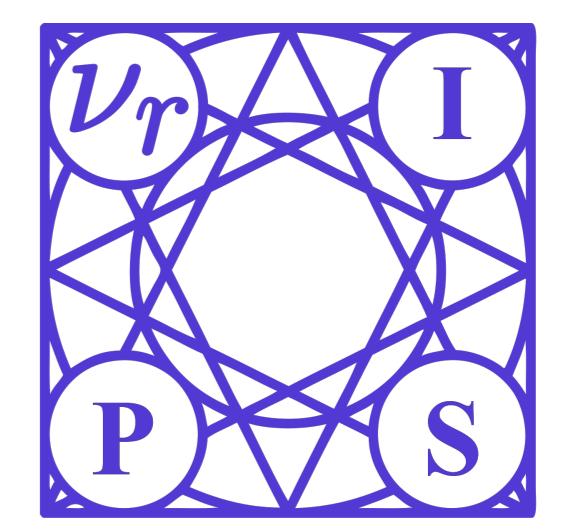




# Transferable Calibration with Lower Bias and Variance in Domain Adaptation

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



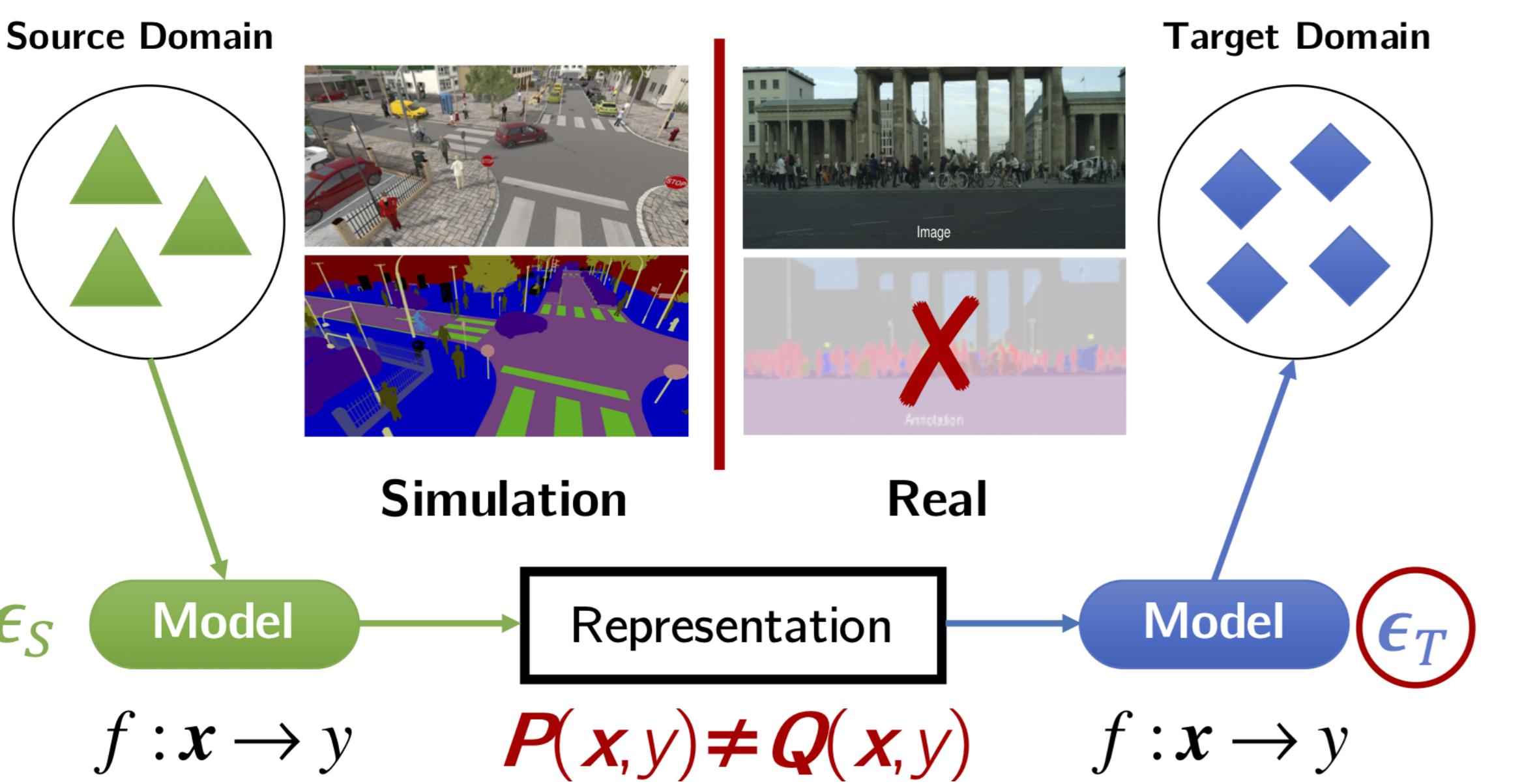
School of Software, KLiss, BNRIst, Tsinghua University <sup>#</sup>University of California, Berkeley

