

Learning to Transfer Learn: Reinforcement Learning-Based Selection for Adaptive Transfer Learning

Linchao Zhu^{1,2}, Sercan Ö. Arık¹, Yi Yang², and Tomas Pfister¹

¹ Google Cloud AI, Sunnyvale, CA

² University of Technology Sydney, Sydney, Australia
{soarik,tpfister}@google.com; {linchao.zhu,yi.yang}@uts.edu.au

Abstract. We propose a novel adaptive transfer learning framework, learning to transfer learn (L2TL), to improve performance on a target dataset by careful extraction of the related information from a source dataset. Our framework considers cooperative optimization of shared weights between models for source and target tasks, and adjusts the constituent loss weights adaptively. The adaptation of the weights is based on a reinforcement learning (RL) selection policy, guided with a performance metric on the target validation set. We demonstrate that L2TL outperforms fine-tuning baselines and other adaptive transfer learning methods on eight datasets. In the regimes of small-scale target datasets and significant label mismatch between source and target datasets, L2TL shows particularly large benefits.

Keywords: Transfer Learning, Visual Understanding, Reinforcement Learning

1 Introduction

Deep neural networks excel at understanding images [15, 47], text [8] and audio [36, 1]. The performance of deep neural networks improves significantly with more training data [16]. As the applications diversify and span 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 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 [41], semantic segmentation [30], language understanding [8], speech synthesis [2], audio-visual recognition [33] and language translation [53].

Towards the motivation of pushing the performance of transfer learning, recent studies [35, 31, 26, 29] have explored the direction of matching the source

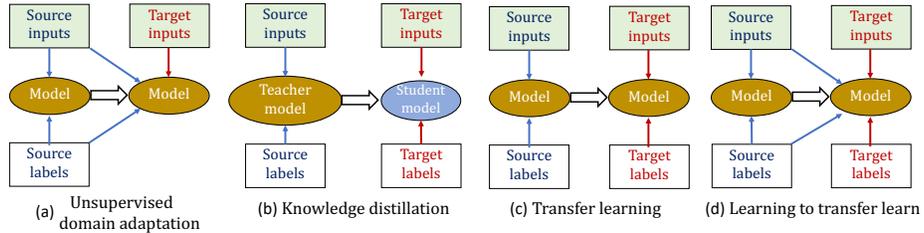


Fig. 1. L2TL and other adaptation settings. (a) Unsupervised domain adaptation incorporates source data and source labels for target domain adaptation, where the target labels are not provided. (b) Knowledge distillation aims to distill source knowledge from a teacher model to a student model, and the student model is usually more lightweight. (c) Conventional transfer learning transfers knowledge via weights of the model pre-trained on a source dataset, to obtain better performance on the target dataset. (d) Our L2TL adaptively infers the importance weights of the source examples based on the feedback from the target objective on the target validation set. With adaptive assignment, more relevant examples get higher weights to extract the information from the source dataset.

and target dataset distributions. Even simple methods to encourage domain similarity, such as prior class distribution matching in Domain Adaptive Transfer Learning (DATL) [35], are shown to be effective – indeed, in some cases, more important than the scale of the source dataset. Such adaptive transfer learning approaches, as in L2TL, typically assume the availability of the labeled source dataset for training on the labeled target dataset (that also differentiates the setting from unsupervised domain adaptation [13] or knowledge distillation [17], see Fig. 1), along with the pre-trained model. Given the increasing availability of very-large scale public datasets for various data types and the demand for cutting-edge deep learning on highly-specialized target tasks with small training datasets, this setting is indeed getting very common in practice [7, 10, 35].

In this paper, our goal is to push this direction further by introducing a novel *reinforcement learning (RL)-based framework*. Our framework, learning to transfer learn (L2TL), adaptively infers the beneficial source samples directly from the performance on the target task. There are cases that source samples could have features that are implicitly relevant to the target samples and would benefit the learning process, but they may belong to different classes. For example, consider the classification problem for bird images. The source dataset may not contain bird images, but may have airplane images with similar visual patterns that would aid the training of the bird classifier as they share similar visual patterns to learn valuable representations of the raw data. *L2TL framework is designed to automatically handle such cases with its policy learning, and can push the performance further in ways that manual source dataset selection or fixed domain similarity methods may not be able to.* L2TL considers cooperative optimization of models for source and target tasks, while using adaptive weights for scaling of constituent loss terms. L2TL leverages 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. Overall, L2TL does not utilize an explicit similarity metric as in [7, 35], but learns source class weights to directly optimize the target dataset performance.

