network and the transplant network to have similar gradients *w.r.t.* the output feature of the category module for distillation. In experiments, our back-distillation algorithm exhibited superior performance to baseline methods in the learning of the adapter. Our method without any training samples even significantly outperformed the baseline with 100 training samples in Table 1.

We can summarize contributions of this study as follows. 1) We propose a theoretical solution to a new problem, *i.e.* asynchronously merging specific neural networks into a generic, distributed transplant network in a weakly-supervised or unsupervised manner. 2) We develop a back-distillation algorithm as a theoretical solution to the challenge of learning space projections between the output of the category module and the input of the task module. 3) Experiments showed that our method significantly outperformed baselines in different applications.

2 Related work

Network structures: Because network transplanting is a new concept in machine learning, we would like to discuss its connections to different state-of-the-art algorithms. Firstly, we can consider the distributed transplant network as a new distributed structure for networks, which disentangles a black-box network into different meaningful modules. Similarly, some studies have explored new representation structures instead of neural networks, such as forests and decision trees [12, 29, 6, 24]. Automatic learning of optimal network structures [30, 17, 31, 28] has also received increasing attention in recent years. Shazeer *et al.* [20] learned a distributed structure of a large neural network, which contained thousands of sub-networks.

Interpretability: Dividing a network into functionally meaningful modules makes the network interpretable. In comparisons, other studies of enhancing network interpretability mainly focus on either learning disentangled, interpretable representations for filters/capsules in middle layers [27, 25, 19] or unsupervisedly learning meaningful input codes of generative networks [3, 9].

Distillation: Hilton *et al.* [10] proposed the concept of network distillation to transfer representations between networks. Some recent studies [6, 24] distilled network representations into decision trees. Anil *et al.* [2] proposed an online distillation method to efficiently learn distributed networks. Zagoruyko *et al.* [26] distilled the attention distribution from the teacher network to the student network, which is related to our back-distillation technique. In contrast, we design pseudo-gradients, instead of using real gradients for distillation, which reduce effects of incorrect feature maps during

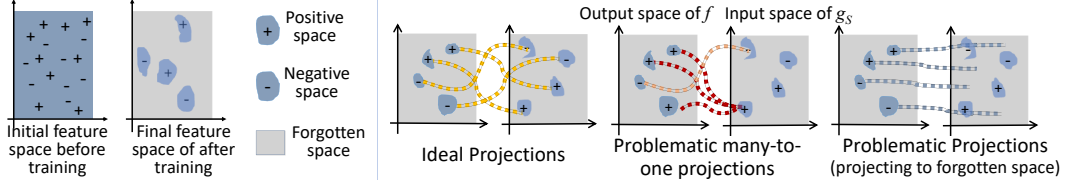


Figure 3: Feature space of a middle layer of a CNN. (left) When we initialize parameters of a CNN, middle-layer features randomly cover all area in the feature space. The learning process forces the CNN to focus on certain positive and negative feature spaces and produces vast forgotten space. (right) We illustrate three examples of space projection of the adapter.

transplantation. We also balance magnitudes of neural activations between two networks and augment pseudo-gradients, which are necessary for network transplanting.

Meta-learning: Meta-learning [5, 1, 16, 22] aims to extract generic knowledge shared by different tasks/categories/models to guide the learning of a specific model. In contrast, the distillation-based learning mainly transfers knowledge between networks to ensure a high learning efficiency.

3 Algorithm

3.1 Generic transplant networks

As shown in Fig. 2(left), we are given a specific network for a task *w.r.t.* a certain category, namely the *teacher net*. We define the bottom m layers of the teacher net as the category module f for feature extraction. We consider the upper layers g_T as a specific task module. We are also provided with a generic transplant network with a generic task module g_S , which has been learned for multiple categories other than the teacher net’s category.

Our task is to transplant g_S to f by learning an adapter h with parameters θ_h , so that the task module g_S can deal with the new category module f . We regard the transplant work as the *student net*. Note that we do not further fine-tune the category module f or the task module g_S to avoid decreasing the generality of these modules. Let x denote the output feature map of the category module f given an image I , i.e. $x = f(I)$. y_T and y_S are referred to as final outputs of the teacher net and the student net, respectively. We use the student net to approximate the teacher net, as follows.

$$y_T = g_T(x), \quad y_S = g_S(h(x)), \quad y_T \approx y_S \Rightarrow g_S(h(\cdot)) \propto g_T(\cdot) \quad (1)$$

This equation shows that the combination of the adapter h of the generic task module g_S needs to approximate feature representations of the specific task module g_T .

