

---

# Learning to Transfer Learn

---

**Linchao Zhu**  
Google Cloud AI  
Sunnyvale, CA  
linchaoz@google.com

**Sercan Ö. Arik**  
Google Cloud AI  
Sunnyvale, CA  
soarik@google.com

**Yi Yang**  
University of Sydney  
Sydney, Australia  
yi.yang@uts.edu.au

**Tomas Pfister**  
Google Cloud AI  
Sunnyvale, CA  
tpfister@google.com

## Abstract

We propose a novel framework, learning to transfer learn (L2TL), to improve transfer learning on a target dataset by judicious extraction of information from a source dataset. Our framework considers joint optimization of strongly-shared weights between models of source and target tasks, and employs adaptive weights for scaling of constituent loss terms. The adaptation of the weights is done using a reinforcement learning (RL)-based policy model, which is guided based on a performance metric on the target validation set. We demonstrate state-of-the-art performance of L2TL given fixed models, consistently outperforming fine-tuning baselines on various datasets. In addition, in the regimes of small-scale target datasets and significant label mismatch between source and target datasets, L2TL outperforms previous methods by a large margin.

## 1 Introduction

Deep neural networks excel at understanding images [He et al., 2016, Simonyan and Zisserman, 2015, Zagoruyko and Komodakis, 2016], text [Conneau et al., 2017, Devlin et al., 2018, Lai et al., 2015] and audio [van den Oord et al., 2016, Amodei et al., 2016, Chiu et al., 2018]. The performance of deep neural networks improves significantly with more training data [Hestness et al., 2017]. As the applications of deep neural networks diversify and extend to use cases with small training datasets, conventional training approaches are often insufficient to yield high performance. It becomes highly beneficial to utilize extra source datasets and “transfer” the relevant information to the target training dataset. Transfer learning, commonly in the form of obtaining a pre-trained model on a large-scale source dataset and then further training it on the target dataset (known as fine-tuning), has become the standard recipe for most real-world artificial intelligence applications. Compared to training from random initialization, fine-tuning yields considerable performance improvements and convergence speedup, as demonstrated for object recognition [Razavian et al., 2014], semantic segmentation [Long et al., 2015], language understanding [Devlin et al., 2018], speech synthesis [Arik et al., 2018], audio-visual recognition [Moon et al., 2014] and language translation [Zoph et al., 2016].

Towards the motivation of pushing the performance of transfer learning, recent studies [Ngiam et al., 2018, Mahajan et al., 2018, Lee et al., 2019] have explored the direction of matching the source and target dataset distributions. Even simple methods to encourage domain similarity, such as prior class distribution matching [Ngiam et al., 2018], are shown to be effective – indeed, in some cases, the distribution match can be much more important than the scale of the source dataset. In this paper, our goal is to push this direction further by introducing a novel RL-based meta learning framework. Our framework, named learning to transfer learn (L2TL), adaptively infers the domain similarity

directly from the performance obtained on the test objective of the target validation set. There are numerous cases that source dataset samples could have features that are highly relevant to the target dataset samples and would benefit the learning process, but they may belong to very different classes. For example, consider the visual classification problem for bird images. The source dataset may not contain bird images, but it may have airplane images that contain similar visual patterns that would aid the training of the bird classifier model. The L2TL framework is designed to automatically handle these cases with its policy training and can push the performance further in ways that manual source dataset selection or fixed domain similarity methods may not be able to. We demonstrate state-of-the-art transfer learning results given fixed models in wide range of datasets:

- *Source and target datasets from similar domains:* We consistently outperform the fine-tuning baseline with a 0.6%-1.3% relative accuracy gain on five fine-grained datasets, and consistently outperform DATL with a 0.3%-1.5% relative accuracy gain.
- *Low-shot target dataset regime:* Our method consistently outperforms the fine-tuning baseline on five fine-grained datasets, up to 6.5% relative accuracy gain for five samples per class.
- *Source and target datasets from dissimilar domains:* We outperform the fine-tuning baseline, up to 1.7% relative accuracy gain on a texture dataset and 0.7 relative AUC gain on a Chest X-Ray dataset [Irvin et al., 2019].

In addition to achieving state-of-the-art performance, our framework is capable of ranking the values of source data points according to their contributions to a target performance. Such capability was considered using indirect measures like domain similarity previously, but L2TL directly relates them to the target performance for the first time.

## 2 Related Work

**Adaptive transfer learning:** There is a long history of transfer learning for neural networks [Yosinski et al., 2014] [Girshick et al., 2014]. Various directions have been considered to improve standard fine-tuning. One direction is carefully choosing which portion of the network to adapt while optimizing the information extraction from the source dataset. In [Guo et al., 2019], a policy network is used to make routing decisions on whether to pass the input through the fine-tuned or the pre-trained layers. In [Li et al., 2018], a regularization scheme is proposed to promote the similarity of the fine-tuned model with the pre-trained model as a favorable inductive bias. Another exploration direction is to carefully decide which input samples are relevant to the target task, as in our paper. [Ge and Yu, 2017] uses filter bank responses to select nearest neighbor source examples and demonstrates improved performance. In [Cui et al., 2018], domain similarity between source and target datasets is quantified using Earth Mover’s Distance (EMD). Transfer learning is demonstrated to benefit from pre-training on a source domain that is similar by this measure. With a simple greedy subset creation selection criteria, promising results are shown for improving the target test set performance. Domain adaptive transfer learning (DATL) in [Ngiam et al., 2018] employs probabilistic shaping, where the value is proportional to the ratios of estimated label prior probabilities. Our method does not use an indirect similarity metric like proximity of filter bank responses, EMD or prior class probabilities. Instead, it proposes a framework with the goal of optimizing the target set metric directly.