We demonstrate promising transfer learning results given fixed models in a wide range of scenarios:

- *Source and target datasets from similar domains:* L2TL consistently outperforms the fine-tuning baseline with a 0.6%-1.3% relative accuracy gain on five fine-grained datasets, and DATL [35] with a 0.3%-1.5% relative accuracy gain. When the similarities become very apparent between some source and target classes, e.g. as in MNIST to SVHN digit recognition transfer case, the relative accuracy gain is 3.5% compared to fine-tuning baseline.
- *Large-scale source datasets:* When very large-scale source datasets are used, the selection of the relevant source classes become more important, the gain of L2TL is up to **7.5%**.
- *Low-shot target dataset regime:* L2TL significantly outperforms fine-tuning on fine-grained target datasets, up to **6.5%** accuracy gain with five samples per class.
- *Source and target datasets from dissimilar domains:* While other advanced transfer learning (based on explicit similarity measures) cannot be readily applied for this scenario, L2TL outperforms the fine-tuning baseline, up to **1.7%** accuracy gain on a texture dataset and **0.7 AUC** gain on Chest X-Ray dataset.

In addition, L2TL yields ranking of the source data samples according to their contributions to the target task, that can open horizons for new forms of interpretable insights.

2 Related Work

Adaptive transfer learning: There is a long history of transfer learning for neural networks, particularly in the form of fine-tuning [12]. Various directions were recently 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 [14], 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 [27], 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 direction is carefully choosing which input samples are relevant to the target task, as in our paper. [10] uses filter bank responses to select nearest neighbor source samples and demonstrates improved performance. In [7], domain similarity between source and target datasets is quantified using Earth Mover’s Distance (EMD). Transfer learning is shown to benefit from pre-training on a source domain that is similar in EMD. With a simple greedy subset creation selection criteria, promising results are shown for improving the target

test set performance. Domain adaptive transfer learning (DATL) [35] employs probabilistic shaping, where the value is proportional to the ratios of estimated label prior probabilities. L2TL does not use a similarity metric like proximity of filter bank responses, EMD or prior class probabilities. Instead, it aims to assign weights to optimize the target set metric directly.

Reweighting training examples: Reweighting of constituent training terms has been considered for various performance goals. [42] applies gradient descent-based meta-learning to update the weights of the input data, with the goal of providing more noise-robust and class-balanced learning. [20] formulates reweighting as a bilevel optimization problem, such that higher weights are encouraged for the training samples with more agreement of the gradients on the validation set. Focal loss [28] is another soft weighting scheme that emphasizes on harder examples. In [21], a student-teacher training framework is proposed 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. [11] studies the value of examples via Shapley values, and it shows that downweighting examples with low values might even improve performance. Reweighting of examples is also used in self-paced learning [24, 46] where the weights are optimized to learn easier examples first. In [50], an RL agent is used to adaptively sample relevant frames from videos. In this paper, unlike the above, we focus on transfer learning – L2TL formulates the transfer learning problem with a new loss function, including class-relevant weights and a dataset-relevant weights. L2TL learns the weight assignments with RL, in a setting where actions (source data selection) are guided with the rewards (target validation performance). Unlike gradient-descent based reweighting, RL-based rewarding is also applicable to scenarios where the target evaluation objective is non-differentiable, e.g., area under the curve (AUC).

Meta learning: Meta-learning broadly refers to learning to learn frameworks [45] whose goal is to improve the adaptation to a new task with the information extracted from other tasks. Meta learners are typically based on inspirations from known learning algorithms like gradient descent [9] or derived from black box neural networks [44]. As the notable meta learning application, in few-shot learning [9, 51], the use of validation loss as a meta-objective has been explored [40]. However, for optimization problems with non-differentiable objectives like neural architecture search, RL-based meta-learning is shown to be a promising approach [52, 39]. RL-based optimization has successfully been applied to other applications with enormously-large search spaces, e.g. learning a data augmentation policy [6]. The specific form of RL 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 many meta learning methods, e.g. those for few-shot learning, we consider a common real-world scenario where a very large-scale source dataset is integrated to extract information from. We do not employ any episodic training, hence L2TL is practically feasible to employ on very large-scale source datasets.

3 Learning From Source and Target Datasets

We consider a general-form training objective function $\mathcal{L}(\mathbf{\Omega}, \zeta_{\mathbf{S}}, \zeta_{\mathbf{T}}, \lambda, \alpha_s, \alpha_t)$ ³ jointly defined on a source dataset D_S and a target dataset D_T :

$$\begin{aligned} \mathcal{L} = & \alpha_s[i] \cdot \sum_{j=1}^{B_S} \lambda(x_j, y_j; \Phi) \cdot L_S(f_S(x_j; \mathbf{\Omega}, \zeta_{\mathbf{S}}), y_j) \\ & + \alpha_t[i] \cdot \sum_{k=1}^{B_T} L_T(f_T(x'_k; \mathbf{\Omega}, \zeta_{\mathbf{T}}), y'_k), \end{aligned} \quad (1)$$

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 sizes⁴, $\alpha_s[i]$ and $\alpha_t[i]$ are the scaling coefficients at i^{th} iteration, λ is the importance weighing function, $f_S(\cdot; \mathbf{\Omega}, \zeta_{\mathbf{S}})$ and $f_T(\cdot; \mathbf{\Omega}, \zeta_{\mathbf{T}})$ are encoding functions for the source and the target datasets with trainable parameters $\mathbf{\Omega}$, $\zeta_{\mathbf{S}}$ and $\zeta_{\mathbf{T}}$ ⁵.

