



# Do We Really Need to Access the Source Data? Source Hypothesis Transfer for Unsupervised Domain Adaptation

**Jian Liang**   Dapeng Hu   Jiashi Feng

Learning and Vision Lab, National University of Singapore (NUS)

International Conference on Machine Learning, 2020



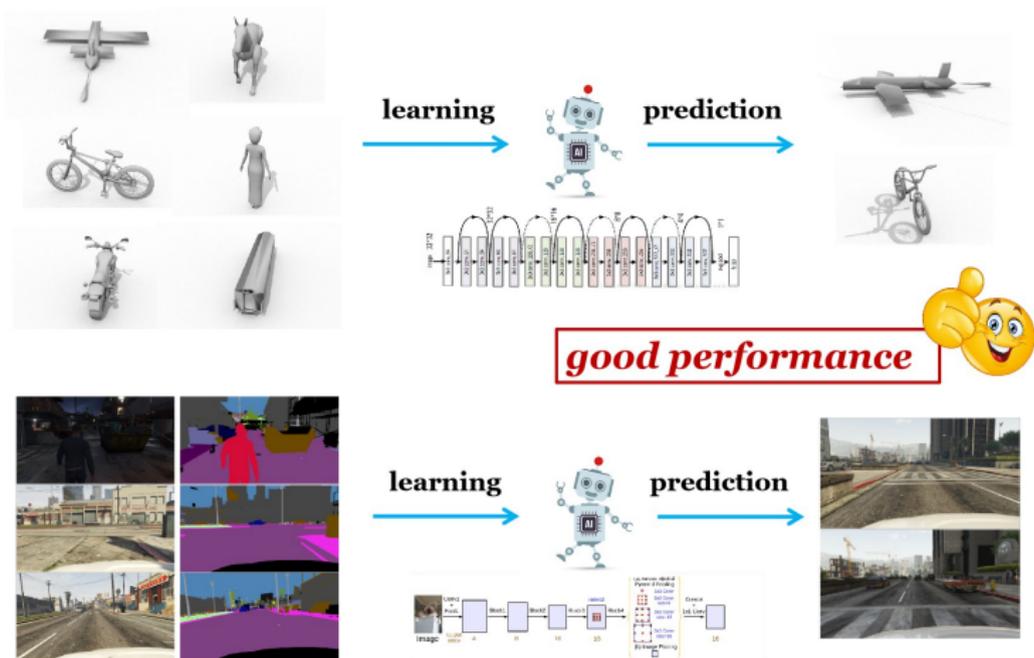
# Outline

1 Background

2 Method

3 Experiments

# Situation 1



**Figure 1:** Test data comes from the same distribution as training data!

# Situation 2

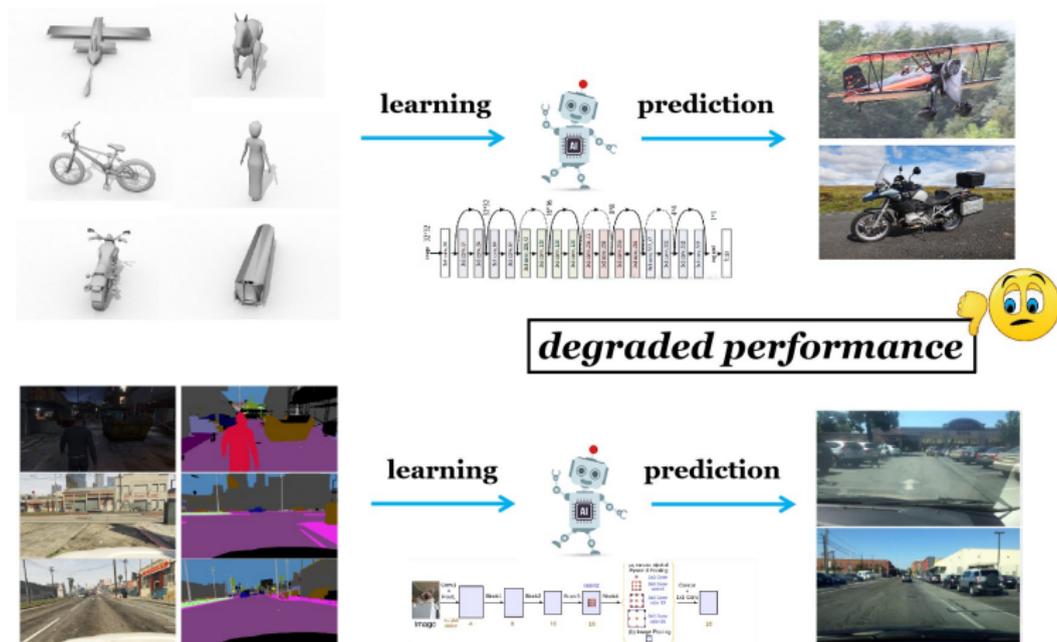


Figure 2: Test data and training data comes from different distributions!