3.2 Problem of space projection for network transplanting

It is a challenge to let an adapter h project the output feature space of f to the input feature space of g_S . The information bottleneck theory [23] shows that a neural network selectively omits (*forgets*) certain space of middle-layer features and gradually focuses on discriminative features during the learning process (see Fig. 3(left)).

Thus, there is vast *forgotten space* for both the output of f and the input of g_S . The forgotten input space of g_S contains input features that cannot pass through all ReLU layers in g_S and reach y_S . The forgotten output space of f contains features that f cannot generate.

Fig. 3(right) illustrates ideal and problematic projections between the input of g_S and the output of f that an adapter may encode. A typical problematic projection is to project a feature x to a forgotten input space of g_S , which makes x ’s information blocked by a ReLU layer. Another typical problematic projection is many-to-one projection, which limits the diversity of features and decreases the representation power of the student net. More crucially, when we only optimize the task loss, initial many-to-one space projections may significantly affect the further learning process, because the back-propagation uses the current space projection as anchors to modify the network.

To ensure ideal projections, we force the attention/gradient of the student net *w.r.t.* x to approximate that of the teacher net, which is a necessary condition of $g_S(h(\cdot)) \approx g_T(\cdot)$.

$$\left. \begin{aligned} y_S &= g_S(h(x)) \\ y_T &= g_T(x) \\ y_T &\approx y_S \end{aligned} \right\} \Rightarrow g_S(h(\cdot)) \approx g_T(\cdot) \Rightarrow D_S \propto D_T, \quad \text{s.t.} \quad \begin{aligned} D_S([\mathbf{X}_h, \mathbf{X}_S], [\theta_h, \theta_S]) &\stackrel{\text{def}}{=} \frac{\partial L}{\partial x} \\ D_T(\mathbf{X}_T, \theta_T) &\stackrel{\text{def}}{=} \frac{\partial L}{\partial x} \end{aligned} \quad (2)$$

where L denotes the loss of the target task. $D_T(\cdot)$ and $D_S(\cdot)$ are the attention/gradient *w.r.t.* x of the teacher net and that of the student net, respectively. θ_T , θ_S , and θ_h denote parameters of g_T , g_S , and h . \mathbf{X}_T , \mathbf{X}_S , and \mathbf{X}_h denote sets of feature maps of conv-layers in g_T , g_S , and h .

Consequently, we formulate a new network-distillation loss, namely *back distillation*, as follows.

$$\min_{\theta_h} \text{Loss}, \quad \text{Loss} = L(y_S, y^*) + \lambda \cdot \|\alpha D_S - D_T\|^2 \quad (3)$$

where y^* denotes the ground-truth label, and α is a scaling scalar.

3.3 Learning via back distillation

Gradients D_S and D_T defined in the above loss are functions of feature maps $[\mathbf{X}_h, \mathbf{X}_S]$ and \mathbf{X}_T . However, in early epochs of learning, it is difficult to obtain the optimal $[\mathbf{X}_h, \mathbf{X}_S]$ that minimizes the distillation loss. Thus, to ease the learning process, we design pseudo-gradients D'_S, D'_T following the paradigm in Equation (4b), which are agnostic with regard to feature maps, to replace D_S and D_T in the loss, respectively. We assume $g_S(h(\cdot)) \approx g_T(\cdot) \Rightarrow D'_S \propto D'_T$.

$$D(\mathbf{X}, \theta) \stackrel{\text{def}}{=} \frac{\partial L}{\partial y} \frac{\partial y}{\partial x^{(n)}} \cdots \frac{\partial x^{(m+1)}}{\partial x^{(m)}} \Big|_{x=x^{(m)}} = f'_{\text{conv}} \circ f'_{\text{relu}} \circ f'_{\text{pool}}^{\text{max}} \circ \cdots \circ f'_{\text{conv}} \left(\frac{\partial L}{\partial y} \right) \quad (4a)$$

$$D'(\theta) \stackrel{\text{def}}{=} G_{y'} \frac{\partial y'}{\partial x^{(n)}} \cdots \frac{\partial x^{(m+1)}}{\partial x^{(m)}} \Big|_{x=x^{(m)}} = f'_{\text{conv}} \circ f'_{\text{dummy}} \circ f'_{\text{pool}}^{\text{avg}} \circ \cdots \circ f'_{\text{conv}}(G_{y'}) \quad (4b)$$

where $f'_1 \circ f'_2(\cdot) \stackrel{\text{def}}{=} f'_1(f'_2(\cdot))$, each f' is the derivative of a function f . of a specific layer for back-propagation. $x^{(m)} \in \mathbf{X}$ denotes the feature map of the m -th layer. Just like $\frac{\partial x^{(m+1)}}{\partial x^{(m)}}$, $G_{y'}$ represents a random initial gradient *w.r.t.* the pseudo-output y' . Please see below paragraphs of data augmentation for details of $G_{y'}$.

