

# Transferable Calibration with Lower Bias and Variance in Domain Adaptation

Ximei Wang, Mingsheng Long (✉), Jianmin Wang, and Michael I. Jordan<sup>#</sup>

School of Software, Tsinghua University  
National Engineering Laboratory for Big Data Software  
<sup>#</sup>University of California, Berkeley

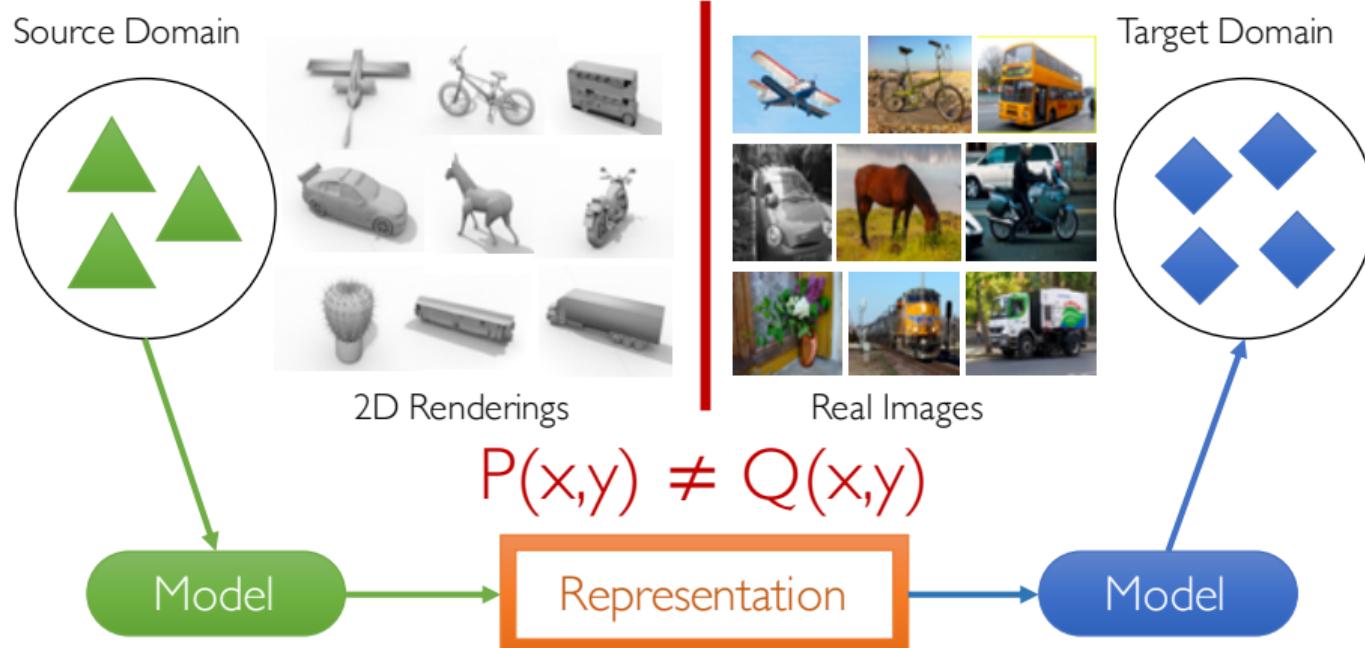
wxm17@mails.tsinghua.edu.cn

<https://w xm17.github.io/>

Neural Information Processing Systems (NeurIPS), 2020

# Domain Adaptation (DA)

Transfer from a labeled source domain to an unlabeled target one.



# Confidence Calibration in Deep Learning<sup>1</sup>

- **Calibration:** A model should output a probability reflecting the true frequency:

$$\mathbb{P}(\hat{Y} = Y | \hat{P} = c) = c, \forall c \in [0, 1] \quad (1)$$

where  $\hat{Y}$  is the class prediction and  $\hat{P}$  is its associated confidence.