# Unsupervised Domain Adaptation (DA)

## Vanilla setting

- **Source Domain**  $\mathcal{D}_s$ :  $n_s$  labeled samples  $\{x_s^i, y_s^i\}_{i=1}^{n_s}$  from  $P_S(X, Y)$ ;
- **Target Domain**  $\mathcal{D}_t$ :  $n_t$  unlabeled samples  $\{x_t^i, ?\}_{i=1}^{n_t}$  from  $P_T(X, Y)$ ;
- **Goal**: Use  $\{x_t^i\}_{i=1}^{n_t}$  during training (transductive) and learn a good classifier to get the values of ? under **domain shift** (i.e.,  $P_S \neq P_T$ ).



Classification



Re-identification



Detection



Control



Segmentation



Visual Localization

Credit to Gabriela Csurka, TaskCV-2019 talk.

# Previous DA Methods

## (I) Input-level Pixel Transfer

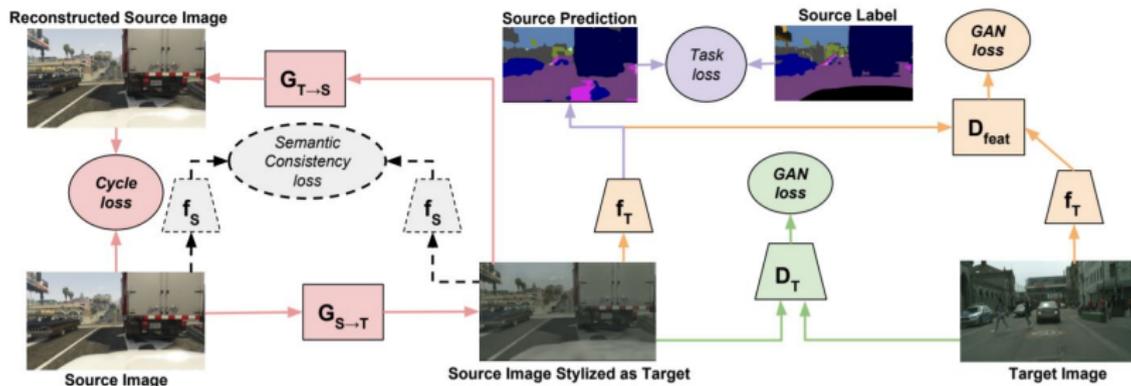


Figure 3: Cycle-consistent adversarial adaptation (CyCADA) <sup>1</sup> overview.



<sup>1</sup>Hoffman, Judy, et al. "CyCADA: Cycle-Consistent Adversarial Domain Adaptation." In ICML 2018.

# Previous DA Methods

## (II) Feature-level Alignment

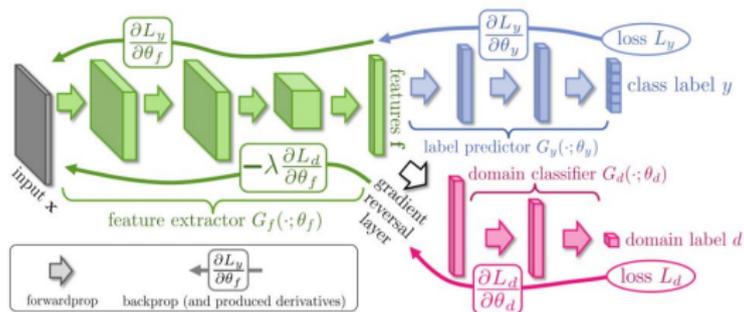


Figure 4: Unsupervised Domain Adaptation by Backpropagation (DANN)<sup>2</sup> overview.

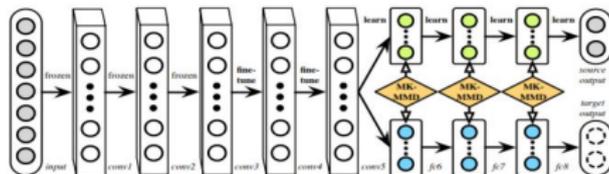


Figure 5: Deep Adaptation Networks (DAN)<sup>3</sup> overview.

<sup>2</sup> Ganin, Yaroslav, and Victor Lempitsky. "Unsupervised Domain Adaptation by Backpropagation." In ICML 2015.

<sup>3</sup> Long, Mingsheng, et al. "Learning Transferable Features with Deep Adaptation Networks." In ICML 2015.