**Reweighting training examples:** Reweighting of constituent training terms has been considered for various performance goals. [Ren et al., 2019] applies gradient descent-based meta learning with the goal of providing noise robustness and class balance in learning. Focal loss [Lin et al., 2017] is another soft weighting scheme that emphasizes harder examples. In [Jiang et al., 2018], a student-teacher training framework is utilized such that the teacher model provides a curriculum via a sample weighting scheme for the student model to focus on samples whose labels are likely to be correct. [Ghorbani and Zou, 2019] studies the value of examples by estimating the data Shapley value, and it shows that removing examples with low values would not harm performance. Reweighting of examples is also used in self-paced learning [Kumar et al., 2010] where the weights are optimized to learn easier examples first. Different from these, our paper focuses on optimizing the weights for each class, with the purpose of improving the transfer learning performance.

**Meta learning:** Meta learning generally refers to learning to learn frameworks [Schmidhuber et al., 1997] whose goal is to improve the adaptation to a new task with the information extracted from

other tasks. Meta learners are typically based on inspiration from known learning algorithms like gradient descent [Finn et al., 2017] or derived from black box neural networks [Santoro et al., 2016]. In few-shot learning, the use of validation loss as a meta-objective has been explored [Ravi and Larochelle, 2017]. However, for optimization problems with non-differentiable objectives like neural architecture search, RL-based meta learning is shown to be a promising approach [Zoph and Le, 2017, Pham et al., 2018]. RL-based optimization has successfully been applied in many search spaces, e.g. learning a data augmentation policy [Cubuk et al., 2018]. The specific form of meta learning application in L2TL is novel – it employs guidance on the source dataset information extraction with the reward from the target validation dataset performance. Different from few-shot learning, we consider a more common real-world scenario where an easily accessible source dataset is integrated.

### 3 Learning From Source and Target Datasets

We consider a training objective function  $\mathcal{L}(\Omega, \zeta_S, \zeta_T, \lambda, \alpha_s, \alpha_t)$ <sup>1</sup> that is jointly optimized for a source dataset  $D_S$  and a target dataset  $D_T$  in the general form:

$$\mathcal{L} = \alpha_s[i] \cdot \sum_{j=1}^{B_S} \lambda(x_j, y_j; \Phi) \cdot L_S(f_S(x_j; \Omega, \zeta_S), y_j) + \alpha_t[i] \cdot \sum_{k=1}^{B_T} L_T(f_T(x'_k; \Omega, \zeta_T), y'_k), \quad (1)$$

where  $(x, y)$  are the input and output pairs  $(x_j, y_j \sim D_S, x'_k, y'_k \sim D_T)$ ,  $B_S$  and  $B_T$  are the source and target batch sizes<sup>2</sup>,  $\alpha_s[i]$  and  $\alpha_t[i]$  are the scaling coefficients at  $i^{th}$  iteration,  $\lambda$  is the importance weighing function,  $f_S(\cdot; \Omega, \zeta_S)$  and  $f_T(\cdot; \Omega, \zeta_T)$  are encoding functions for the source and the target datasets with trainable parameters  $\Omega$ ,  $\zeta_S$  and  $\zeta_T$ <sup>3</sup>. To maximally benefit from the source dataset, a vast majority of the trainable parameters should be shared. If we consider the decompositions,  $f_S(\cdot; \Omega, \zeta_S) = h_S(\cdot; \zeta_S) \circ g(\cdot; \Omega)$  and  $f_T(\cdot; \Omega, \zeta_T) = h_T(\cdot; \zeta_T) \circ g(\cdot; \Omega)$ ,  $g$  can be a high capacity function with large number of trainable parameters that can be represented with a deep neural network, and  $h_T$  and  $h_S$  are low capacity functions with small number of parameters that can be represented with very shallow neural networks.<sup>4</sup>

The learning goal of Eq. 1 is to generalize to unseen samples from a held-out target validation dataset, and maximize an evaluation metric on it:

$$\sum_{x', y' \sim D'_T} R(f_T(x'; \hat{\Omega}, \hat{\zeta}_T), y')), \quad (2)$$

where  $R$  is the target performance metric (that may or may not be differentiable with respect to  $x$  and  $y$ ) such as the top-1 accuracy or area under the curve (AUC) for classification.  $\hat{\Omega}, \hat{\zeta}_T$  are the pre-trained weights optimized in Eq. 1.

Without transfer learning, i.e., when only the target dataset is considered,  $\alpha_s[i] = 0$  and  $\alpha_t[i] = 1$  for all  $i$ . In fine-tuning, the optimization is first considered for the source dataset for  $N_S$  steps with uniform weighing of the samples  $\lambda(x, y) = 1$ , and then for the target dataset using the pre-trained weights  $\hat{\Omega}, \hat{\zeta}_T$ , i.e.:

$$(\alpha_s[i], \alpha_t[i]) = \begin{cases} (1, 0), & i < N_S \\ (0, 1), & i > N_S \end{cases} \quad (3)$$

Next, we describe our framework towards optimal learning from source and target datasets.

### 4 Learning to Transfer Learn Framework

The learning to transfer learn (L2TL) framework (shown in Fig. 1) aims to learn the weight assignment adaptively, rather than using a fixed weight assignment function  $\lambda(x, y; \Phi)$  to measure the relatedness

<sup>1</sup>Function arguments are not often shown in the paper for notational convenience.

<sup>2</sup>Batch approximations may be optimal for different batch sizes for source and target dataset and thus may employ different batch normalization parametrization.

<sup>3</sup>In  $f(\cdot; \mathbf{W})$  representation,  $\mathbf{W}$  denote the trainable parameters.

<sup>4</sup>Source datasets are typically much larger and contain more classes, hence  $h_S$  may have higher number of parameters than  $h_T$ .