- **Calibration Metric:** Expected Calibration Error (ECE)

$$\begin{aligned}\mathcal{L}_{\text{ECE}} &= \sum_{m=1}^B \frac{|B_m|}{n} |\mathbb{A}(B_m) - \mathbb{C}(B_m)| \\ \mathbb{A}(B_m) &= |B_m|^{-1} \sum_{i \in B_m} 1(\hat{y}_i = y_i) \quad (\textbf{Accuracy}) \\ \mathbb{C}(B_m) &= |B_m|^{-1} \sum_{i \in B_m} \max_k p(\hat{y}_i^k | x_i, \theta) \quad (\textbf{Confidence})\end{aligned} \quad (2)$$

- Deep networks learn high accuracy at the expense of over-confidence.

<sup>1</sup> Guo et al. On Calibration of Modern Neural Networks. ICML 2017.

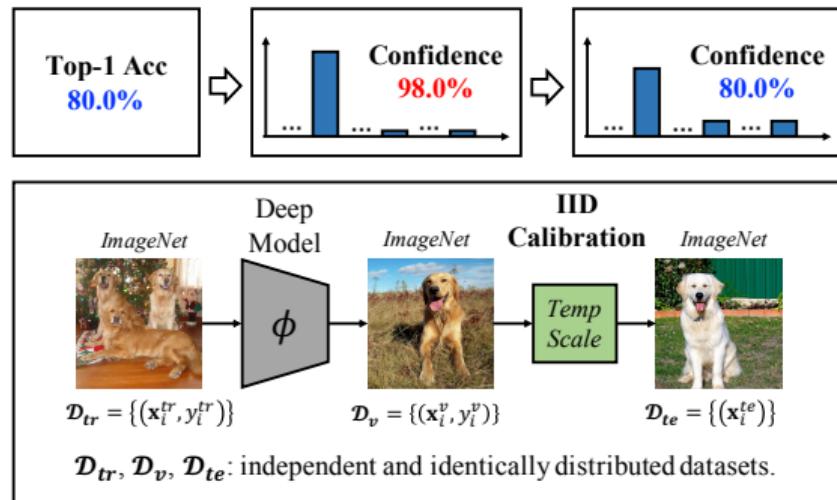
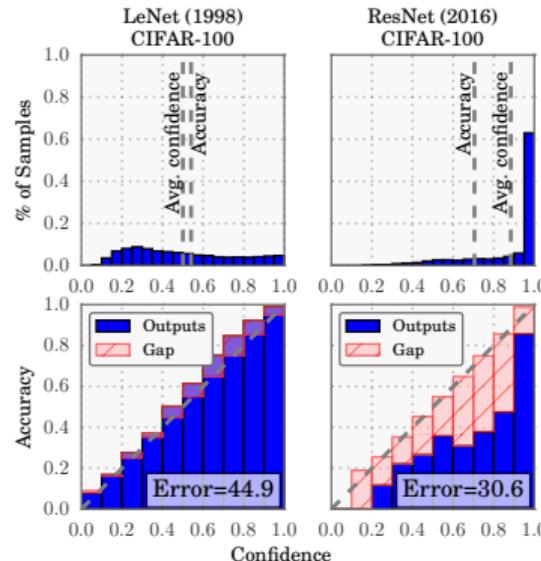
# Temperature Scaling for IID Calibration

- IID Calibration: Temperature Scaling

$$T^* = \arg \min_T E_{(x_v, y_v) \in \mathcal{D}_v} \mathcal{L}_{\text{NLL}}(\sigma(z_v/T), y_v) \quad (3)$$