# Previous DA Methods

## (III) Output-level Regularization

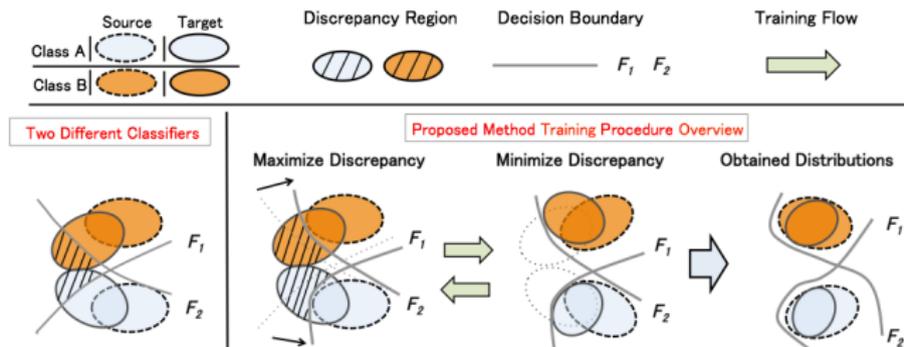


Figure 6: Maximum classifier discrepancy (MCD)<sup>4</sup> overview.

Or exploit the low-density separation principle:

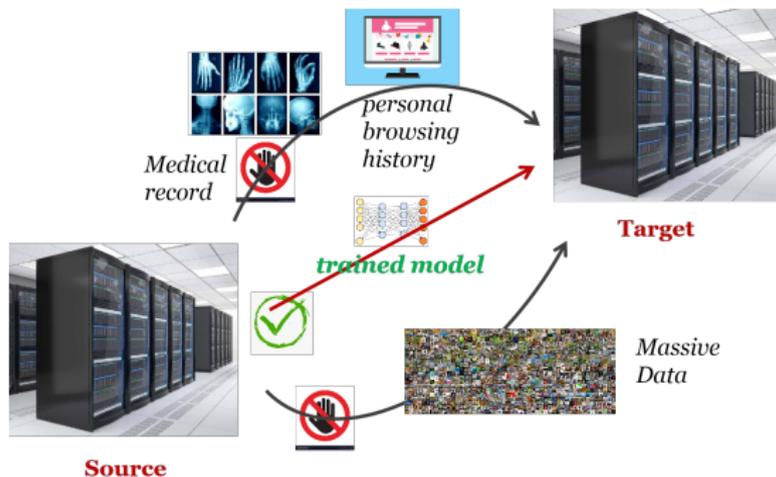
- entropy minimization
- pseudo-labeling / self-training
- virtual adversarial training
- consistency regularization

<sup>4</sup>Saito, Kuniaki, et al. "Maximum classifier discrepancy for unsupervised domain adaptation." In CVPR 2018.

# So is DA solved?

## Limitation of existing DA methods

- **Not Secure**: the full access to source data is required.
- **Concentrated**: processing different domains in the same machine.



# Outline

1 Background

**2 Method**

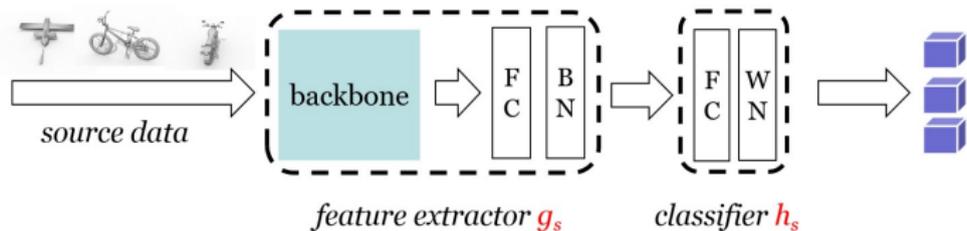
3 Experiments

# Unsupervised Domain Adaptation

## Model Adaptation Setting

- **Source Model**  $f_s : \mathcal{X}_s \rightarrow \mathcal{Y}_s$  trained on  $\mathcal{D}_s$ ;
- **Target Domain**  $\mathcal{D}_t$ :  $n_t$  unlabeled samples  $\{x_t^i, ?\}_{i=1}^{n_t}$ ;
- **Goal**: learn a good classifier  $f_t : \mathcal{X}_t \rightarrow \mathcal{Y}_t$  to get the values of ?.

# How to generate a good source model $f_s$ ?



**Figure 7:** The network of source model for object recognition.  $f_s(x) = h_s(g_s(x))$ , where  $g_s : \mathcal{X}_s \rightarrow \mathbb{R}^d$  and  $h_s : \mathbb{R}^d \rightarrow \mathbb{R}^K$ .