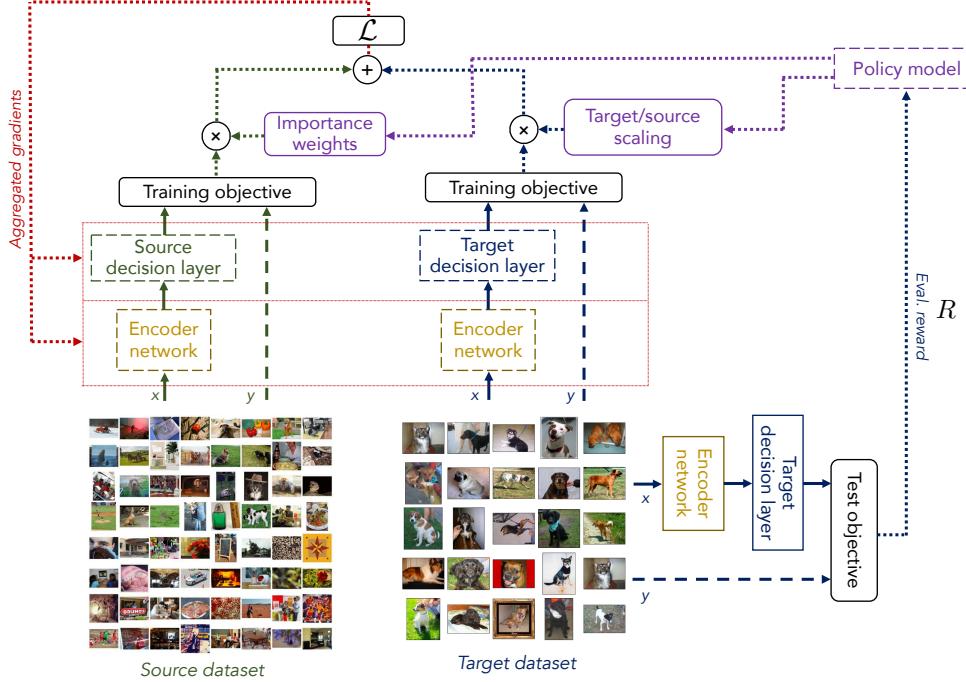


Figure 1: Overall diagram of the L2TL framework. Dashed boxes correspond to trainable functions.

between the source domain and the target domain. Learning of the adaptive weights in L2TL is guided by the performance on a held-out target validation dataset (quantified using the given metric  $R$ ). Thus, beyond targeting general relatedness, the framework directly targets relatedness for the specific goal of improvement in target evaluation performance.

While optimizing for  $\lambda(x, y; \Phi)$ , one straightforward option for the choice of scaling coefficients would be alternating them between  $(1, 0)$  and  $(0, 1)$  – i.e. training the source dataset until convergence with optimized  $\hat{\Phi}$  and then training the target dataset until convergence with the pre-trained weights from the source dataset. Yet, the approach may potentially require many alternating update steps and the computational cost may become prohibitively high. Instead, the policy model in L2TL is designed to output  $(\alpha_s[i], \alpha_t[i])$  along with  $\lambda$ . Without loss of generality, we can optimize a single weight  $\alpha_s[i]$  (setting  $\alpha_t[i] = 1$ ) as the optimization is scale invariant. The policy optimization step is decoupled from the gradient-descent based optimization for  $\Omega$ ,  $\zeta_S$  and  $\zeta_T$ . Updates are reflected to the policy model via the information embodied in  $\Omega$  and  $\zeta_T$ .

In the first phase of a learning iteration, we apply gradient decent-based optimization to learn the encoder weights  $\Omega$ , and the classifier layer weights  $\zeta_S$ ,  $\zeta_T$  to minimize the loss function  $\mathcal{L}$ :

$$\hat{\Omega}, \hat{\zeta}_S, \hat{\zeta}_T = \operatorname{argmin}_{\Omega, \zeta_S, \zeta_T} \mathcal{L}(\hat{\Phi}; \Omega, \zeta_S, \zeta_T). \quad (4)$$

In this learning phase, the policy model is fixed, and its actions are sampled to determine weights. The loss might be skewed when most of source dataset samples in a batch are unrelated, while some batches contain more related examples. To ease this problem, we sample a larger batch and dynamically select more related examples. At each iteration, we sample a training batch of size  $M_S \cdot B_S$ , and use the top  $B_S$  of them with the highest weights for training updates. This approach also yields computational benefits as the gradients would not be computed for most source dataset samples until convergence.

In the second phase of a learning iteration, the goal is to optimize policy weight  $\Phi$  that maximizes the evaluation metric  $R_{D'_T}$  on the target validation set with given encoder weights from the first phase:

$$\max_{\Phi} R_{D'_T}(\hat{\Omega}, \hat{\zeta}_S, \hat{\zeta}_T; \Phi). \quad (5)$$

We treat this step as an RL problem, such that the policy model outputs the actions for  $\lambda(x, y; \Phi)$  and  $\alpha$  towards optimization of a reward. In the most general form  $\lambda(x, y; \Phi)$  may yield a very high dimensionality for optimization of  $\Phi$ . For simplicity, we consider a sample-independent modeling of  $\lambda(x, y; \Phi)$ , similar to [Ngiam et al., 2018], i.e.,  $\lambda(x, y; \Phi) = \lambda(y; \Phi)$ . For more efficient optimization, we discretize the possible values of  $\lambda(y; \Phi)$  into pre-defined number of actions, in the range  $\lambda(y) \in [0, 1]$ . We define  $n$  actions, such that each action  $k \in [0, n - 1]$  corresponds to a weight value  $k/(n - 1)$ . For example, when  $n = 11$ , the weight values are  $[0, 0.1, 0.2, \dots, 1.0]$ . We also discretize the possible values for  $\alpha$ , using  $n'$  actions. Each action  $k'$  corresponds to  $\beta k'/(n' - 1)$ , where  $\beta$  is a hyperparameter to constrain the value range of  $\alpha$ . The search space has  $n' \times (c_S)^n$  possibilities, where  $c_S$  is the number of classes in the source dataset. When training the policy model, we use policy gradient to maximize the reward on the target dataset  $D_{T'}$ , using a batch size of  $B_P$ . We use the evaluation mode for the convolutional encoder. At update iteration  $t$ , we denote the advantage  $A_t = R_t - b_t$ , where  $b_t$  is the baseline. Following [Pham et al., 2018], we use the moving average baseline to reduce variance, i.e.,  $b_t = (1 - \gamma)b_t + \gamma R_t$ , where  $\gamma$  is the decay rate. The policy gradient is computed by REINFORCE [Williams, 1992] and optimized using Adam optimizer [Kingma and Ba, 2014].