To maximally benefit from the source dataset, a vast majority of the trainable parameters should be shared. If we consider the decompositions, $f_S(\cdot; \mathbf{\Omega}, \zeta_{\mathbf{S}}) = h_S(\cdot; \zeta_{\mathbf{S}}) \circ g(\cdot; \mathbf{\Omega})$ and $f_T(\cdot; \mathbf{\Omega}, \zeta_{\mathbf{T}}) = h_T(\cdot; \zeta_{\mathbf{T}}) \circ g(\cdot; \mathbf{\Omega})$, g shall be a high capacity encoder with large number of trainable parameters that can be represented with a deep neural network, and h_T and h_S are low capacity mapping functions with small number of parameters that can be represented with very shallow neural networks.⁶ The learning goal of Eq. 1 is generalizing to unseen target validation dataset, via maximization of the performance metric R :

$$\sum_{x', y' \sim D_T} R(f_T(x'; \hat{\mathbf{\Omega}}, \hat{\zeta}_{\mathbf{T}}), y'). \quad (2)$$

R does not have to be differentiable with respect to x and y and may include metrics like the top-1 accuracy or area under the curve (AUC) for classification. $\hat{\mathbf{\Omega}}$ and $\hat{\zeta}_{\mathbf{T}}$ are the pre-trained weights optimized in Eq. 1.

Without transfer learning, i.e., training with only target dataset, $\alpha_s[i] = 0$ 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 $\hat{\mathbf{\Omega}}, \hat{\zeta}_{\mathbf{T}}$, 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.

³ Function arguments are not often shown in the paper for notational convenience.

⁴ Batch approximations may be optimal for different batch sizes for source and target dataset and thus may employ different batch normalization parametrization.

⁵ In $f(\cdot; \mathbf{W})$ representation, \mathbf{W} denote the trainable parameters.

⁶ Source datasets are typically much larger and contain more classes, hence h_S may have higher number of parameters than h_T .

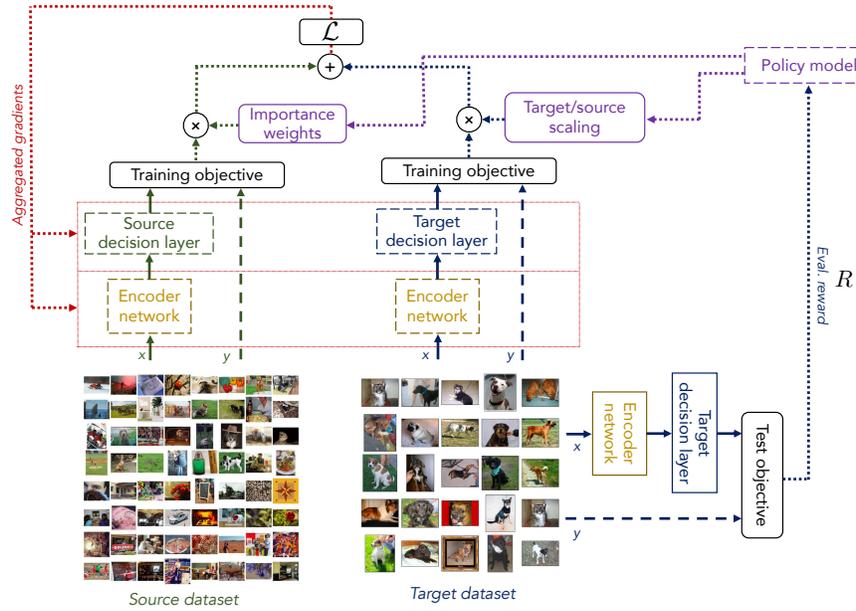


Fig. 2. Overall diagram of the L2TL framework. Dashed boxes correspond to trainable functions. L2TL employs a policy model to determine weights of the source dataset samples, to extract the information in a careful way to maximize the target dataset test objective. The models on source and target datasets are shared, via the encoder network.

4 Learning to Transfer Learn Framework

We propose learning to transfer learn (L2TL) framework (shown in Fig. 2) to learn the weight assignment adaptively, rather than using a fixed weight assignment function $\lambda(x, y; \Phi)$ to measure the relatedness between the source domain and the target domain. Learning of the adaptive weights in L2TL is guided by the performance metric R on a held-out target validation dataset. 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 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, we design the policy model in L2TL to output $(\alpha_s[i], \alpha_t[i])$ along with λ .⁷ The policy optimization step is decoupled from the gradient-descent based optimization for Ω, ζ_S

⁷ Without loss of generality, we can optimize a single weight $\alpha_s[i]$ (setting $\alpha_t[i] = 1$) as the optimization is scale invariant.

Algorithm 1: L2TL – Learning to Transfer Learn

```

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

```

and ζ_T . Updates are reflected to the policy model via the information embodied in $\mathbf{\Omega}$ and ζ_T . Algorithm. 1 overviews the training updates steps.