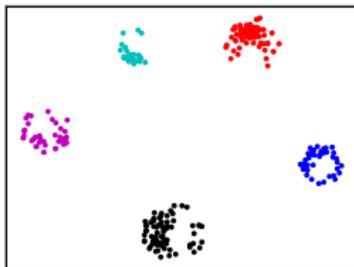
## Classification loss:

$$\mathcal{L}_{src}^{ls}(f_s; \mathcal{X}_s, \mathcal{Y}_s) = -\mathbb{E}_{(x_s, y_s) \in \mathcal{X}_s \times \mathcal{Y}_s} \sum_{k=1}^K q_k^{ls} \log(\delta_k(f_s(x_s))), \quad (1)$$

where  $q_k^{ls} = (1 - \alpha)q_k + \alpha/K$  is the smoothed label and  $\alpha$  is the smoothing parameter which is empirically set to 0.1.

## What we can learn from the model $f_s$ ?

$\forall y_s = k$ , maximizing  $f_s^{(k)}(x_s) = \frac{\exp(w_k^\top g_s(x_s))}{\sum_i \exp(w_i^\top g_s(x_s))}$  means minimizing the distance between  $g_s(x_s)$  and  $w_k$ , where  $w_k$  is the  $k$ -th weight vector in  $h_s$ .

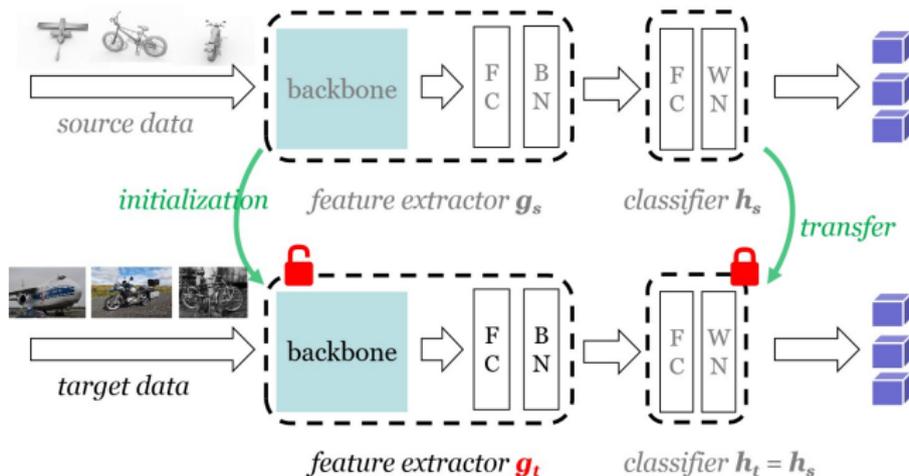


**Figure 8:** t-SNE visualizations of source features  $g_s(x)$ ,  $x \in \mathcal{X}_s$ . Each color denotes one class.

Fortunately, even we have no access to the source data or features directly, we may still estimate the distribution of source features via  $h_s$ .

# Source Hypothesis Transfer

## Framework



**Figure 9:** The proposed Source Hypothesis Transfer (SHOT) framework.

Ideally, we expect the feature extractor  $g_t$  can produce source-like features for target data, that is to say, the corresponding outputs of  $h_s$  are also close to one-hot encoding like those of source features.

# Source Hypothesis Transfer

## Information Maximization (IM)

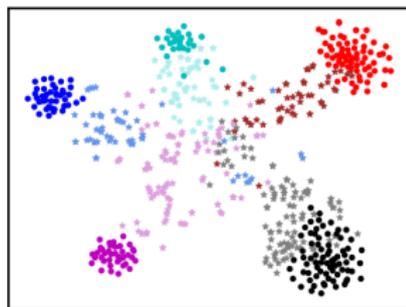
In practice, we minimize the following  $\mathcal{L}_{ent}$  and  $\mathcal{L}_{div}$  that together constitute the IM loss:  $[f_t(x) = h_s(g_t(x))]$

$$\begin{aligned}\mathcal{L}_{ent}(f_t; \mathcal{X}_t) &= -\mathbb{E}_{x_t \in \mathcal{X}_t} \sum_{k=1}^K \delta_k(f_t(x_t)) \log(\delta_k(f_t(x_t))), \\ \mathcal{L}_{div}(f_t; \mathcal{X}_t) &= \sum_{k=1}^K \hat{p}_k \log \hat{p}_k = D_{KL}(\hat{p}_k, \frac{1}{K} \mathbf{1}_K) - \log K,\end{aligned}\tag{2}$$

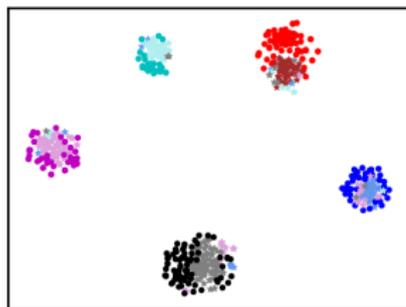
where  $f_t(x) = h_t(g_t(x))$  is the  $K$ -dimensional output of each target sample,  $\hat{p} = \mathbb{E}_{x_t \in \mathcal{X}_t}[\delta(f_t^{(k)}(x_t))]$  is the mean output embedding of the whole target domain, and  $\mathbf{1}_K$  is a  $K$ -dimensional vector with all ones.