## 5 Experiments

We demonstrate the performance of L2TL in various scenarios. As the source dataset, we use the ImageNet dataset [Russakovsky et al., 2015] containing 1.28M images, and 1K classes, and a much larger JFT-300M dataset (containing  $\sim$ 300M images that are labeled with 18,291 classes) [Sun et al., 2017]. For JFT-300M, the images are automatically annotated using user feedback from web pages. Target datasets are chosen based on the scenarios. Hyperparameters of the encoder models are chosen from the published baselines and policy model parameters are cross-validated on a validation set. All hyperparameters are presented in Supplementary Material A. The policy model does not take inputs and it is parameterized by a random variable. For datasets that the testing accuracy is reported using the model trained on training and validation examples, L2TL is first trained on the training split using the reward from the validation set. Then, the learned control variables are then used to train the joint model on the training and validation examples. We do not use the test set during model training. For the fine-tune experiment, we use best set of hyperparameters evaluated on the validation set. The results are averaged using three runs.

### 5.1 Source and target datasets from similar domains

We initially consider the scenario of target datasets with classes that mostly exist in the source dataset. In this scenario, DATL [Ngiam et al., 2018] was proven to be effective. We evaluate L2TL on five datasets focusing on different subsets. The detailed dataset splits are presented in Table 1. The reward is measured on the validation set using top-1 accuracy.

Table 1: Details for five fine-grained datasets: Birdsnap (Birds) [Berg et al., 2014], Oxford-IIIT Pets (Pets) [Parkhi et al., 2012], Stanford Cars (Cars) [Krause et al.], FGVC Aircraft (Aircraft) [Maji et al., 2013], and Food-101 (Food) [Bossard et al., 2014].

	Birds	Pets	Cars	Aircraft	Food
# of classes	500	37	196	100	101
# of train examples	42,405	2,940	6,494	3,334	68,175
# of valid examples	4,981	740	1,650	3,333	7,575
# of test examples	2,443	3,669	8,041	3,333	25,250

The results of L2TL along with fine-tuning and DATL benchmarks are shown in Table 2. As can be seen, L2TL outperforms fine-tuning across all the datasets with 0.6%-1.3% relative accuracy difference, which demonstrates the strength of L2TL in selecting related source examples across various domains. DATL performs worse than fine-tuning on Birdsnap and Aircraft, unlike L2TL. This underlines the importance of leveraging the visual similarity in the ways beyond the label match as in DATL. We also evaluate results when using JFT-300M. The results are shown in Table 3. L2TL preserves its benefits in learning relatedness for much larger scale datasets.

Table 2: Transfer learning performance with ImageNet source dataset, compared to MixDCNN [Wang et al., 2015], EMD [Cui et al., 2018], OPAM [Peng et al., 2017] and DATL [Ngiam et al., 2018] (\* denotes the results fine-tuning baseline from [Ngiam et al., 2018]).

Method	Target dataset test accuracy (%)				
	Birds	Pets	Cars	Aircraft	Food
MixDCNN	74.8			82.5	
EMD	-	-	91.3	85.5	88.7
OPAM	-	93.8	92.2	-	-
Fine-tuning*	77.2	93.3	91.5	88.8	88.7
DATL *	76.6	94.1	92.1	87.8	88.9
Fine-tuning L2TL	77.1	93.1	92.0	88.2	88.4
L2TL	<b>78.1</b>	<b>94.4</b>	<b>92.6</b>	<b>89.1</b>	<b>89.2</b>

Table 3: Results on the test set for Birdsnap using images from JFT-300M.

Method	Accuracy
Fine-tuning	74.9
DATL	81.7
L2TL	<b>82.4</b>

Table 4: Top chosen classes from ImageNet source dataset that are related to the target datasets.

Target dataset	Source classes (weights)
Birds	bee eater (1.0); aepycoerus melampus (0.95); sea cradle (0.92); barracouta (0.91) valley (0.90); sombrero (0.86); rosehip (0.84); Scottish deerhound (0.82); black swan (0.77); bell pepper (0.77); coyote (0.77); American robin (0.76);
Pets	coral reef (0.88); prayer rug (0.88); koala (0.84); fire salamander (0.81) Irish setter (0.79); Arabian camel (0.78); Irish terrier (0.74); leaf beetle (0.72); Rottweiler (0.71); giant schnauzer (0.70); English setter (0.69)
Cars	desktop computer (0.80); butternut squash (0.76); barrel, cask (0.65); weevil (0.60); pool table (0.56); clumber (0.54); passenger car (0.50); race car (0.49); washer (0.46); terrapin (0.45); seat belt (0.32); fire truck (0.29);
Aircraft	bagel, beigel (0.97); ballplayer, baseball player (0.84); freight car (0.80); teapot (0.83); crate (0.78); velvet (0.74); electric locomotive (0.68); pirate, pirate ship (0.45); amphibious vehicle (0.28); airliner (0.10);
Food	caldron, cauldron (0.84); menu (0.81); seashore, coast, seacoast (0.74) acorn squash (0.73); dining table, board (0.67); globe artichoke (0.67); mushroom (0.65); chocolate sauce, chocolate syrup (0.62); plate (0.54);