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

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

In this phase, the policy model is fixed, and its actions are sampled to determine weights. Although most batches would contain relevant source dataset samples, the loss might be skewed if most of source dataset samples in a batch are irrelevant (and would ideally get lower weights). To ease this problem, we use a larger batch size and dynamically select the most relevant 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, given encoder weights from the first phase, our goal is to optimize policy weights $\mathbf{\Phi}$ and maximize the evaluation metric $R_{D'_{T}}$ on the target validation set:

$$\max_{\mathbf{\Phi}} R_{D'_{T}}(\hat{\mathbf{\Omega}}, \hat{\zeta}_S, \hat{\zeta}_T; \mathbf{\Phi}). \quad (5)$$

D'_{T} is the held-out dataset to compute the reward. We treat this phase as an RL problem, such that the policy model outputs the action of value assignment for

$\lambda(x, y; \Phi)$ and α towards optimization of the reward (with the environment being training and evaluation setup, and the state being the encoder weights as the consequence of the first phase). In its general form $\lambda(x, y; \Phi)$ may yield a very high dimensionality for optimization of Φ . For simplicity and computational-efficiency, we consider sample-independent modeling of $\lambda(x, y; \Phi)$, similar to [35], i.e., $\lambda(x, y; \Phi) = \lambda(y; \Phi)$.⁸

For more efficient optimization via efficient systematic exploration of the very large action space, 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 the 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 α , using n' actions. Each action k' corresponds to $\beta k'/(n'-1)$, where β is a hyperparameter to constrain the value range of α . 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 of B_P samples. At iteration t , we denote the advantage $A_t = R_t - b_t$, where b_t is the baseline. Following [39], we use the moving average baseline to reduce variance, i.e., $b_t = (1 - \gamma)b_t + \gamma R_t$, where γ is the decay rate. The policy gradient is computed using REINFORCE [49] and optimized using Adam [22].

5 Experiments

5.1 Datasets and Implementation Details

We demonstrate the performance of L2TL in various scenarios. As the source dataset, we use the ImageNet dataset [43] containing 1.28M images from 1K classes, and also a much larger dataset, i.e., JFT-300M, containing ~ 300 M images from 18,291 classes to demonstrate the scalability of our approach. As the target datasets, we evaluate on **five** fine-grained image datasets (summarized in Table 1): **Birdsnap** [3], **Oxford-IIIT Pets** [37], **Stanford Cars** [23], **FGVC Aircraft** [32], and **Food-101** [4]. In addition, we also consider transfer learning scenario from MNIST to SVHN [34] to assess the effectiveness of L2TL for small-scale source datasets.

We also consider two target datasets with classes that do not exist in the source datasets: Describable Textures Dataset (DTD) [5] and Chest X-Ray Dataset CheXpert [19]. **Describable Textures Dataset:** DTD 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. **Chest X-Ray Dataset:** The CheXpert medical dataset is used for chest radiograph interpretation task. It consists of

⁸ A search space with a higher optimization granularity is expected to improve the results, albeit accompanied by significantly increased computational complexity for meta learning of x -dependent $\lambda(x, y; \Phi)$.

Table 1. Details of the five fine-grained datasets: Birdsnap (Birds) [3], Oxford-IIIT Pets (Pets) [37], Stanford Cars (Cars) [23], FGVC Aircraft (Air) [32], and Food-101 (Food) [4].

	Birdsnap	Oxford-IIIT Pets	Stanford Cars	Aircraft	Food-101
# 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

Table 2. Transfer learning performance with ImageNet source dataset. * indicates our implementation.

Method	Target dataset test accuracy (%)				
	Birdsnap	Oxford-IIIT Pets	Stanford Cars	Aircraft	Food-101
Fine-tuning [35]	77.2	93.3	91.5	88.8	88.7
Fine-tuning*	77.1	93.1	92.0	88.2	88.4
MixDCNN [48]	74.8	-	-	82.5	-
EMD [7]	-	-	91.3	85.5	88.7
OPAM [38]	-	93.8	92.2	-	-
DATL [35]	76.6	94.1	92.1	87.8	88.9
Our L2TL	78.1	94.4	92.6	89.1	89.2

224,316 chest radiographs of 65,240 patients labeled for 14 observations as positive, negative, or uncertain. Following [19], we report AUC on five classes and we regard “uncertain” examples as positive.⁹ For L2TL, we use the mean AUC as the reward.

Implementation Details. When the source dataset is ImageNet, we use a 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 α is 100.

We use Inception-V3 for all the experiments except CheXpert. For target dataset, we search the initial learning rate from $\{0.001, 0.005, 0.01, 0.05, 0.1, 0.15, 0.2, 0.4\}$, and weight decay from $\{0, 4 \times 10^{-5}\}$. All the datasets are optimized using SGD with a momentum of 0.9, trained for 20,000 iterations. We use the single central crop during evaluation. The learning rate is cosine decayed after first 2,000 iterations warmup. 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 $\gamma = 0.05$. We follow 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 [18] and follow the evaluation protocol specified in [19], where ten crops are used for evaluation and 30 checkpoints are ensembled to obtain the final results. We cross validate

⁹ Our reproduced results are matched with [19] on mean AUC. However, there are variances as we can see that for some classes, we achieve slightly worse than [19]. This may be because of the small number of validation examples (200) used.