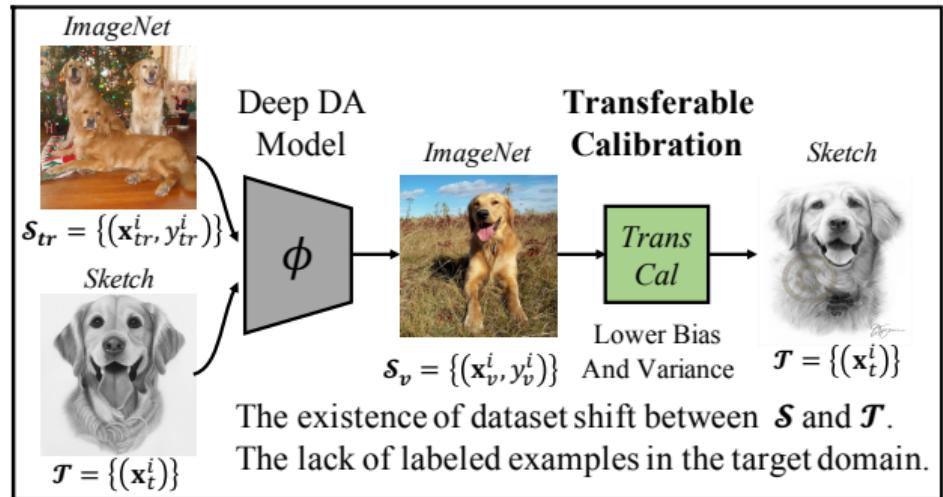
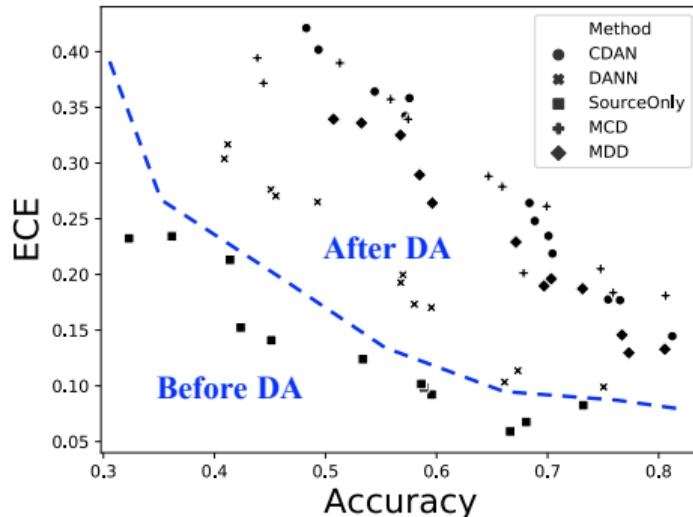
$\sigma$  is the softmax function,  $\mathcal{L}_{\text{NLL}}$  is Negative Log-Likelihood loss.

- Transform logits  $z_{te}$  into calibrated probabilities  $p_{te} = \sigma(z_{te}/T^*)$ .



# Dilemma of Accuracy vs Confidence in DA

- Transfer models yield high accuracy at the expense of over-confidence.



- Calibration in transfer learning is challenging due to the coexistence:
  - Domain shift — ECE should be unbiased to the target domain
  - Unlabeled target — ECE on the target domain is incomputable
- ECE-Accuracy-Dilemma of confidence calibration in Transfer Learning

# Transferable Calibration Framework

- **Transferable Calibration:** Attains on the source data but transfers to the target domain
- $E_{x \sim p} [w(x)\mathcal{L}_{(\cdot)}(\phi(x), y)]$  is an **unbiased** estimator of the target calibration error  $\mathbb{E}_q$

$$\begin{aligned} E_{x \sim q} [\mathcal{L}_{(\cdot)}(\phi(x), y)] &= \int_q \mathcal{L}_{(\cdot)}(\phi(x), y) q(x) dx \\ &= \int_p \frac{q(x)}{p(x)} \mathcal{L}_{(\cdot)}(\phi(x), y) p(x) dx = E_{x \sim p} [w(x)\mathcal{L}_{(\cdot)}(\phi(x), y)] \end{aligned} \tag{4}$$

- Discriminative density ratio estimation method: **LogReg**
- Use Bayesian formula to derive  $\hat{w}(x)$  from a logistic regression classifier

$$\hat{w}(x) = \frac{q(x)}{p(x)} = \frac{v(x|d=0)}{v(x|d=1)} = \frac{P(d=1)}{P(d=0)} \frac{P(d=0|x)}{P(d=1|x)} \tag{5}$$