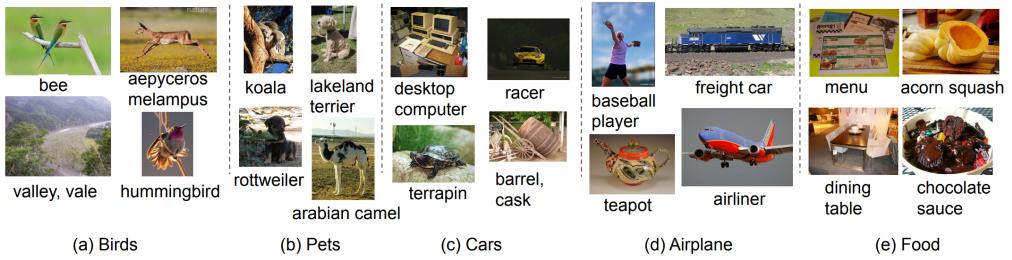


Figure 2: Representative examples from the source datasets that are related to the target dataset.

**Visualization:** To better understand the quality of the learned weights, we sample  $10k$  actions from the policy and rank the source labels according to their weights. The classes with high weights are shown in Table 4. For Birdsnap, we can see that the most related class is “bea eater” which is one of the bird species in ImageNet. “Aepycoerus melampus” is an antelope that has narrow mouth, which is similar to some birds with sharp spout. The “valley” class matches the background in some examples in Birdsnap. These results show that L2TL can extract related classes from the source based on the major object as well as background scenes. For Stanford Cars, we surprisingly observe that “barrel, cask” has high probability to have high weight, which indeed includes wheels and car-looking body types in many images. “terrapin” is a reptile that crawls on the ground with four legs, whose shape looks like vehicles in some way. For Food-101, we can see the top classes are the least relevant intuitively but still contain classes with visually-relevant patterns. For example, “caldron, cauldron”

may contain images with food inside, but most “seashore, coast” are less related to food. We visualize a few representative examples in Figure 2.

### 5.1.1 Low-shot examples regime

In this section, we study the performance changes when we have a smaller number of training examples. In the extreme low-shot regime, generalizing to unseen examples is most challenging as the model can be prone to overfitting. Fig. 3 shows performance in the small data regime for the same datasets. In most cases, we observe significant increase in performance when the number of examples per class is smaller. For five examples per class, the gap is as high as 6.5% (for Stanford Car). We observe that the gap between the L2TL and the fine-tune baseline often becomes smaller when more examples are used, but still remains as high as 1.5% when 60 examples per class are used (for Birdsnap). The results show that our L2TL framework can yield significant improvements in real-world tasks where the number of training examples are limited.

### 5.1.2 Ablation studies on the effectiveness of RL

In section, we study a few elements that could affect the performance of L2TL framework. In all plots, we show the classification accuracy during the whole training process. The curves are oscillatory at the beginning, but become stable later during the training.

**Learning to assign weights:** We study the effectiveness of RL training comparing it to two baselines: (i) random search: where the policy model is not optimized and random actions are chosen as the policy output, (ii) uniform weights: a constant importance weight is assigned to all training samples. Note in these baselines,  $\alpha$  is still optimized via policy gradient. The results are shown in Fig. 4. For all datasets, we observe that L2TL clearly outperforms both baselines after sufficient number of iterations, underlining the importance of reweighting policy learning via policy gradient.

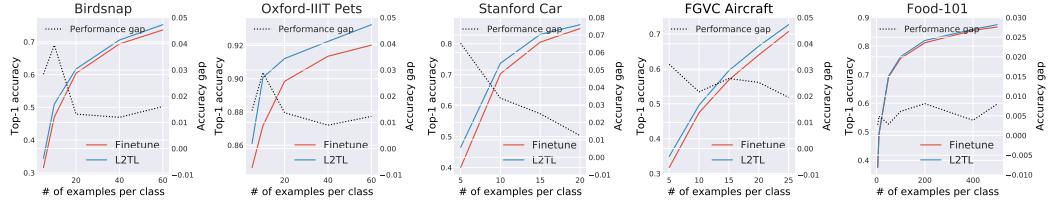


Figure 3: Number of examples per class vs. top-1 accuracy for L2TL and fine-tuning

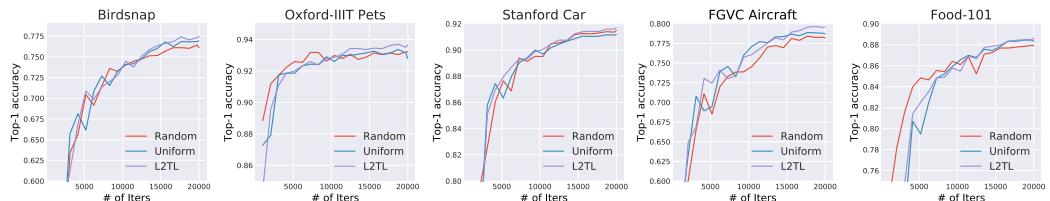


Figure 4: Performance comparison between L2TL, random search and uniform weights.