Table 3. Results on Birdsnap using JFT-300M as the source dataset. The performance is reported on the test set.

Method	Birdsnap accuracy (%)
Fine-tuning	74.9
DATL [35]	81.7
Our L2TL	82.4

Table 4. Transfer learning from MNIST to SVHN. As shown in [29], fine-tuning shows gains over other training cases due to the inherent similarity between datasets. L2TL efficiently exploits this further by emphasizing on some MNIST classes more than others, and improves the transfer learning gains significantly.

Method	SVHN accuracy (%)
Random initialization	64.8
Fine-tuning	71.7
Our L2TL	75.2

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>.

Hyperparameters of the encoder models are chosen from the published baselines and the policy model parameters are cross-validated on a validation set. For datasets that the testing accuracy is reported using the model trained on training and validation samples, L2TL is first trained on the training set using the reward from the validation set. Then, the learned control variables are used to train the joint model on the combined set of training and validation samples – we completely exclude the test set during training. For the fine-tuning experiments, we use the best set of hyperparameters evaluated on the validation set. We present the results averaged over three runs. We observe that the standard deviation for the L2TL accuracy to be around 0.1%, much smaller than the gap between different methods.

5.2 Similar domain transfer learning

We initially consider the scenario of target datasets with classes that mostly exist in the source dataset.

Performance and comparison to other transfer learning methods. We first evaluate L2TL on five fine-grained datasets focusing on different subsets, with the reward of validation set top-1 accuracy. Table 2 shows the results of L2TL along with fine-tuning and state-of-the-art transfer learning benchmarks. With a well-optimized network architecture and learning rate scheduling, fine-tuning is already a solid baseline for the datasets in Table 1 [7, 35]. Yet, L2TL outperforms fine-tuning across all the datasets with 0.6%-1.3% accuracy difference, which demonstrates the strength of L2TL in selecting related source

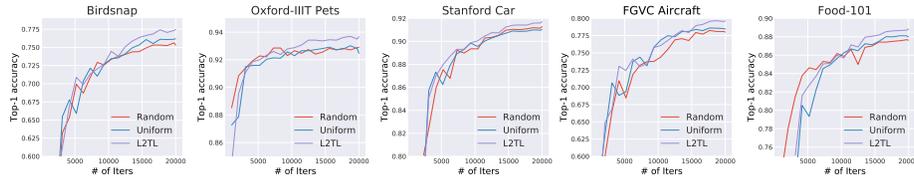


Fig. 3. Performance comparisons between L2TL, random search and uniform weights. The curves are oscillatory at the beginning, but become stable later during the training. L2TL outperforms the baselines when the training converges.

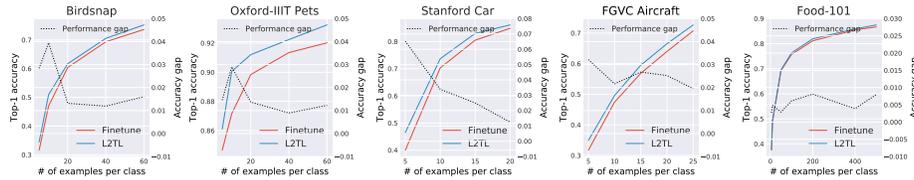


Fig. 4. Number of examples per class vs. top-1 accuracy for L2TL and fine-tuning.

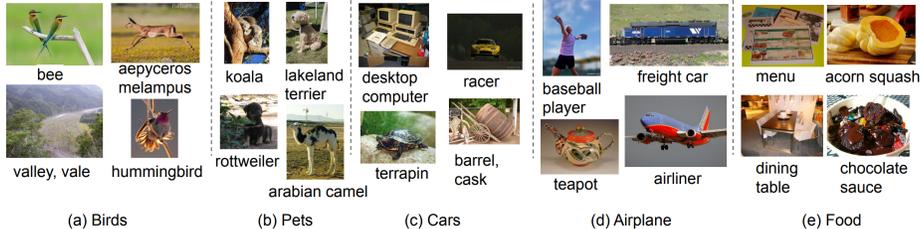
examples across various domains. DATL performs worse than fine-tuning on Birdsnap and Aircraft, unlike L2TL. This underlines the importance of leveraging the visual similarity in the ways beyond label matching as in DATL. When the much larger JFT-300M source dataset is considered, Table 3 shows that L2TL shows even greater benefits in learning related samples to extract knowledge from despite the large-scale of options, demonstrating **7.5%** improvement over the fine-tuning baseline.