Both D'_S and D'_T follow the paradigm in Equation (4b), in order to make the gradient agnostic with regard to \mathbf{X} . The derivative of a convolution operation *w.r.t.* the feature map, f'_{conv} , is independent with the input feature map of the conv-layer¹. We remove all dropout layers and replace max-pooling layers $f'_{\text{pool}}^{\text{max}}$ with average-pooling layers $f'_{\text{pool}}^{\text{avg}}$. We also revise the derivative of the ReLU layer as either $f'_{\text{dummy}}^{\text{1st}}(\frac{\partial L}{\partial x^{(k)}}) = \frac{\partial L}{\partial x^{(k)}}$ or $f'_{\text{dummy}}^{\text{2nd}}(\frac{\partial L}{\partial x^{(k)}}) = \frac{\partial L}{\partial x^{(k)}} \odot \mathbf{1}(x_{\text{rand}} > 0)$, where x_{rand} is a random feature map. For each input image, we set the same random tensor $x_{\text{rand}} \in [-1, +1]^{s_1 \times s_2 \times s_3}$ for both g_S and g_T to make D'_S and D'_T comparable with each other. Please see Section 4 for detailed settings.

Note that we exclusively used all above modifications of ReLU, pooling, and dropout layers for back distillation. Besides D'_S, D'_T , we still compute original gradients D_S and D_T to optimize $L(y, y^)$.*

In this way, we conduct the back-distillation algorithm by minimizing $\min_{\theta_h} \text{Loss} = L(y_S, y^*) + \lambda \cdot \|\alpha D'_S - D'_T\|^2$. The distillation loss can also be directly optimized by further propagating gradients of gradient maps to the upper layers, and we may consider this as *back-back-propagation*¹.

Data augmentation: We augment different values of $G_{y'}$. We use each $G_{y'}$ to produce a pair of D'_S and D'_T for back distillation. Given a network for object segmentation, its output is a tensor $y_S \in \mathbb{R}^{H \times W \times C}$, where H and W denote the height and width of the output image, and C indicates the number of segmentation labels. We randomly sample $G_{y'} \in [-1, +1]^{H \times W \times C}$ for each image.

Given a network for single-category classification, its output y_S is a scalar. Nevertheless, we can also generate a random matrix $G_{y'} \in [-1, +1]^{S \times S \times 1}$ for each input image ($S=7$ in experiments), which produces two enlarged pseudo-gradient maps D'_S, D'_T for back distillation¹.

¹Please see supplementary materials for details.

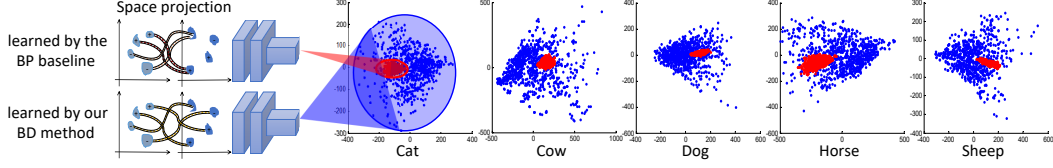


Figure 4: Comparison of the projected feature spaces. For each category, blue points indicate 4096-d $fc8$ features of different images, when we used our method to learn the adapter. Red points correspond to $fc8$ features of different images, when the adapter was learned by only using the task loss, *i.e.* the OL baseline. We visualize the first two principal components of $fc8$ features. Because the OL baseline usually learned problematic space projections (*e.g.* many-to-one projections and projections to forgotten spaces in Fig. 3), most information in $h(x)$ could not pass through ReLU layers to the $fc8$ layer. Therefore, when the adapter was learned based on the OL baseline, units in the $fc8$ layer could only be weakly triggered, and many-to-one projections decreased the diversity of $fc8$ features.