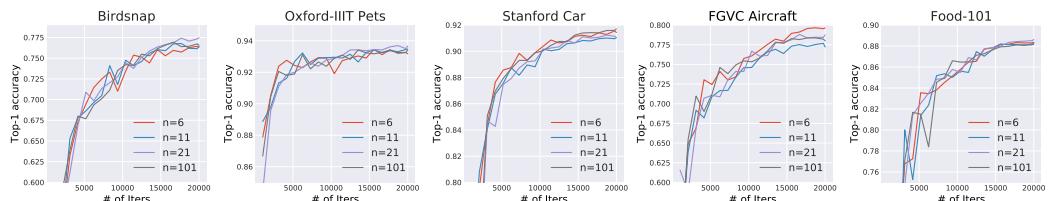


Figure 5: Performance comparisons when using different number of actions –  $n \in \{6, 11, 21, 101\}$ .

**The effect of the number of policy actions:** We study the impact of the number actions by tuning  $n$ . The results are shown in Fig. 5. As can be seen, sensitivity to  $n$  differs across datasets, potentially

due to the number of classes in the dataset and their inherent similarities. Food-101 has the least sensitivity to  $n$  – all the cases converge around 88.4%, whereas FGVC Aircraft dataset has the highest sensitivity. When  $n$  is larger, the difference of class weight outputs between possible actions is smaller. When  $n$  is smaller, there are fewer choices in the action space, and it can be easier to optimize the policy. However, the differences between the action values are large, which can fail to distinguish the weight differences between the classes. Overall, our experimental results suggest that the number of actions is an important hyperparameter for some datasets.

### 5.1.3 Computational cost of training

We present the computational cost for fine-tuning, DATL and L2TL in Table 5. In DATL, given a new target dataset, a new model has to be trained on the resampled data until convergence. This step is time-consuming for large-scale source datasets. In L2TL, the transfer learning step is more expensive than fine-tuning, as it requires the computation on both the source and the target datasets. Yet, it only requires a single training pass on the source dataset, and hence yields much lower training time compared to DATL.

Table 5: Computational cost of training using Inception-V3 on Cloud TPU v2 when the source dataset is ImageNet. The last column of total time assumes availability of a pre-trained source model, which is a common scenario.

Method	Number of iterations		Time per iterations		Total time	
	Pre-training	Transfer learning	Pre-training	Transfer learning	From scratch	With pre-trained source model
Fine-tuning	213,000	20,000	0.14s	0.21s	9.5h	1.2h
DATL	713,000	20,000	0.14s	0.21s	28.9h	20.6h
L2TL	213,000	20,000	0.14s	0.75s	12.5h	4.2h

## 5.2 Source and target datasets from dissimilar domains

In this section, we consider two target datasets with classes that do not exist in source dataset:

- (i) Describable Textures Dataset (DTD) dataset [Cimpoi et al., 2014]: It contains textural images in the wild from 47 classes such as striped and matted. The dataset has 20 splits and we evaluate the testing results on the first split. Each training, validation, and testing split has 1,880 images.
- (ii) Chest X-Ray Dataset CheXpert [Irvin et al., 2019]: This medical dataset has been recently introduced for chest radiograph interpretation. It consists of 224,316 chest radiographs of 65,240 patients labeled for 14 observations as positive, negative, or uncertain. Following [Irvin et al., 2019], we report AUC on five classes and we regard “uncertain” examples as positive.

The results for DTD are shown in Table 6. We observe that ImageNet fine-tuning greatly improves the classification results. L2TL further improves the fine-tuning baseline by 1.5%, demonstrating the significance of L2TL enabling the use of related source classes instead of using all classes. For the low-shot example case, we observe a more than 5% gain when only 10 examples are provided per class, validating the effectiveness of L2TL in more accurate transfer learning. We show the top related classes from ImageNet in Fig. 6. These results demonstrates that L2TL is able to utilize visual similarities between the source and the target classes. These similarities occur in the form of texture pattern for most DTD classes, for example, “prairie chicken” images from ImageNet typically contain patterns very relevant to “lined” from DTD.

Table 6: Results on the test set for DTD on split 1.

Method	Acc
Random Init	57.4
Fine-tuning	70.3
L2TL	<b>72.0</b>
10-shot, Fine-tuning	55.0
10-shot, L2TL	<b>60.1</b>



Figure 6: Top four related classes from ImageNet while transferring to DTD target. Representative images from each ImageNet class are shown along with related examples from DTD.

On the CheXpert dataset, we use the mean AUC as the reward, and the convolutional architecture is DenseNet-121 [Huang et al., 2017]. The results are shown in Table 7. We use the same evaluation protocol as in [Irvin et al., 2019]. We reproduce the fine-tuning baseline in [Irvin et al., 2019].<sup>5</sup> L2TL performs better than the fine-tuning baseline by 0.7%. There are not many straightforward visual similarities to humans to identify from ImageNet to relate to CheXpert, but L2TL is still capable of discovering hidden patterns in the images.

Table 7: AUC comparisons on the CheXpert dataset. \* denotes the results reported in [Irvin et al., 2019]. “Effusion” denotes Pleural Effusion.

	Atelectasis	Cardiomegaly	Consolidation	Edema	Effusion	Mean
Fine-tuning *	85.8	83.2	89.9	94.1	93.4	89.3
Fine-tuning	85.2	83.8	90.0	94.5	92.8	89.3
L2TL	<b>86.1</b>	<b>84.4</b>	<b>91.5</b>	<b>94.8</b>	93.2	<b>90.0</b>

## 6 Conclusions and Future Work

We propose a novel RL-based framework, L2TL, to improve transfer learning on a target dataset by careful extraction of information from a source dataset. L2TL considers joint optimization of models for source and target tasks, while using adaptive weights for scaling of constituent loss terms. We use the performance metric on the target validation set as the reward to train the policy model, which outputs the weights for each source class adaptively. We demonstrate state-of-the-art performance of L2TL for various datasets. The performance benefit of L2TL typically gets more significant as the target dataset size gets smaller. Our framework does not utilize an explicit similarity metric, but learns source class weights to directly optimize the target dataset performance. In cases where source and target datasets come from substantially different domains, L2TL still yields clear improvements. This improvement often comes from utilizing the source dataset classes that have relevant visual patterns despite belonging to a substantially different class.