We additionally conduct the experiments on transfer learning from MNIST to SVHN [34], in a setting similar to [29]. Although both datasets correspond to the same content of digits, the font style of digits are quite distinct, with varying degree of differences among individual digits. Following [29], we construct the SVHN dataset by randomly sampling 60 images per class, resulting in 600 images in the training split. We use the pretrained LeNet [25] for transferring source dataset of MNIST to the target dataset of SVHN. Table 4 shows that comparing to fine-tuning, L2TL obtains more than 3.5% improvement in performance, via upweighing of the relevant digit images from MNIST. This also validates the effectiveness of L2TL even with small-scale source datasets.

Learning importance weights. We study the effectiveness of learning importance weights in L2TL by comparing to two baselines: (i) random search: the policy model is not optimized and random actions are chosen as the policy output, and (ii) uniform weights: a constant importance weight is assigned to all training samples. Note that for these baselines, α is still optimized via policy gradient. We show the best results of the baselines, after optimizing the hyperparameters on the validation set. As shown in Fig. 3, L2TL outperforms both after sufficient number of iterations, demonstrating the importance of reweighting via policy gradient. L2TL converges to the final result in the last few thousand it-

Table 5. Target dataset test accuracy (%) on Stanford Car with different number of target samples per class.

Method	Number of samples per class			
	5	10	15	20
Fine-tuning	40.0	70.3	80.5	84.9
Our L2TL	46.5	73.7	83.0	86.1

**Fig. 5.** Representative examples from the source datasets with high weight from L2TL for different target datasets. In most cases, we observe the selected examples to be highly-related to the target dataset.

erations, although a larger variance is observed at the beginning of training. We do not observe large variations in the final performance with different runs (e.g. the standard deviation of performance is around 0.1 % over 3 runs). As the classifiers of both the source and the target dataset converge, the small variations in the weights for each class would not heavily affect the final performance.

Small target dataset regime. In the extreme regime of very small number of training examples, generalizing to unseen examples is particularly challenging as the model can be prone to overfitting. Fig. 4 shows that 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) (see Table 5). 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% with 60 examples per class (for Birdsnap). These underline the potential of L2TL for significant performance improvements in real-world tasks where the number of training examples are limited.

High-weight source samples. To build insights on the learned weights, we sample $10k$ actions from the policy and rank the source labels according to their weights. For Birds, the top source class is “bee eater” which is one of the bird species in ImageNet. The second top “aepyceros melampus” is an antelope that has narrow mouth, which is similar to some birds with sharp spout. The “valley” also matches the background in some images. For Cars, we interestingly observe the high-weight class “barrel, cask”, which indeed include 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, the high-weight classes seem relevant in a more subtle way – e.g., “caldron” might

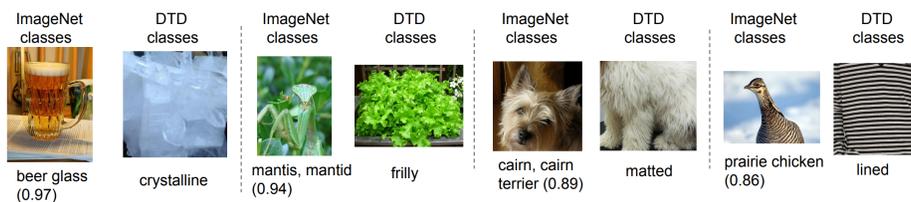


Fig. 6. Top source classes with the highest weight from ImageNet while transferring to DTD target. Representative images from each ImageNet class are shown along with related examples from DTD.

Table 6. Results on the test set on DTD, split 1.

Number of Training examples	Method	Accuracy (%)
Full training set	Random initialization	57.4
	Fine-tuning	70.3
	L2TL	72.0
10 examples per category	Fine-tuning	55.0
	Our L2TL	60.1

have images with food inside. Fig. 5 visualizes a few representative examples from each dataset. More classes can be found in our supplementary material. These demonstrate that L2TL can carefully extract the related classes from the source based on the pattern/shape of the objects, or background scenes. L2TL yields ranking of the source data samples, which can be utilized as new forms of interpretable insights for model developers.

5.3 Dissimilar domain transfer learning

We evaluate L2TL on datasets that are dissimilar to the source dataset, where alternative methods like DATL cannot be readily applied. Table 6 shows the results on DTD. We observe that ImageNet fine-tuning greatly improves the classification results compared to training from random initialization. L2TL further improves the fine-tuning baseline by 1.5%, demonstrating the strength of L2TL selectively using related source classes instead of all classes. For the low-shot target dataset regime, with 10 examples per class, the improvement is more than 5%, suggesting the premise of L2TL even more strongly. Fig. 6 shows that L2TL is able to utilize visually-similar patterns between the source and the target classes. The 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. Training with such visually-similar patterns especially helps the low layers of the networks as they can reuse most of the relevant representations when transferring knowledge [29].

Similarly, Table 7 shows the results on CheXpert, using target validation AUC as a L2TL reward. L2TL performs better than the fine-tuning baseline

Table 7. AUC comparisons on the CheXpert dataset. We followed the the same evaluation protocol in [19] .