4 Experiments

We designed three experiments to evaluate the proposed method, including 1) learning toy transplant networks by inserting adapters between middle layers of pre-trained CNNs, 2) learning transplant networks for object classification, and 3) learning transplant networks for object segmentation.

Theoretically, our back-distillation strategy decreases the demand for training samples. We learned adapters in transplant networks with limited numbers of samples (*i.e.* 10, 20, 50, and 100 samples) without fine-tuning the category module or the task module. Learning adapters without changing the upper task module presents the core technical challenge. We even tested the performance of learning without any training samples in Experiment 1 by optimizing the distillation loss without considering the classification loss.

We compared our back-distillation (BD) method with two baselines in the scenario of weakly-/unsupervised learning. The first baseline only optimized the task loss without distillation, namely *ordinary learning* (OL). The second baseline is the traditional distillation method [10] (TD), which transferred the output of the specific network to the transplant network. We tested TD in object segmentation, because unlike single-category classification, segmentation outputs had rich correlations between soft output labels. All the three methods only learned adapters to enable a fair comparison.

Datasets and implementation details: We learned transplant networks for object classification in the first two experiments, and learned transplant networks for object segmentation in the third experiment. For the learning of object segmentation, we learned fully convolutional networks (FCNs) proposed in [18]. We merged FCNs for different categories to construct the transplant network. We followed experimental settings in [18] to learn an FCN for each category, which used the Pascal VOC 2011 dataset (with segmentation labels on 8498 PASCAL training images that were collected by Hariharan *et al.* [7] for object segmentation). Note that we only learned and merged five FCNs for five animal categories in the Pascal VOC dataset, *i.e.* the cat, cow, dog, horse, and sheep categories. It is because these five categories share similar object structures, which make middle-layer features of an FCN for a category were transferable to other categories. For the learning of object classification, we followed experimental settings in [27] that used the PASCAL-Part Dataset [4] to learn CNNs for single-category classification. Similarly, we also merged five specific CNN that were pre-trained for the above five animal categories to construct a generic transplant network.

The adapter contained several conv-layers, and a ReLU layer follows each conv-layer. For robust network transplanting, we further inserted a “re-scaling” layer¹ between the category module and the adapter, which normalized the scale of category features x . We formulated the re-scaling layer as $x^{\text{out}} = \beta \cdot x$, where $\beta = \mathbb{E}_{I \in \mathbf{I}_S} [\|f_S(I)\|_F] / \mathbb{E}_{I \in \mathbf{I}_T} [\|f(I)\|_F]$. \mathbf{I}_T and \mathbf{I}_S denote the image set of the target category and the image set of categories that had been already modeled by the transplant network, respectively². $f_S(I)$ denotes the input feature of g_S when we input I to the transplant network for inference. $\|\cdot\|_F$ denotes the Frobenius norm. Similarly, we also normalized the gradient *w.r.t.* x by setting $\alpha = \mathbb{E}_{I \in \mathbf{I}_T} [D'_T] / \mathbb{E}_{I \in \mathbf{I}_S} [D'_S]$. In addition, to further enlarge the dissimilarity between output

²Because we used the task module in the dog network as the generic task module g_S , we got $\mathbf{I}_S = \mathbf{I}_{\text{dog}}$.

Table 1: Error rate of classification when we insert one conv-layer and one ReLU layer to a pre-trained CNN as the adapter.

# of samples		cat	cow	dog	horse	sheep	Avg.
100	OL	12.89	3.09	12.89	10.82	9.28	9.79
	BD	1.55	0.52	3.61	1.55	1.03	1.65
50	OL	13.92	15.98	12.37	16.49	15.46	14.84
	BD	1.55	0.52	3.61	1.55	1.03	1.65
20	OL	16.49	26.80	28.35	32.47	25.77	25.98
	BD	1.55	0.52	3.09	1.55	1.03	1.55
10	OL	39.18	39.18	35.05	41.75	38.66	38.76
	BD	1.55	0.52	3.61	1.55	1.03	1.65
0	OL	—	—	—	—	—	—
	BD	1.55	0.52	4.12	1.55	1.03	1.75

Table 2: Error rate of classification when we insert three conv-layers and three ReLU layers to a pre-trained CNN as the adapter.