Our general L2TL framework can be pushed beyond the approximations in this paper. A search space with a higher optimization granularity is expected to improve the results. Modeling the  $x$  dependence of the  $\lambda(x, y; \Phi)$  function using a policy gradient can be a promising step towards that approach (particularly for the datasets with high intra-class variance, such a model can help to select particular samples within the class) albeit likely accompanied by increased computational complexity.

## Acknowledgement

Discussions with Zizhao Zhang, Jiquan Ngiam, Boqing Gong, Wei Wei are gratefully acknowledged.

---

<sup>5</sup>Our reproduced results are matched with [Irvin et al., 2019] on mean AUC. However, there are variances as we can see that for some classes, we achieve slightly worse than [Irvin et al., 2019]. This may because of the small number of validation examples (200) used.

## References

- D. Amodei, R. Anubhai, E. Battenberg, C. Case, J. Casper, et al. Deep speech 2: End-to-end speech recognition in english and mandarin. In *ICML*, 2016.
- S. Ö. Arik, J. Chen, K. Peng, W. Ping, and Y. Zhou. Neural voice cloning with a few samples. In *NIPS*, 2018.
- T. Berg, J. Liu, S. Woo Lee, M. L. Alexander, D. W. Jacobs, and P. N. Belhumeur. Birdsnap: Large-scale fine-grained visual categorization of birds. In *CVPR*, 2014.
- L. Bossard, M. Guillaumin, and L. Van Gool. Food-101–mining discriminative components with random forests. In *ECCV*, 2014.
- C. Chiu, T. N. Sainath, Y. Wu, R. Prabhavalkar, P. Nguyen, et al. State-of-the-art speech recognition with sequence-to-sequence models. In *ICASSP*, 2018.
- M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, , and A. Vedaldi. Describing textures in the wild. In *CVPR*, 2014.
- A. Conneau, H. Schwenk, L. Barrault, and Y. LeCun. Very deep convolutional networks for natural language processing. In *ACL*, 2017.
- E. D. Cubuk, B. Zoph, D. Mane, V. Vasudevan, and Q. V. Le. Autoaugment: Learning augmentation policies from data. *arXiv preprint arXiv:1805.09501*, 2018.
- Y. Cui, Y. Song, C. Sun, A. Howard, and S. Belongie. Large scale fine-grained categorization and domain-specific transfer learning. In *CVPR*, 2018.
- J. Devlin, M. Chang, K. Lee, and K. Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. *arxiv:1810.04805*, 2018.
- C. Finn, P. Abbeel, and S. Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *ICML*, 2017.
- W. Ge and Y. Yu. Borrowing treasures from the wealthy: Deep transfer learning through selective joint fine-tuning. In *CVPR*, 2017.
- A. Ghorbani and J. Zou. Data shapley: Equitable valuation of data for machine learning. In *ICML*, 2019.
- R. B. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *CVPR*, 2014.
- Y. Guo, H. Shi, A. Kumar, K. Grauman, T. Rosing, and R. S. Feris. Spottune: Transfer learning through adaptive fine-tuning. In *CVPR*, 2019.
- K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- J. Hestness, S. Narang, N. Ardalani, G. F. Diamos, H. Jun, H. Kianinejad, M. M. A. Patwary, Y. Yang, and Y. Zhou. Deep learning scaling is predictable, empirically. *arXiv:1712.00409*, 2017.
- G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger. Densely connected convolutional networks. In *CVPR*, 2017.
- J. Irvin, P. Rajpurkar, M. Ko, Y. Yu, S. Ciurea-Ilcus, C. Chute, H. Marklund, B. Haghgoo, R. Ball, K. Shpanskaya, et al. CheXpert: A large chest radiograph dataset with uncertainty labels and expert comparison. In *AAAI*, 2019.
- L. Jiang, Z. Zhou, T. Leung, L. Li, and L. Fei-Fei. Mentornet: Regularizing very deep neural networks on corrupted labels. In *ICML*, 2018.
- D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *ICLR*, 2014.
- J. Krause, J. Deng, M. Stark, and L. Fei-Fei. Collecting a large-scale dataset of fine-grained cars. *The Second Workshop on Fine-Grained Visual Categorization*, 2013.
- M. P. Kumar, B. Packer, and D. Koller. Self-paced learning for latent variable models. In *NIPS*, 2010.
- S. Lai, L. Xu, K. Liu, and J. Zhao. Recurrent convolutional neural networks for text classification. In *AAAI*, 2015.