Method	Atelectasis	Cardiomegaly	Consolidation	Edema	Pleural Effusion	Mean
Fine-tuning [19]	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
Our L2TL	86.1	84.4	91.5	94.8	93.2	90.0

Table 8. Computational cost of training on Cloud TPU v2. We use Inception-V3 as the backbone. The last column (“With PT model”) assumes availability of a pre-trained source model. “TL” denotes transfer learning from source to target.

Method	Number of iterations		Time per iterations		Total time	
	Pre-training	TL	Pre-training	TL	From scratch	With PT model
Fine-tuning	213,000	20,000	0.14s	0.21s	9.5h	1.2h
DATL [35]	713,000	20,000	0.14s	0.21s	28.9h	20.6h
Our L2TL	213,000	20,000	0.14s	0.75s	12.5h	4.2h

with an AUC improvement of 0.7. There are not many straightforward visual similarities to humans between ImageNet and CheXpert, but L2TL is still capable of discovering them to improve performance.

5.4 Computational cost of training

L2TL uses both the source and target data for training, and the source data can be potentially very large, but the excess computational overhead of L2TL is indeed not large. Table 8 presents the computational cost for fine-tuning, DATL and L2TL with Imagenet as the source dataset. 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 source and target datasets. Yet, it only requires a single training pass on the source dataset, and thus the training time is much lower compared to DATL, and only $\sim 30\%$ higher than fine-tuning when the whole training is considered.

6 Conclusions

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. We demonstrate the effectiveness of L2TL for various cases. L2TL consistently improves fine-tuning across all datasets. The performance benefit of L2TL is more significant for small-scale target datasets or large-scale source datasets. Even for the cases where source and target datasets come from substantially-different domains, L2TL still yields clear improvements.

References

1. Amodei, D., Anubhai, R., Battenberg, E., Case, C., Casper, J., et al.: Deep speech 2: End-to-end speech recognition in english and mandarin. In: ICML (2016)
2. Arik, S.Ö., Chen, J., Peng, K., Ping, W., Zhou, Y.: Neural voice cloning with a few samples. In: NeurIPS (2018)
3. Berg, T., Liu, J., Woo Lee, S., Alexander, M.L., Jacobs, D.W., Belhumeur, P.N.: Birdsnap: Large-scale fine-grained visual categorization of birds. In: CVPR (2014)
4. Bossard, L., Guillaumin, M., Van Gool, L.: Food-101—mining discriminative components with random forests. In: ECCV (2014)
5. Cimpoi, M., Maji, S., Kokkinos, I., Mohamed, S., , Vedaldi, A.: Describing textures in the wild. In: CVPR (2014)
6. Cubuk, E.D., Zoph, B., Mane, D., Vasudevan, V., Le, Q.V.: Autoaugment: Learning augmentation policies from data. arXiv preprint arXiv:1805.09501 (2018)
7. Cui, Y., Song, Y., Sun, C., Howard, A., Belongie, S.: Large scale fine-grained categorization and domain-specific transfer learning. In: CVPR (2018)
8. Devlin, J., Chang, M., Lee, K., Toutanova, K.: BERT: pre-training of deep bidirectional transformers for language understanding. arxiv:1810.04805 (2018)
9. Finn, C., Abbeel, P., Levine, S.: Model-agnostic meta-learning for fast adaptation of deep networks. In: ICML (2017)
10. Ge, W., Yu, Y.: Borrowing treasures from the wealthy: Deep transfer learning through selective joint fine-tuning. In: CVPR (2017)
11. Ghorbani, A., Zou, J.: Data shapley: Equitable valuation of data for machine learning. In: ICML (2019)
12. Girshick, R.B., Donahue, J., Darrell, T., Malik, J.: Rich feature hierarchies for accurate object detection and semantic segmentation. In: CVPR (2014)
13. Gong, B., Shi, Y., Sha, F., Grauman, K.: Geodesic flow kernel for unsupervised domain adaptation. In: CVPR (2012)
14. Guo, Y., Shi, H., Kumar, A., Grauman, K., Rosing, T., Feris, R.S.: Spottune: Transfer learning through adaptive fine-tuning. In: CVPR (2019)
15. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR (2016)
16. Hestness, J., Narang, S., Ardalani, N., Diamos, G.F., Jun, H., Kianinejad, H., Patwary, M.M.A., Yang, Y., Zhou, Y.: Deep learning scaling is predictable, empirically. arXiv:1712.00409 (2017)
17. Hinton, G., Vinyals, O., Dean, J.: Distilling the knowledge in a neural network. arXiv preprint arXiv:1503.02531 (2015)
18. Huang, G., Liu, Z., Van Der Maaten, L., Weinberger, K.Q.: Densely connected convolutional networks. In: CVPR (2017)
19. Irvin, J., Rajpurkar, P., Ko, M., Yu, Y., Ciurea-Ilcus, S., Chute, C., Marklund, H., Haghgoo, B., Ball, R., Shpanskaya, K., et al.: Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison. In: AAAI (2019)
20. Jenni, S., Favaro, P.: Deep bilevel learning. In: ECCV (2018)
21. Jiang, L., Zhou, Z., Leung, T., Li, L., Fei-Fei, L.: Mentornet: Regularizing very deep neural networks on corrupted labels. In: ICML (2018)
22. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: ICLR (2014)
23. Krause, J., Deng, J., Stark, M., Fei-Fei, L.: Collecting a large-scale dataset of fine-grained cars. The Second Workshop on Fine-Grained Visual Categorization, 2013