## Summary

- A Transferable Calibration (TransCal) method, achieving more accurate calibration with lower bias and variance in a unified hyperparameter-free optimization framework.
- 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*.
- Extensive experiments on various DA methods, datasets, and calibration metrics, while the effectiveness of our method has been justified both theoretically and empirically.
- Code available @ [github.com/thuml/TransCal](https://github.com/thuml/TransCal)

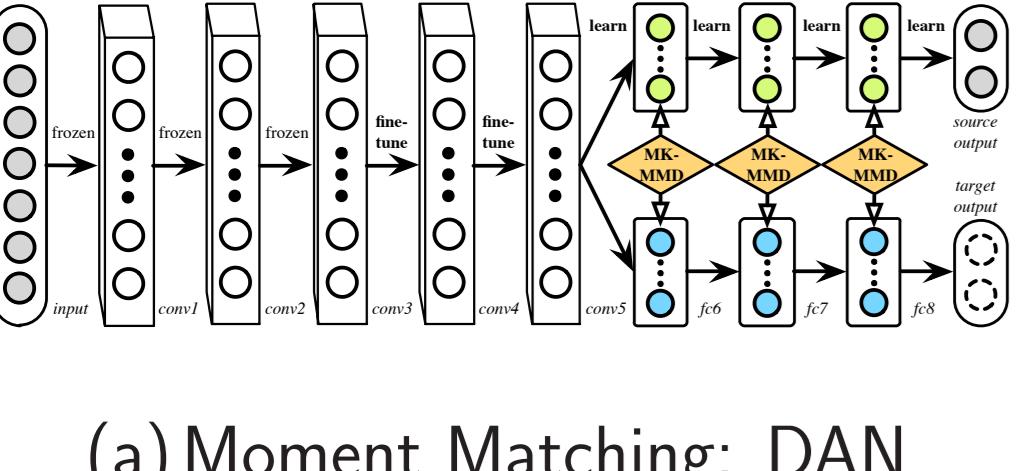
## Domain Adaptation (DA)

- Deep learning across domains:  $(P \neq Q)$
- Non independent and identically distributed distributions (Non-IID)

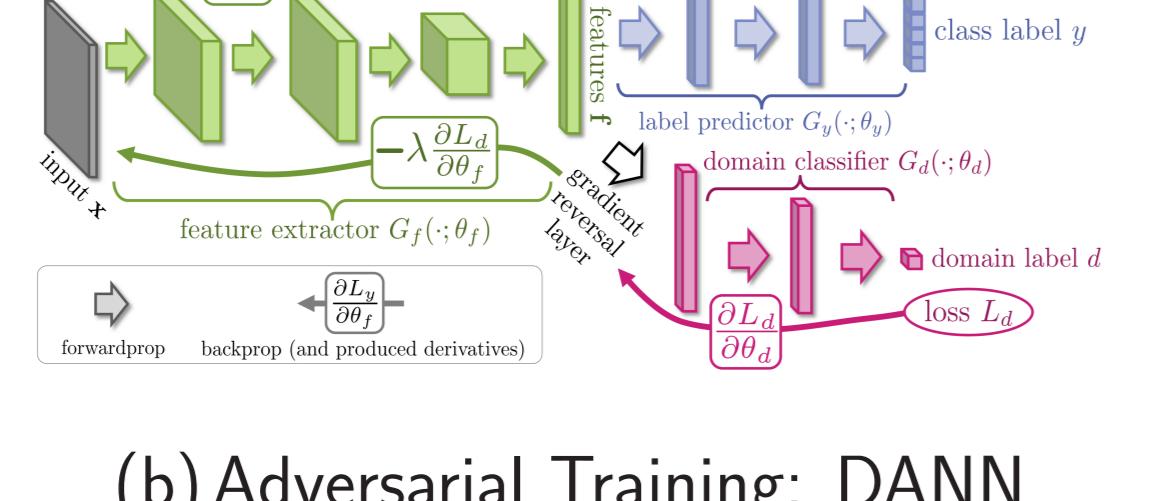


## Mainstream Approaches to DA

- Numerous deep DA methods can be mainly grouped into two categories: moment matching and adversarial training.
- Most of DA methods focus on improving the accuracy in the target domain but fail to estimate the predictive uncertainty, falling short of a miscalibration problem.