- J. Lee, W. Yoon, S. Kim, D. Kim, S. Kim, C. H. So, and J. Kang. Biobert: a pre-trained biomedical language representation model for biomedical text mining. *arXiv:1901.08746*, 2019.
- X. Li, Y. Grandvalet, and F. Davoine. Explicit inductive bias for transfer learning with convolutional networks. In *ICML*, 2018.
- T. Lin, P. Goyal, R. B. Girshick, K. He, and P. Dollár. Focal loss for dense object detection. In *ICCV*, 2017.
- J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015.
- D. Mahajan, R. B. Girshick, V. Ramanathan, K. He, M. Paluri, Y. Li, A. Bharambe, and L. van der Maaten. Exploring the limits of weakly supervised pretraining. In *ECCV*, 2018.
- S. Maji, E. Rahtu, J. Kannala, M. Blaschko, and A. Vedaldi. Fine-grained visual classification of aircraft. *arXiv:1306.5151*, 2013.
- S. Moon, S. Kim, and H. Wang. Multimodal transfer deep learning for audio visual recognition. *arXiv:1412.3121*, 2014.
- J. Ngiam, D. Peng, V. Vasudevan, S. Kornblith, Q. V. Le, and R. Pang. Domain adaptive transfer learning with specialist models. *arXiv preprint arXiv:1811.07056*, 2018.
- O. M. Parkhi, A. Vedaldi, A. Zisserman, and C. Jawahar. Cats and dogs. In *CVPR*, 2012.
- Y. Peng, X. He, and J. Zhao. Object-part attention model for fine-grained image classification. *TIP*, 2017.
- H. Pham, M. Y. Guan, B. Zoph, Q. V. Le, and J. Dean. Efficient neural architecture search via parameter sharing. In *ICML*, 2018.
- S. Ravi and H. Larochelle. Optimization as a model for few-shot learning. In *ICLR*, 2017.
- A. S. Razavian, H. Azizpour, J. Sullivan, and S. Carlsson. CNN features off-the-shelf: an astounding baseline for recognition. In *CVPR*, 2014.
- M. Ren, W. Zeng, B. Yang, and R. Urtasun. Learning to reweight examples for robust deep learning. In *ICML*, 2019.
- O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. *IJCV*, 2015.
- A. Santoro, S. Bartunov, M. Botvinick, D. Wierstra, and T. Lillicrap. Meta-learning with memory-augmented neural networks. In *ICML*, 2016.
- J. Schmidhuber, J. Zhao, and M. Wiering. Shifting inductive bias with success-story algorithm, adaptive levin search, and incremental self-improvement. *Machine Learning*, 1997.
- K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.
- C. Sun, A. Shrivastava, S. Singh, and A. Gupta. Revisiting unreasonable effectiveness of data in deep learning era. In *ICCV*, 2017.
- A. van den Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, et al. Wavenet: A generative model for raw audio. *arXiv:1609.03499*, 2016.
- D. Wang, Z. Shen, J. Shao, W. Zhang, X. Xue, and Z. Zhang. Multiple granularity descriptors for fine-grained categorization. In *ICCV*, 2015.
- R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 1992.
- J. Yosinski, J. Clune, Y. Bengio, and H. Lipson. How transferable are features in deep neural networks? In *NIPS*, 2014.
- S. Zagoruyko and N. Komodakis. Wide residual networks. In *BMVC*, 2016.
- B. Zoph and Q. V. Le. Neural architecture search with reinforcement learning. In *ICLR*, 2017.
- B. Zoph, D. Yuret, J. May, and K. Knight. Transfer learning for low-resource neural machine translation. In *ACL*, 2016.

---

**Algorithm 1:** L2TL – Learning to Transfer Learn

---

```
N ← number of training iterations
for  $i \leftarrow 1$  to  $N$  do
|    $l_s \leftarrow 0, l_t \leftarrow 0$ 
|   for  $j \leftarrow 1$  to  $B_S$  do
|   |   Sample  $x_j, y_j$  from  $D_S$ 
|   |   Calculate classification loss  $L_S(x_j, y_j; \Omega, \zeta_S)$ 
|   |   Calculate example weight  $\lambda(x_j, y_j; \Phi)$ 
|   |    $l_s = l_s + \lambda \cdot L_S$ 
|   |    $l_s = \alpha_i \cdot l_s$ 
|   |   for  $k \leftarrow 1$  to  $B_T$  do
|   |   |   Sample  $x'_k, y'_k$  from  $D_T$ 
|   |   |   Calculate classification loss  $L_T(x'_k, y'_k; \Omega, \zeta_T)$ 
|   |   |    $l_t = l_t + L_T$ 
|   |   Update  $\Omega, \zeta_S, \zeta_T$  using stochastic gradient descent
|   |    $r \leftarrow 0$ 
|   |   for  $k \leftarrow 1$  to  $B_P$  do
|   |   |   Sample  $x'_k, y'_k$  from  $D_{T'}$ 
|   |   |   Calculate reward  $R(f_T(x'_k), y'_k)$ 
|   |   |    $r = r + R$ 
|   |   Update  $\Phi$  with reward  $r$  using policy gradient
```

---

## A Hyperparameters

When the source dataset is ImageNet, we use batch size  $B_S = 256$ ,  $B_T = 256$ ,  $B_P = 1024$  and a batch multiplier  $M_S = 5$  for all the experiments. For the JFT-300M dataset, to reduce the number of training iterations, we use  $B_S = 1,024$ . The number of actions  $n'$  for  $\alpha$  is 100.

We use the Inception-V3 architecture for all the experiments except CheXpert. For target dataset, we searched the initial learning rate from  $\{0.001, 0.005, 0.01, 0.05, 0.1, 0.15, 0.2, 0.4\}$ , weight decay from  $\{0, 4 \times 10^{-5}\}$ . All the datasets are optimized by SGD with momentum 0.9. We use the single central crop during evaluation. The learning rate is cosine decayed after first 2,000 iterations warmup. The number of training iterations is 20,000. When optimizing our policy model, we use the Adam optimizer with a fixed learning rate 0.0001. As policy model parameters, we set  $\beta = 0.5$  and  $\gamma = 0.05$  for all the experiments. We followed the standard image preprocessing procedure for Inception-V3 on both the source images and the target images.

For CheXpert, we use the DenseNet-121 architecture Huang et al. [2017] and following the evaluation protocol specified in Irvin et al. [2019], where ten crops are used for evaluation and 30 checkpoints are ensemble to obtain the final results. We cross validate weight decay and initial learning rate, where the weight decay is searched in  $[0, 0.0001]$  and the learning rate searched in range  $[0.5, 0.8, 1.0, 1.3, 1.5, 2.0]$ . All other hyperparameters are same as above. We use the same input preprocessing as described in <https://github.com/zoogzog/chexnet>.

## B Algorithm

We illustrate the training updates steps of L2TL in Algorithm. 1.