24. Kumar, M.P., Packer, B., Koller, D.: Self-paced learning for latent variable models. In: *NeurIPS* (2010)
25. LeCun, Y., Bottou, L., Bengio, Y., Haffner, P., et al.: Gradient-based learning applied to document recognition. *Proceedings of the IEEE* **86**(11), 2278–2324 (1998)
26. Lee, J., Yoon, W., Kim, S., Kim, D., Kim, S., So, C.H., Kang, J.: Biobert: a pre-trained biomedical language representation model for biomedical text mining. *arXiv:1901.08746* (2019)
27. Li, X., Grandvalet, Y., Davoine, F.: Explicit inductive bias for transfer learning with convolutional networks. In: *ICML* (2018)
28. Lin, T., Goyal, P., Girshick, R.B., He, K., Dollár, P.: Focal loss for dense object detection. In: *ICCV* (2017)
29. Liu, H., Long, M., Wang, J., Jordan, M.I.: Towards understanding the transferability of deep representations (2019)
30. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: *CVPR* (2015)
31. Mahajan, D., Girshick, R.B., Ramanathan, V., He, K., Paluri, M., Li, Y., Bharambe, A., van der Maaten, L.: Exploring the limits of weakly supervised pre-training. In: *ECCV* (2018)
32. Maji, S., Rahtu, E., Kannala, J., Blaschko, M., Vedaldi, A.: Fine-grained visual classification of aircraft. *arXiv:1306.5151* (2013)
33. Moon, S., Kim, S., Wang, H.: Multimodal transfer deep learning for audio visual recognition. *arXiv:1412.3121* (2014)
34. Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B., Ng, A.Y.: Reading digits in natural images with unsupervised feature learning (2011)
35. Ngiam, J., Peng, D., Vasudevan, V., Kornblith, S., Le, Q.V., Pang, R.: Domain adaptive transfer learning with specialist models. *arXiv preprint arXiv:1811.07056* (2018)
36. van den Oord, A., Dieleman, S., Zen, H., Simonyan, K., Vinyals, O., et al.: Wavenet: A generative model for raw audio. *arXiv:1609.03499* (2016)
37. Parkhi, O.M., Vedaldi, A., Zisserman, A., Jawahar, C.: Cats and dogs. In: *CVPR* (2012)
38. Peng, Y., He, X., Zhao, J.: Object-part attention model for fine-grained image classification. *TIP* (2017)
39. Pham, H., Guan, M.Y., Zoph, B., Le, Q.V., Dean, J.: Efficient neural architecture search via parameter sharing. In: *ICML* (2018)
40. Ravi, S., Larochelle, H.: Optimization as a model for few-shot learning. In: *ICLR* (2017)
41. Razavian, A.S., Azizpour, H., Sullivan, J., Carlsson, S.: CNN features off-the-shelf: an astounding baseline for recognition. In: *CVPR* (2014)
42. Ren, M., Zeng, W., Yang, B., Urtasun, R.: Learning to reweight examples for robust deep learning. In: *ICML* (2019)
43. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al.: Imagenet large scale visual recognition challenge. *IJCV* (2015)
44. Santoro, A., Bartunov, S., Botvinick, M., Wierstra, D., Lillicrap, T.: Meta-learning with memory-augmented neural networks. In: *ICML* (2016)
45. Schmidhuber, J., Zhao, J., Wiering, M.: Shifting inductive bias with success-story algorithm, adaptive levin search, and incremental self-improvement. *Machine Learning* (1997)
46. Shu, J., Xie, Q., Yi, L., Zhao, Q., Zhou, S., Xu, Z., Meng, D.: Meta-weight-net: Learning an explicit mapping for sample weighting. In: *NeurIPS* (2019)

47. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. In: ICLR (2015)
48. Wang, D., Shen, Z., Shao, J., Zhang, W., Xue, X., Zhang, Z.: Multiple granularity descriptors for fine-grained categorization. In: ICCV (2015)
49. Williams, R.J.: Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning* (1992)
50. Wu, Z., Xiong, C., Ma, C.Y., Socher, R., Davis, L.S.: Adafame: Adaptive frame selection for fast video recognition. In: CVPR (2019)
51. Zhu, L., Yang, Y.: Label independent memory for semi-supervised few-shot video classification. *TPAMI* (2020). <https://doi.org/10.1109/TPAMI.2020.3007511>
52. Zoph, B., Le, Q.V.: Neural architecture search with reinforcement learning. In: ICLR (2017)
53. Zoph, B., Yuret, D., May, J., Knight, K.: Transfer learning for low-resource neural machine translation. In: ACL (2016)