# Source Hypothesis Transfer

## Information Maximization (IM) - (Cont'd)



(a) Source model only



(b) SHOT-IM

**Figure 10:** t-SNE visualizations. Circles in dark colors denote the **unseen** source data and stars in light denote the target data. Different colors represent different classes.

IM loss relies heavily on the initialization and does not fully consider **the structure of target data**. Even features from different domains are well aligned, there still exists **cross-label matching**.

# Source Hypothesis Transfer

## Self-supervised Pseudo-labeling

We exploit **target-specific centroids** to obtain accurate pseudo labels.

### 1 Centroid Initialization & Cluster Assignment.

$$c_k^{(0)} = \frac{\sum_{x_t \in \mathcal{X}_t} \delta(\hat{f}_t^{(k)}(x)) \hat{g}_t(x)}{\sum_{x_t \in \mathcal{X}_t} \delta(\hat{f}_t^{(k)}(x))}, \quad (3)$$
$$\hat{y}_t = \arg \min_k D_f(\hat{g}_t(x_t), c_k^{(0)}),$$

### 2 Centroid Update & Cluster Assignment.

$$c_k^{(1)} = \frac{\sum_{x_t \in \mathcal{X}_t} \mathbb{1}(\hat{y}_t = k) \hat{g}_t(x)}{\sum_{x_t \in \mathcal{X}_t} \mathbb{1}(\hat{y}_t = k)}, \quad (4)$$
$$\hat{y}_t = \arg \min_k D_f(\hat{g}_t(x_t), c_k^{(1)}).$$

\*  $D_f(a, b)$  measures the cosine distance between  $a$  and  $b$ .

# Source Hypothesis Transfer

## Complete objective

$$\mathcal{L}(g_t) = \mathcal{L}_{ent}(h_s \circ g_t; \mathcal{X}_t) + \mathcal{L}_{div}(h_s \circ g_t; \mathcal{X}_t) - \beta \mathbb{E}_{(x_t, \hat{y}_t) \in \mathcal{X}_t \times \hat{\mathcal{Y}}_t} \sum_{k=1}^K \mathbb{1}_{[k=\hat{y}_t]} \log(\delta_k(h_s(g_t(x_t))))). \quad (5)$$

## Difference with prior work.

Both TDA<sup>a</sup> and MCS<sup>b</sup> are shallow methods that ignore feature representation learning, deteriorating the performance. FADA<sup>c</sup> is elegantly designed for multi-source domain adaptation.

<sup>a</sup>Chidlovskii, Boris, Stephane Clinchant, and Gabriela Csurka. "Domain adaptation in the absence of source domain data." In KDD 2016.

<sup>b</sup>Liang, Jian, et al. "Distant supervised centroid shift: A simple and efficient approach to visual domain adaptation." In CVPR 2019.

<sup>c</sup>Peng, Xingchao, et al. "Federated Adversarial Domain Adaptation." In ICLR 2020.

# Outline

1 Background

2 Method

3 Experiments

# Setup

## Data Sets and Various Scenarios

- 1 Digit recognition (**MNIST**, **USPS**, **SVHN**)
- 2 Cross-domain object recognition (**Office**, **Office-Home**, **Office-Caltech**)
- 3 Synthetic-to-real object recognition (**VisDA-C**)

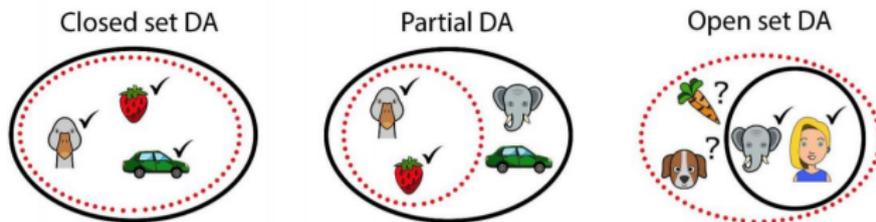
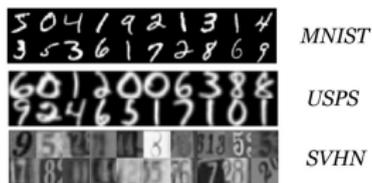


Figure 11: Typical UDA scenarios.