# of samples		cat	cow	dog	horse	sheep	Avg.
100	OL	9.28	6.70	12.37	11.34	3.61	8.66
	BD	1.03	2.58	4.12	1.55	2.58	2.37
50	OL	14.43	13.92	15.46	8.76	7.22	11.96
	BD	3.09	3.09	4.12	2.06	4.64	3.40
20	OL	22.16	25.77	32.99	22.68	22.16	25.15
	BD	7.22	6.70	7.22	2.58	5.15	5.77
10	OL	36.08	32.99	31.96	34.54	34.02	33.92
	BD	8.25	15.46	10.31	13.92	10.31	11.65
0	OL	—	—	—	—	—	—
	BD	50.00	50.00	50.00	49.48	50.00	49.90

features of different category modules, we randomly reordered filters in the top conv-layer of the category module f^1 .

4.1 Experiment 1: Adding adapter layers to pre-trained CNNs

In this experiment, we conducted a toy test, *i.e.* inserting and learning an adapter into a pre-trained CNN, to prove the effectiveness of the proposed method. More challenging tests of transplanting modules between different specific networks will be conducted in Experiments 2 and 3.

We used VGG-16 networks [21] learned for single-category classification on the PASCAL-Part Dataset [4], which achieved error rates of 1.6%, 0.6%, 4.1%, 1.6%, and 1.0% for the classification of the cat, cow, dog, horse, and sheep categories, respectively (these CNNs were strongly supervised using all training samples). We considered the first five layers of the VGG-16 (including two conv-layers, two ReLU layers, and one pooling layer) as the category module and regarded the upper layers as the task module. Then, we designed two types of adapters, *i.e.* the adapter with a single conv-layer and a single ReLU layer and the adapter with three conv-layers followed by three ReLU layers. Each conv-layer in these adapters contained M filters, each of which was a $3 \times 3 \times M$ tensor, where M is the channel number of x . We added zero padding to ensure that the output of the adapter was of the same size as its input.

In this experiment, we simply set $G'_y = 1$ without any further data augmentation. Instead, we used the revised dummy ReLU layer in Equation (4b) to ensure the value diversity of D'_S, D'_T for learning. More specifically, we used $f'_{\text{dummy}} = f_{\text{dummy}}^{\text{2nd}}$ to compute gradients of ReLU layers in the task module, and used $f'_{\text{dummy}} = f_{\text{dummy}}^{\text{1st}}$ to compute gradients of ReLU layers in the adapter³. We set $\lambda = 10.0/\mathbb{E}_{I \in \mathcal{I}_T}[D'_T]$ for object classification in Experiments 1 and 2.

In Fig. 4, we compared the space of $fc8$ features, when we used our method and the OL baseline, respectively, to learn the three-layer adapter. The adapter learned by our method passed much stronger information to the final $fc8$ layer and yielded more diverse features. It demonstrates that our method better avoided problematic projections in Fig. 3 than the OL baseline.

In Tables 1 and 2, we compared our method with the OL baseline when we added the single-layer adapter and the three-layer adapter to the pre-trained CNN, respectively. The back-distillation strategy greatly reduces the demand for training samples. Thus, we tested our method with a few (10–100) or **even without** any training samples. Table 1 shows that compared to the 9.79%–38.76% error rates of the OL baseline, our method yielded a significant lower classification error (1.55%–1.75%). Even **without any training samples**, our method still outperformed the OL methods with 100 training samples. Note that when the adapter contained three conv-layers, our back-distillation method was hampered without given any training samples. When we learned a complex adapter without training samples, our method may produce a biased short-cut solution.

4.2 Experiment 2: Learning transplant network for object classification

In this experiment, we transplanted a generic task module to a specific category module for object classification. We merged VGG-16 networks [21] learned for single-category classification on the PASCAL-Part Dataset [4] to construct the transplant network. We considered the first five layers of

³All derivative functions in Equation (4b) are only used for distillation, which will not affect the back-propagation of $L(y, y^*)$'s gradients.

Table 3: Error rate of single-category classification when we transplant the task module from a pre-trained *dog* network to the network of the target category. The adapter contains three conv-layers.

# of samples		cat	cow	horse	sheep	Avg.	# of samples		cat	cow	horse	sheep	Avg.
100	OL	20.10	12.37	18.56	11.86	15.72	20	OL	31.96	37.11	39.69	35.57	36.08
	BD	9.79	5.67	8.25	4.64	7.09		BD	21.13	35.57	32.47	22.68	27.96
50	OL	22.68	19.59	19.07	14.95	19.07	10	OL	41.75	37.63	44.33	33.51	39.31
	BD	10.82	18.04	13.92	5.15	11.98		BD	34.02	42.27	44.85	33.51	38.66

Table 4: Pixel accuracy of object segmentation when we transplant the task module from a *dog* network to the network of the target category. The adapter contains a conv-layer and a ReLU layer.