# Transferable Calibration: Bias Reduction

- Importance-weighting for an unbiased estimate of **target ECE** if  $\hat{w}(x) = w(x)$
- The **bias** between the estimated ECE and the ground-truth ECE

$$\begin{aligned} & \left| E_{x \sim q} \left[ \mathcal{L}_{\text{ECE}}^{\hat{w}(x)} \right] - E_{x \sim q} \left[ \mathcal{L}_{\text{ECE}}^{w(x)} \right] \right| \\ &= |E_{x \sim p} [\hat{w}(x) \mathcal{L}_{\text{ECE}}(\phi(x), y)] - E_{x \sim p} [w(x) \mathcal{L}_{\text{ECE}}(\phi(x), y)]| \\ &= |E_{x \sim p} [(w(x) - \hat{w}(x)) \mathcal{L}_{\text{ECE}}(\phi(x), y)]|. \end{aligned} \tag{6}$$

- The **bias** of them can be further bounded by

$$\begin{aligned} & |E_{x \sim p} [(w(x) - \hat{w}(x)) \mathcal{L}_{\text{ECE}}(\phi(x), y)]| \\ & \leq \sqrt{E_{x \sim p} [(w(x) - \hat{w}(x))^2] E_{x \sim p} [(\mathcal{L}_{\text{ECE}}(\phi(x), y))^2]} \quad (\text{Cauchy-Schwarz Inequality}) \\ & \leq \frac{1}{2} \left( E_{x \sim p} [(w(x) - \hat{w}(x))^2] + E_{x \sim p} [(\mathcal{L}_{\text{ECE}}(\phi(x), y))^2] \right) \quad (\text{AM/GM Inequality}) \end{aligned} \tag{7}$$

# Transferable Calibration: Bias Reduction

- For any  $x$  s.t.  $P(d = 1|x) \neq 0$ , the following inequality holds:

$$\frac{1}{M+1} \leq P(d = 1|x) \leq 1, \quad \text{since } w(x) = \frac{P(d = 0|x)}{P(d = 1|x)} = \frac{1 - P(d = 1|x)}{P(d = 1|x)} = \frac{1}{P(d = 1|x)} - 1. \quad (8)$$

- The discrepancy between  $\hat{w}(x)$  and  $w(x)$  can be bounded by

$$\begin{aligned} \mathbb{E}_{x \sim p} [(w(x) - \hat{w}(x))^2] &= \mathbb{E}_{x \sim p} \left[ \left( \frac{P(d = 1|x) - \hat{P}(d = 1|x)}{P(d = 1|x)\hat{P}(d = 1|x)} \right)^2 \right] \\ &\leq (M+1)^4 \mathbb{E}_{x \sim p} \left[ (P(d = 1|x) - \hat{P}(d = 1|x))^2 \right]. \end{aligned} \quad (9)$$

- Use  $\lambda$  ( $0 \leq \lambda \leq 1$ ) to control the bound  $M$  of the importance weights

$$T^* = \arg \min_{T, \lambda} \mathbb{E}_{x_v \sim p} [\tilde{w}(x_v) \mathcal{L}_{\text{ECE}}(\sigma(\phi(x_v)/T), y)], \quad \tilde{w}(x_v^i) = [\hat{w}(x_v^i)]^\lambda. \quad (10)$$

## Control Variate Method

- (a) Feature adaptation reduces distribution discrepancy  $d_{\alpha+1}(q||p)$

$$\begin{aligned} \text{Var}_{x \sim p} [\mathcal{L}_{\text{ECE}}^w] &= \mathbb{E}_{x \sim p} [(\mathcal{L}_{\text{ECE}}^w)^2] - (\mathbb{E}_{x \sim p} [\mathcal{L}_{\text{ECE}}^w])^2 \\ &\leq d_{\alpha+1}(q||p) (\mathbb{E}_{x \sim p} \mathcal{L}_{\text{ECE}}^w)^{1-\frac{1}{\alpha}} - (\mathbb{E}_{x \sim p} \mathcal{L}_{\text{ECE}}^w)^2, \quad \forall \alpha > 0. \end{aligned} \quad (11)$$