# Results

## Vanilla Closed-set Domain Adaptation

Method (Source→Target)	Ar→Cl	Ar→Pr	Ar→Re	Cl→Ar	Cl→Pr	Cl→Re	Pr→Ar	Pr→Cl	Pr→Re	Re→Ar	Re→Cl	Re→Pr	Avg.
<b>DANN</b> (ICML 2015)	45.6	59.3	70.1	47.0	58.5	60.9	46.1	43.7	68.5	63.2	51.8	76.8	57.6
<b>DAN</b> (ICML 2015)	43.6	57.0	67.9	45.8	56.5	60.4	44.0	43.6	67.7	63.1	51.5	74.3	56.3
<b>CDAN+E</b> (NeurIPS 2018)	50.7	70.6	76.0	57.6	70.0	70.0	57.4	50.9	77.3	70.9	56.7	81.6	65.8
<b>CDAN+BSP</b> (ICML 2019)	52.0	68.6	76.1	58.0	70.3	70.2	58.6	50.2	77.6	72.2	<b>59.3</b>	81.9	66.3
<b>SAFN</b> (ICCV 2019)	52.0	71.7	76.3	64.2	69.9	71.9	63.7	51.4	77.1	70.9	57.1	81.5	67.3
<b>CDAN+TransNorm</b> (NeurIPS 2019)	50.2	71.4	77.4	59.3	72.7	73.1	61.0	53.1	79.5	71.9	59.0	82.9	67.6
<b>Source model only</b>	44.6	67.3	74.8	52.7	62.7	64.8	53.0	40.6	73.2	65.3	45.4	78.0	60.2
<b>SHOT-IM (ours)</b>	55.4	76.6	80.4	66.9	74.3	75.4	65.6	54.8	80.7	73.7	58.4	83.4	70.5
<b>SHOT (ours)</b>	<b>57.1</b>	<b>78.1</b>	<b>81.5</b>	68.0	<b>78.2</b>	<b>78.1</b>	67.4	<b>54.9</b>	<b>82.2</b>	73.3	58.8	84.3	<b>71.8</b>
<b>SRDC</b> <sup>5</sup> (CVPR 2020)	52.3	76.3	81.0	<b>69.5</b>	76.2	78.0	<b>68.7</b>	53.8	81.7	<b>76.3</b>	57.1	<b>85.0</b>	71.3

Table 1: Accuracies (%) on Office-Home dataset (ResNet-50).

Method (Source→Target)	S→M	U→M	M→U	Avg.
<b>ADDA</b> (CVPR 2017)	76.0±1.8	90.1±0.8	89.4±0.2	85.2
<b>ADR</b> (ICLR 2018)	95.0±1.9	93.1±1.3	93.2±2.5	93.8
<b>CDAN+E</b> (NeurIPS 2018)	89.2	98.0	95.6	94.3
<b>CyCADA</b> (ICML 2018)	90.4±0.4	96.5±0.1	95.6±0.4	94.2
<b>rRevGrad+CAT</b> (ICCV 2019)	98.8±0.0	96.0±0.9	94.0±0.7	96.3
<b>SWD</b> (CVPR 2019)	<b>98.9±0.1</b>	97.1±0.1	<b>98.1±0.1</b>	<b>98.0</b>
<b>Source model only</b>	67.1±0.9	87.8±2.3	89.6±0.4	81.5
<b>SHOT-IM (ours)</b>	89.6±5.0	96.8±0.4	91.9±0.4	92.8
<b>SHOT (ours)</b>	<b>98.9±0.0</b>	<b>98.4±0.6</b>	<b>98.0±0.2</b>	<b>98.4</b>
<b>STAR</b> <sup>6</sup> (CVPR 2020)	98.8±0.1	97.7±0.1	97.8±0.1	98.1
<b>Target-supervised (Oracle)</b>	99.4±0.0	99.4±0.0	98.0±0.1	98.9

Table 2: Accuracies (%) on Digits dataset. S: SVHN, M: MNIST, U: USPS.

<sup>5</sup>Tang, Hui, et al. "Unsupervised Domain Adaptation via Structurally Regularized Deep Clustering." In CVPR 2020.

<sup>6</sup>Lu, Zhihe, et al. "Stochastic Classifiers for Unsupervised Domain Adaptation." In CVPR 2020.

# Results

## Multi-source and Multi-target Domain Adaptation

Multi-source (R $\rightarrow$ )	R $\rightarrow$ A	R $\rightarrow$ C	R $\rightarrow$ D	R $\rightarrow$ W	Avg.
<b>FADA</b> (ICLR 2020)	84.2	88.7	87.1	88.1	87.1
<b>DAN</b> (ICML 2015)	91.6	89.2	99.1	99.5	94.8
<b>DCTN</b> (CVPR 2018)	92.7	90.2	99.0	99.4	95.3
<b>MCD</b> (CVPR 2018)	92.1	91.5	99.1	99.5	95.6
<b>M<sup>3</sup>SDA-<math>\beta</math></b> (ICCV 2019)	94.5	92.2	<b>99.2</b>	99.5	96.4
<b>Source model only</b>	95.4	93.7	98.9	98.3	96.6
<b>SHOT-IM (ours)</b>	96.2	96.1	98.5	99.7	97.6
<b>SHOT (ours)</b>	<b>96.4</b>	<b>96.2</b>	98.5	<b>99.7</b>	<b>97.7</b>