# of samples		cat	cow	horse	sheep	Avg.	# of samples		cat	cow	horse	sheep	Avg.
100	OL	76.54	74.60	81.00	78.37	77.63	20	OL	71.13	74.82	76.83	77.81	75.15
	TD	74.7	80.2	78.1	80.5	78.3		TD	71.2	74.8	76.1	78.1	75.0
	BD	85.17	90.04	90.13	86.53	87.97		BD	84.03	88.37	89.22	85.01	86.66
50	OL	71.30	74.76	76.83	78.47	75.34	10	OL	70.46	74.74	76.49	78.25	74.99
	TD	68.3	76.5	78.6	80.6	76.0		TD	70.5	74.7	76.8	78.3	75.1
	BD	83.14	90.02	90.46	85.58	87.30		BD	82.32	89.49	85.97	83.50	85.32

the VGG-16 as the category module and regarded the upper layers as the task module. We used the task module of the CNN of the *dog*⁴ category as a generic task module, and we transplanted category modules of other four categories to this generic task module. The adapter contained three conv-layers followed by three ReLU layers. Each conv-layer in the adapter contained M filters, each of which was a $1 \times 1 \times M$ tensor, which made its input and output feature maps have the same size.

In this experiment, we used data augmentation techniques in Section 3.3 to generate G'_y . Following Equation (4b), we used the $f'_{\text{dummy}} = f_{\text{dummy}}^{\text{1st}}$ operation to compute gradients of ReLU layers in both the task module and the adapter². The only exception was the lowest ReLU layer of the task module, for which we applied $f'_{\text{dummy}} = f_{\text{dummy}}^{\text{2nd}}|_{x_{\text{rand}}}$. We generated $x_{\text{rand}} = [x', x', \dots, x']$ for each input image by concatenating s_3 matrices x' along the third dimension, where $x' \in [-1, +1]^{s_1 \times s_2 \times 1}$ contained 20%/80% positive/negative elements.

Table 3 evaluates our method when we learned adapters of transplant networks for object classification. We tested our method with a few (10–100) training samples. When there were more than 50 training samples, our method yielded about a half classification error of the OL baseline.

4.3 Experiment 3: Learning transplant network for object segmentation

In this experiment, we transplanted a generic task module to a specific category module for object segmentation. We learned five FCNs based on the VGG-16 [21] for single-category segmentation following experimental settings in [18]. These FCNs achieved pixel-level segmentation accuracies (defined in [18]) of 95.0%, 94.7%, 95.8%, 94.6%, and 95.6% for the cat, cow, dog, horse, and sheep categories, respectively (these FCNs were strongly supervised by all training samples). We considered the first five layers of the FCN as the category module and regarded the upper layers as the task module. We used the task module of the CNN of the *dog* category as a generic task module. We transplanted category modules of other four categories to this generic task module. The adapter contained one conv-layer with M filters and a ReLU layer. Each filter was a $1 \times 1 \times M$ tensor, which made its input and output feature maps have the same size. We also used data augmentation techniques in Section 3.3 to generate G'_y . Following Equation (4b), we used the $f'_{\text{dummy}} = f_{\text{dummy}}^{\text{1st}}$ operation to compute gradients of ReLU layers in both the task module and the adapter². We set $\lambda = 1.0/\mathbb{E}_T[D'_T]$ for all categories in this experiment.

Table 4 compares pixel-level segmentation accuracy between our method and the OL baseline when we learned transplant networks for object segmentation. We tested our method with a few (10–100) training samples. Compared to the OL baseline, our method exhibited 10%–12% higher accuracy.

⁴Because the dog category contained more training samples, the CNN for the dog category was believed to be better learned. Thus, we used the task module learned using dog images as a generic task module.

5 Conclusions and discussion

In this paper, we focused on a new task, *i.e.* merging pre-trained specific networks into a generic, distributed transplant network in a weakly-/un-supervised manner. We discussed the importance and core technical challenges of this task and developed the back-distillation algorithm as a theoretical solution to the challenging space-projection problem.