- (b) Control variate explicitly reduces the variance  $\sigma^2$

- Given two unbiased estimators:  $\mathbb{E}[z] = \zeta, \mathbb{E}[t] = \tau$
  - Construct a new estimator:  $z^* = z + \eta(t - \tau)$
  - $z^*$  is still unbiased:  $\mathbb{E}[z^*] = \mathbb{E}[z] + \eta\mathbb{E}[t - \tau] = \zeta + \eta(\mathbb{E}[t] - \mathbb{E}[\tau]) = \zeta$
  - $\text{Var}[z^*] = \text{Var}[z + \eta(t - \tau)] = \eta^2\text{Var}[t] + 2\eta\text{Cov}(z, t) + \text{Var}[z]$
  - $\min \text{Var}[z^*] = (1 - \rho_{z,t}^2)(\text{Var}[z], \text{ when } \hat{\eta} = -\frac{\text{Cov}(z,t)}{\text{Var}[t]})$
  - Since  $0 \leq \rho_{z,t}^2 \leq 1$ ,  $\text{Var}[z^*] \leq \text{Var}[z]$ , the variance is reduced.

<sup>2</sup>Lemieux. Control variates. In Wiley StatsRef: Statistics Reference Online, American Cancer Society, 2017.

# Transferable Calibration: Variance Reduction

- **Serial Control Variate:**  $\text{Var}[u^{**}] \leq \text{Var}[u^*] \leq \text{Var}[u]$

$$\begin{aligned} u^* &= u + \eta_1(t_1 - \tau_1) \\ u^{**} &= u^* + \eta_2(t_2 - \tau_2) \end{aligned} \tag{12}$$

- First, use importance weight  $\tilde{w}(x_s)$  as a control covariate

$$\mathbb{E}_q^*(\hat{y}, y) = \tilde{\mathbb{E}}_q(\hat{y}, y) - \frac{1}{n_s} \frac{\text{Cov}(\mathcal{L}_{\text{ECE}}^{\tilde{w}}, \tilde{w}(x))}{\text{Var}[\tilde{w}(x)]} \sum_{i=1}^{n_s} [\tilde{w}(x_s^i) - 1]. \tag{13}$$

- Second, use the prediction correctness  $r(x_s)$  as another control variate

$$\mathbb{E}_q^{**}(\hat{y}, y) = \mathbb{E}_q^*(\hat{y}, y) - \frac{1}{n_s} \frac{\text{Cov}(\mathcal{L}_{\text{ECE}}^{\tilde{w}*}, r(x))}{\text{Var}[r(x)]} \sum_{i=1}^{n_s} [r(x_s^i) - c], \tag{14}$$