Multi-target ( $\rightarrow$ R)	A $\rightarrow$ R	C $\rightarrow$ R	D $\rightarrow$ R	W $\rightarrow$ R	Avg.
<b>SE</b> (ICLR 2018)	90.3	94.7	88.5	85.3	89.7
<b>MCD</b> (CVPR 2018)	91.7	95.3	89.5	84.3	90.2
<b>DANN</b> (ICML 2015)	91.5	94.3	90.5	86.3	90.7
<b>DADA</b> (ICML 2019)	92.0	95.1	91.3	93.1	92.9
<b>Source model only</b>	90.7	96.1	90.2	90.9	92.0
<b>SHOT-IM (ours)</b>	95.7	97.2	<b>96.3</b>	96.1	96.3
<b>SHOT (ours)</b>	<b>96.2</b>	<b>97.3</b>	<b>96.3</b>	<b>96.2</b>	<b>96.5</b>

**Table 3:** Accuracies (%) on **Office-Caltech** dataset (ResNet-101). [**\*R** denotes the **rest** three domains except the single source / target.]

# Results

## Partial-set and Open-set Domain Adaptation

Partial-set DA (Source→Target)	Ar→Cl	Ar→Pr	Ar→Re	Cl→Ar	Cl→Pr	Cl→Re	Pr→Ar	Pr→Cl	Pr→Re	Re→Ar	Re→Cl	Re→Pr	Avg.
<b>IWAN</b> (CVPR 2018)	53.9	54.5	78.1	61.3	48.0	63.3	54.2	52.0	81.3	76.5	56.8	82.9	63.6
<b>SAN</b> (ECCV 2018)	44.4	68.7	74.6	67.5	65.0	77.8	59.8	44.7	80.1	72.2	50.2	78.7	65.3
<b>ETN</b> (CVPR 2019)	59.2	77.0	79.5	62.9	65.7	75.0	68.3	55.4	84.4	75.7	57.7	84.5	70.5
<b>SAFN</b> (ICCV 2019)	58.9	76.3	81.4	70.4	73.0	77.8	72.4	55.3	80.4	75.8	60.4	79.9	71.8
Source model only	44.6	67.3	74.8	52.7	62.7	64.8	53.0	40.6	73.2	65.3	45.4	78.0	60.2
SHOT-IM (ours)	57.9	83.6	88.8	72.4	74.0	79.0	76.1	60.6	<b>90.1</b>	<b>81.9</b>	<b>68.3</b>	<b>88.5</b>	76.8
SHOT (full, ours)	<b>64.8</b>	<b>85.2</b>	<b>92.7</b>	<b>76.3</b>	<b>77.6</b>	<b>88.8</b>	<b>79.7</b>	<b>64.3</b>	89.5	80.6	66.4	85.8	<b>79.3</b>
<b>RTNet<sub>adv</sub></b> <sup>7</sup> (CVPR 2020)	63.2	80.1	80.7	66.7	69.3	77.2	71.6	53.9	84.6	77.4	57.9	85.5	72.3
Open-set DA (Source→Target)	Ar→Cl	Ar→Pr	Ar→Re	Cl→Ar	Cl→Pr	Cl→Re	Pr→Ar	Pr→Cl	Pr→Re	Re→Ar	Re→Cl	Re→Pr	Avg.
<b>ATI-λ</b> (ICCV 2017)	55.2	52.6	53.5	69.1	63.5	74.1	61.7	64.5	70.7	79.2	72.9	75.8	66.1
<b>OSBP</b> (ECCV 2018)	56.7	51.5	49.2	67.5	65.5	74.0	62.5	64.8	69.3	80.6	74.7	71.5	65.7
<b>OpenMax</b> (CVPR 2016)	56.5	52.9	53.7	69.1	64.8	74.5	64.1	64.0	71.2	80.3	73.0	76.9	66.7
<b>STA</b> (CVPR 2019)	58.1	53.1	54.4	<b>71.6</b>	69.3	<b>81.9</b>	63.4	<b>65.2</b>	74.9	<b>85.0</b>	<b>75.8</b>	80.8	69.5
Source model only	36.3	54.8	69.1	33.8	44.4	49.2	36.8	29.2	56.8	51.4	35.1	62.3	46.6
SHOT-IM (ours)	62.5	77.8	83.9	60.9	73.4	79.4	64.7	58.7	83.1	69.1	62.0	82.1	71.5
SHOT (full, ours)	<b>64.5</b>	<b>80.4</b>	<b>84.7</b>	63.1	<b>75.4</b>	81.2	65.3	59.3	<b>83.3</b>	69.6	64.6	<b>82.3</b>	<b>72.8</b>
<b>Inheritune</b> <sup>8</sup> (CVPR 2020)	60.1	70.9	83.2	64.0	70.0	75.7	<b>66.1</b>	54.2	81.3	74.9	56.2	78.6	69.6

Table 4: Accuracies (%) on Office-Home dataset (ResNet-50).