The back-distillation strategy significantly decreases the demand for training samples. In experiments, we learned the adapter to connect a pre-trained task module and a generic task module with a few or even without training samples. Experimental results demonstrated the superior efficiency of our method in scenarios of weakly-supervised learning and unsupervised learning. In particular, our method without any training samples even significantly outperformed the OL baseline with 100 training samples in Table 1.

When all adapters have been learned, we may fine-tune a task module using training samples of multiple categories. Note that our back-distillation loss is agnostic with regard to training samples, while the fine-tuning performance is sensitive to the number of training samples. Thus, with a few training samples, whether an additional fine-tuning operation will increase or decrease the generality of the transplant network is a difficult question, which requires much more sophisticated analysis and experiments to answer in the future.

References

- [1] M. Andrychowicz, M. Denil, S. G. Colmenarejo, M. W. Hoffman, D. Pfau, T. Schaul, B. Shillingford, and N. de Freitas. Learning to learn by gradient descent by gradient descent. *In NIPS*, 2016.
- [2] R. Anil, G. Pereyra, A. Passos, R. Ormandi, G. E. Dahl, and G. E. Hinton. Large scale distributed neural network training through online distillation. *In ICLR*, 2018.
- [3] X. Chen, Y. Duan, R. Houthoofd, J. Schulman, I. Sutskever, and P. Abbeel. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. *In NIPS*, 2016.
- [4] X. Chen, R. Mottaghi, X. Liu, S. Fidler, R. Urtasun, and A. Yuille. Detect what you can: Detecting and representing objects using holistic models and body parts. *In CVPR*, 2014.
- [5] Y. Chen, M. W. Hoffman, S. G. Colmenarejo, M. Denil, T. P. Lillicrap, M. Botvinick, and N. de Freitas. Learning to learn without gradient descent by gradient descent. *In ICML*, 2017.
- [6] N. Frosst and G. Hinton. Distilling a neural network into a soft decision tree. *In arXiv:1711.09784*, 2017.
- [7] B. Hariharan, P. Arbelaez, L. Bourdev, S. Maji, and J. Malik. Semantic contours from inverse detectors. *In ICCV*, 2011.
- [8] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *In CVPR*, 2016.
- [9] I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, and A. Lerchner. β -vae: learning basic visual concepts with a constrained variational framework. *In ICLR*, 2017.
- [10] G. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. *In NIPS Workshop*, 2014.
- [11] G. Huang, Z. Liu, K. Q. Weinberger, and L. van der Maaten. Densely connected convolutional networks. *In CVPR*, 2017.
- [12] P. Kotschieder, M. Fiterau, A. Criminisi, and S. R. Bulò. Deep neural decision forests. *In ICCV*, 2015.
- [13] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. *In NIPS*, 2012.
- [14] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *In Proceedings of the IEEE*, 1998.
- [15] H. Li, Z. Lin, J. Brandt, X. Shen, and G. Hua. A convolutional neural network cascade for face detection. *In CVPR*, 2015.
- [16] K. Li and J. Malik. Learning to optimize. *In arXiv:1606.01885*, 2016.
- [17] C. Liu, B. Zoph, J. Shlens, W. Hua, L.-J. Li, L. Fei-Fei, A. Yuille, J. Huang, and K. Murphy. Progressive neural architecture search. *In arXiv:1712.00559*, 2017.
- [18] J. Long, E. Shelhamer, and T. Darrel. Fully convolutional networks for semantic segmentation. *In CVPR*, 2015.
- [19] S. Sabour, N. Frosst, and G. E. Hinton. Dynamic routing between capsules. *In NIPS*, 2017.
- [20] N. Shazeer, A. Mirhoseini, K. Maziarz, A. Davis, Q. Le, G. Hinton, and J. Dean. outrageously large neural networks: the sparsely-gated mixture-of-experts layer. *In ICLR*, 2017.
- [21] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *In ICLR*, 2015.
- [22] J. Wang, Z. Kurth-Nelson, D. Tirumala, H. Soyer, J. Leibo, R. Munos, C. Blundell, D. Kumaran, and M. Botvinick. Learning to reinforcement learn. *In arXiv:1611.05763v3*, 2017.
- [23] N. Wolchover. New theory cracks open the black box of deep learning. *In Quanta Magazine*, 2017.
- [24] M. Wu, M. C. Hughes, S. Parbhoo, M. Zazzi, V. Roth, and F. Doshi-Velez. Beyond sparsity: Tree regularization of deep models for interpretability. *In NIPS TIML Workshop*, 2017.
- [25] T. Wu, X. Li, X. Song, W. Sun, L. Dong, and B. Li. Interpretable r-cnn. *In arXiv:1711.05226*, 2017.