- Reduce **bias, variance, and shift** all-in-one for Transferable Calibration

# TransCal Algorithm

---

## Algorithm 1 Transferable Calibration in Domain Adaptation

---

- 1: **Input:** Labeled source dataset  $\mathcal{S} = \{(\mathbf{x}_s^i, \mathbf{y}_s^i)\}_{i=1}^{n_s}$  and unlabeled target dataset  $\mathcal{T} = \{(\mathbf{x}_t^i)\}_{i=1}^{n_t}$
  - 2: **Parameter:** Temperature  $T$  and learnable meta parameter  $\lambda$
  - 3: Partition  $\mathcal{S}$  into  $\mathcal{S}_{tr} = \{(\mathbf{x}_{tr}^i, \mathbf{y}_{tr}^i)\}_{i=1}^{n_{tr}}$  and  $\mathcal{S}_v = \{(\mathbf{x}_v^i, \mathbf{y}_v^i)\}_{i=1}^{n_v}$
  - 4: Train a DA model  $\phi(\mathbf{x}) = G(F(\mathbf{x}))$  on  $\mathcal{S}_{tr}$  and  $\mathcal{T}$  via any DA method until converge
  - 5: Randomly upsample the source or the target dataset to make  $n_{tr} = n_t$
  - 6: Fix the DA model and compute features  $\mathcal{F}_{tr} = \{\mathbf{f}_{tr}^i\}_{i=1}^{n_{tr}}$ ,  $\mathcal{F}_v = \{\mathbf{f}_v^i\}_{i=1}^{n_v}$ ,  $\mathcal{F}_t = \{\mathbf{f}_t^i\}_{i=1}^{n_t}$
  - 7: Train a logistic regression model  $H$  to discriminate the features  $\mathcal{F}_{tr}$  and  $\mathcal{F}_t$  until converge
  - 8: Compute  $\hat{w}(\mathbf{x}_v^i) = [1 - H(\mathbf{f}_v^i)] / H(\mathbf{f}_v^i)$  and  $\tilde{w}(\mathbf{x}_v^i) = [\hat{w}(\mathbf{x}_v^i)]^\lambda$
  - 9: Compute  $E_{\mathbf{x} \sim p} \mathcal{L}_{ECE}^{\tilde{w}}$ ,  $\mathbb{E}_q^*(\hat{\mathbf{y}}, \mathbf{y})$  and  $\mathbb{E}_q^{**}(\hat{\mathbf{y}}, \mathbf{y})$  as in Eq. 9, Eq. 11 and Eq. 13 respectively
  - 10: Jointly optimize the transferable calibration objective as  $T^* = \arg \min_{T, \lambda} \mathbb{E}_q^{**}(\sigma(\phi(\mathbf{x}_v)/T), \mathbf{y}_v)$
  - 11: Calibrate the logit vectors on the target domain by  $\hat{\mathbf{y}}_t = \sigma(\phi(\mathbf{x}_t)/T^*)$
-

# Experiments and Results

Table 2: ECE (%) vs. Acc (%) via various calibration methods on *Office-Home* with CDAN

Metric	Cal. Method	A→C	A→P	A→R	C→A	C→P	C→R	R→A	R→C	R→P	Avg
Acc	Before Cal.	49.4	68.4	75.5	57.6	70.1	70.4	68.9	54.4	81.2	<b>68.3</b>
	MC-dropout [12]	47.2	66.2	71.4	57.1	65.7	70.6	68.3	53.6	80.7	66.7
	TransCal (ours)	49.4	68.4	75.5	57.6	70.1	70.4	68.9	54.4	81.2	<b>68.3</b>
ECE	Before Cal.	40.2	26.4	17.8	35.8	23.5	21.9	24.8	36.4	14.5	26.8
	MC-dropout [12]	33.1	21.3	15.0	24.2	20.5	13.2	25.6	14.2	22.4	19.6
	Matrix Scaling	44.7	28.8	19.7	36.1	25.4	24.1	38.1	15.7	29.5	29.1
	Vector Scaling	34.7	18.0	11.3	23.4	15.4	11.5	27.3	8.5	20.0	18.9
	Temp. Scaling	28.3	17.6	10.1	<b>21.2</b>	<u>13.2</u>	8.2	26.0	8.8	18.1	16.8
	CPCS [38]	35.0	29.4	8.3	<u>21.3</u>	29.0	<b>5.6</b>	<b>19.9</b>	9.1	20.3	19.8
	TransCal (w/o Bias)	<b>21.7</b>	<u>10.8</u>	<u>5.8</u>	27.6	<b>9.2</b>	<u>6.0</u>	27.4	5.2	<u>16.9</u>	<u>14.5</u>
	TransCal (w/o Variance)	31.2	16.4	<u>6.5</u>	31.1	14.7	<u>16.1</u>	27.5	<b>4.1</b>	20.0	18.6
	TransCal (ours)	<u>22.9</u>	<b>9.3</b>	<b>5.1</b>	21.7	14.0	6.4	<u>21.6</u>	<u>4.5</u>	<b>15.6</b>	<b>13.5</b>
	Oracle	5.8	8.1	4.8	10.0	7.7	4.2	5.5	3.9	6.2	6.2