<sup>7</sup> Chen, Zhihong, et al. "Selective transfer with reinforced transfer network for partial domain adaptation." In CVPR 2020.

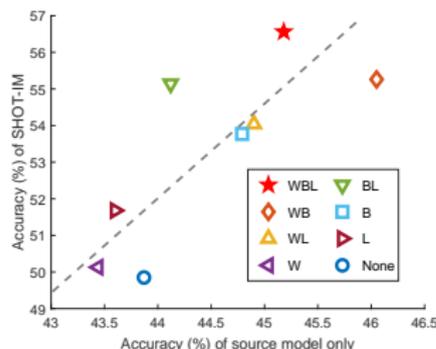
<sup>8</sup> Kundu, Jogendra Nath, et al. "Towards Inheritable Models for Open-Set Domain Adaptation." In CVPR 2020.   

# Analysis

## Ablation study

Methods / Datasets	Office	Office-Home	VisDA-C
Source model only	79.3	60.2	46.6
naive pseudo-labeling (PL) <sup>9</sup>	83.0	64.1	76.6
Self-supervised PL (ours)	87.6	68.9	80.7
$\mathcal{L}_{ent}$	83.5	55.5	63.3
$\mathcal{L}_{ent} + \mathcal{L}_{div}$	87.3	70.5	80.4
$\mathcal{L}_{ent} + \mathcal{L}_{div} + \text{naive PL}$	87.5	70.3	82.9
$\mathcal{L}_{ent} + \mathcal{L}_{div} + \text{Self-supervised PL}$	88.6	71.8	82.9

Table 5: Average accuracies on three closed-set UDA datasets.



Accuracies (%) on the Ar→CI task for *closed-set UDA*. [Weight normalization/ Batch normalization/ Label smoothing]

<sup>9</sup> Lee, Dong-Hyun. "Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks." In ICML Workshop 2013.

# If we cannot train the source model by ourselves?

## An interesting example

To find the answer, we utilize the most popular off-the-shelf pre-trained ImageNet model ResNet-50 and consider a PDA task (**ImageNet**  $\rightarrow$  **Caltech**) to evaluate the effectiveness of SHOT below.

Methods	ResNet-50	ETN <sup>†10</sup>	SHOT-IM	SHOT
Accuracy	69.7 $\pm$ 0.0	83.2 $\pm$ 0.2	81.7 $\pm$ 0.5	<b>83.3 <math>\pm</math> 0.1</b>

**Table 6:** Results of a PDA task (**ImageNet**  $\rightarrow$  **Caltech**). <sup>†</sup>utilizes the training set of ImageNet besides the off-the-shelf pre-trained ResNet-50 model.

<sup>10</sup>Cao, Zhangjie, et al. "Learning to transfer examples for partial domain adaptation." In *CVPR 2019*. 

# Summary

## Take Home Message

- 1 With only the source model provided, **our approach** achieves competitive and even state-of-the-art performance.
- 2 Feature alignment can be achieved **implicitly with output-level regularization** like entropy minimization and information maximization.
- 3 To combat domain shift, **self-supervision** from the target domain itself is quite critical.

# Additional Discussions with Concurrent Works

- 1 Main Idea of MoA <sup>11</sup>
  - generate pseudo source samples
  - pseudo labeled source data & unlabeled target data (semi-supervised learning)
- 2 Main Idea of USFDA <sup>12</sup>
  - simulate labeled negative samples
  - entropy minimization with fixed decision boundary
- 3 Main Idea of SFDA <sup>13</sup>
  - extra target-specific classifier (prototype based) in addition to source-oriented classifier

Difference: We need no additional components like data generator or classifier within the training algorithm.

<sup>11</sup>Li, Rui, et al. "Model Adaptation: Unsupervised Domain Adaptation without Source Data." In CVPR 2020.

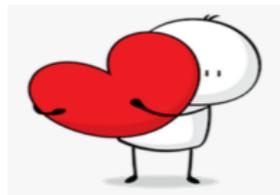
<sup>12</sup>Kundu, Jogendra Nath, Naveen Venkat, and R. Venkatesh Babu. "Universal Source-Free Domain Adaptation." In CVPR 2020.

<sup>13</sup>Kim, Youngeun, et al. "Domain Adaptation without Source Data." Submitted to NeurIPS 2020.

# Thank you!

- 1 Code is available at <https://github.com/tim-learn/SHOT/>.
- 2 If you require any further information, feel free to contact me.

Email: [liangjian92@gmail.com](mailto:liangjian92@gmail.com)



(a) Love



(b) Peace



(c) Health