- [26] S. Zagoruyko and N. Komodakis. Paying more attention to attention: improving the performance of convolutional neural networks via attention transfer. *In arXiv:1612.03928*, 2017.
- [27] Q. Zhang, Y. N. Wu, and S.-C. Zhu. Interpretable convolutional neural networks. *In CVPR*, 2018.
- [28] Z. Zhong, J. Yan, and C.-L. Liu. Practical network blocks design with q-learning. *In AAAI*, 2018.
- [29] Z.-H. Zhou and J. Feng. Deep forest: Towards an alternative to deep neural networks. *In IJCAI*, 2017.
- [30] B. Zoph and Q. V. Le. Neural architecture search with reinforcement learning. *In ICLR*, 2017.
- [31] B. Zoph, V. Vasudevan, J. Shlens, and Q. V. Le. Learning transferable architectures for scalable image recognition. *In arXiv:1707.07012*, 2017.

Appendix: computation of gradients *w.r.t.* the distillation loss

In order to optimize the distillation loss, we need to compute $\frac{\partial D'_S(\theta_h)}{\partial \theta_h}$, *i.e.* $\frac{\partial D'(\theta)}{\partial \theta}$ with respect to the following $D'(\theta)$.

$$D'(\theta) \stackrel{\text{def}}{=} G_{y'} \frac{\partial y'}{\partial x^{(n)}} \cdots \frac{\partial x^{(m+1)}}{\partial x^{(m)}} \Big|_{x=x^{(m)}} = f'_{\text{conv}} \circ f'_{\text{dummy}} \circ f'_{\text{pool}}^{\text{avg}} \circ \cdots \circ f'_{\text{conv}}(G_{y'})$$

Thus, we first explore the close-form formulation of the function $D'(\theta)$. As shown in the above equation, we can transform the back-propagation process for computing $D'(\theta)$ as a number of cascaded functions $f'_{\text{conv}}, \dots, f'_{\text{pool}}^{\text{avg}}, f'_{\text{dummy}}, f'_{\text{conv}}$, which are derivatives of $f_{\text{conv}}, \dots, f_{\text{pool}}^{\text{avg}}, f_{\text{dummy}}, f_{\text{conv}}$. Since we have formulated f'_{dummy} in the manuscript and it is easy to obtain $f'_{\text{pool}}^{\text{avg}}$, we mainly focus on the formulation of f'_{conv} .

In general, it is not difficult to derive the derivative of any convolution operation. Here, we focus on the most common case, *i.e.* the convolution operation with a padding $\pm p$ and a stride of 1. Given a tensor $x \in \mathbb{R}^{M \times M \times D}$ and C convolutional filter with weights $w \in \mathbb{R}^{m \times m \times D \times C}$ and a bias term $b \in \mathbb{R}^C$, the convolution can be written as $y = x \otimes w + b$. For VGG networks, people usually set $m = 2p + 1$. w can further absorb b by adding the $(D + 1) - \text{th}$ channel to w and adding the $(D + 1) - \text{th}$ channel to x with $x_{:, :, D+1, :} = 1$. Thus, we can obtain

$$y = x \otimes w$$

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial y} \otimes W$$

where $W \in \mathbb{R}^{m \times m \times C \times D}$ and $W_{i,j,k,l} = w_{m+1-i, m+1-j, l, k}$. Thus, we can write the derivative of f'_{conv} as

$$f'_{\text{conv}}(G') = G' \otimes W$$

In this way, we obtain the close-form formulation of the function $D'(\theta)$. We can easily compute $\frac{\partial D'(\theta)}{\partial \theta}$ using the chain rule of back propagation. We can consider this process as a *back-back-propagation*.

About the case of using an enlarged $G_{y'}$

When we use an enlarged pseudo-gradient $G_{y'} \in [-1, +1]^{S \times S \times 1}$, we can obtain an enlarged gradient map D'_S . As discussed above, the computation of D'_S can be considered as a number of cascaded functions $f'_{\text{conv}}, \dots, f'_{\text{pool}}^{\text{avg}}, f'_{\text{dummy}}, f'_{\text{conv}}$ with the input $G_{y'}$, which are quite similar to the forward propagation in the neural network.

Appendix: visualization of network structures used in three experiments