(a) Moment Matching: DAN



(b) Adversarial Training: DANN

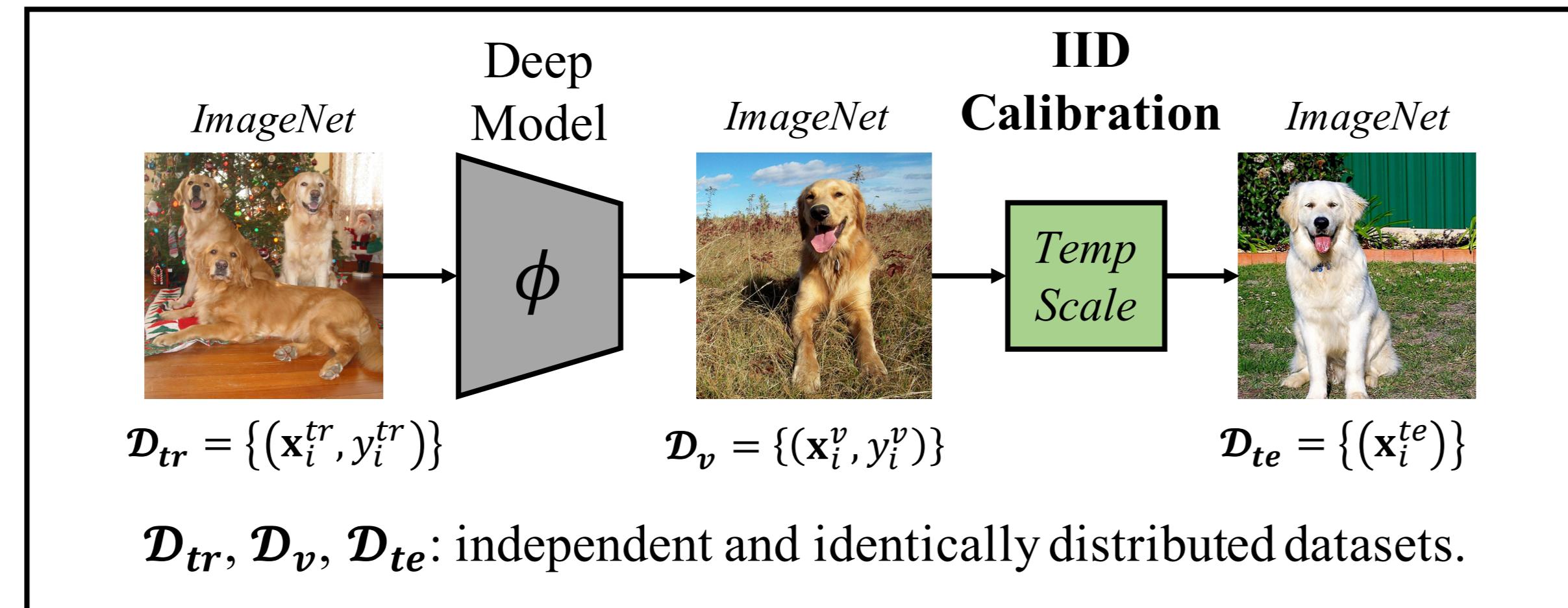
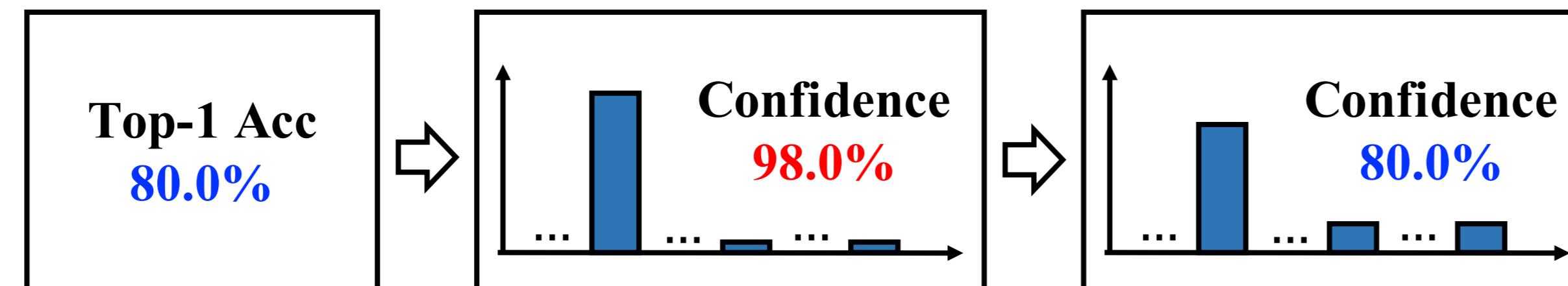
## Confidence Calibration in Deep Learning