# Experiments and Results

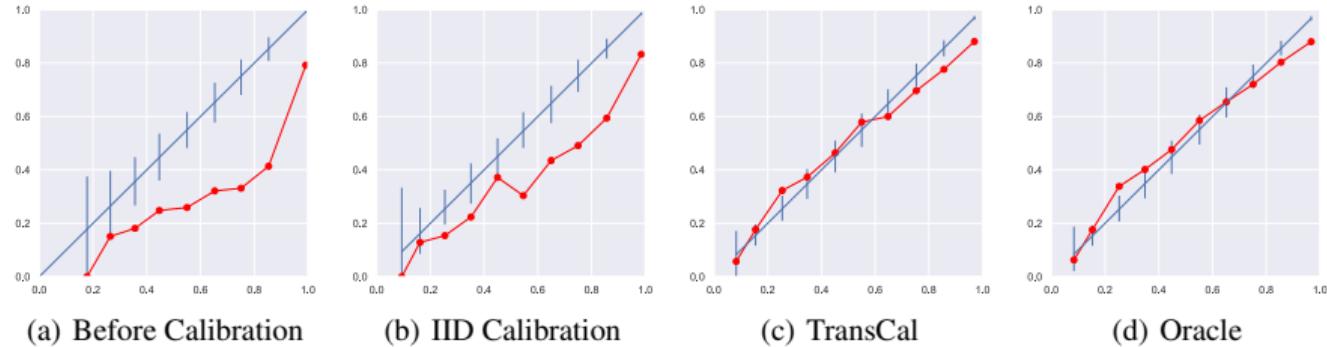


Figure 2: Reliability diagrams from *Clipart* to *Product* with CDAN [25] before and after calibration.

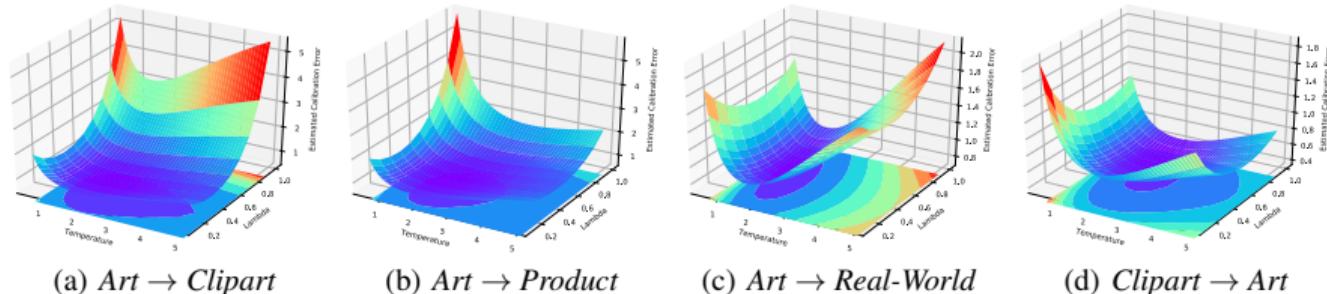


Figure 3: The estimated calibration error with respect to different values of temperature  $T$  and meta parameter  $\lambda$  (both are *learnable*), showing that different models achieve optimal values at different  $\lambda$ .

# Summary

- A dilemma in the open problem of Calibration in DA: existing domain adaptation models learn higher classification accuracy *at the expense of* well-calibrated probabilities.
- A Transferable Calibration (TransCal) method, achieving more accurate calibration with lower bias and variance in a unified hyperparameter-free optimization framework.
- Extensive experiments on various DA methods, datasets, and calibration metrics, while the effectiveness of our method has been justified both theoretically and empirically.
- Code will be available @ [github.com/thuml/TransCal](https://github.com/thuml/TransCal)

## Future Work

- 1. Design DA methods based on our ECE-Accuracy-Dilemma observation
- 2. TransCal may still fall short under the following circumstances:
  - The domain gap is extremely large even after applying domain adaptation methods
  - The source or the target dataset is too small to estimate importance weights
  - TransCal is based on the covariate shift assumption and it remains unclear whether it can still perform well under label shift, especially when we meet with a long-tailed distribution.