- A model should output a prediction probability reflecting the true frequency of an event:

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

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

- DNNs learn high accuracy at the cost of **over-confidence**.



## Calibration Metric

- Expected Calibration Error (ECE)

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

## Temperature Scaling for IID Calibration

- Fix the neural model trained on the training set  $\mathcal{D}_{tr}$
- Attain the optimal temperature  $T^*$  by minimizing

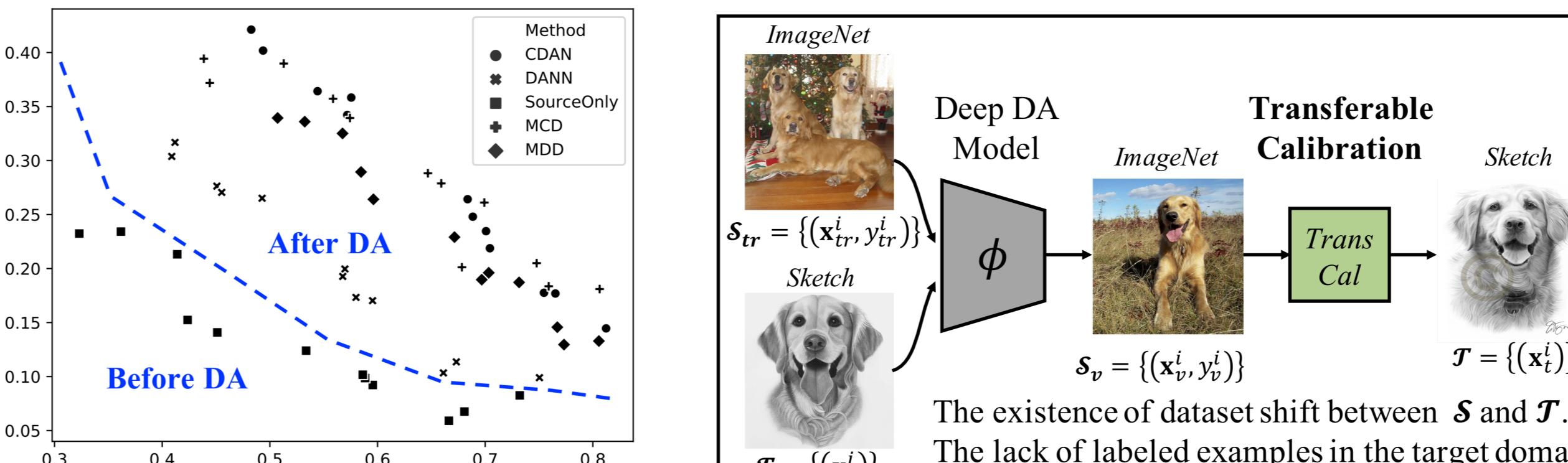
$$T^* = \arg \min_T \mathbb{E}_{(x_v, y_v) \in \mathcal{D}_v} \mathcal{L}_{\text{NLL}}(\sigma(z_v / T), y_v) \quad (3)$$

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

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

## Dilemma of Accuracy vs Confidence in DA

- DA models yield high acc at the cost of **poorly-calibration**.



- Calibration in DA is challenging due to the existence of domain shift and the lack of target label

## Transferable Calibration Framework

- Estimate the target ECE by importance weighting

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

- Estimate density ratio from a logistic regression classifier

$$\widehat{w}(x) = \frac{q(x)}{p(x)} = \frac{v(x|d=0)}{v(x|d=1)} = \frac{P(d=1|x)P(d=0|x)}{P(d=0|x)P(d=1|x)}, \quad (5)$$

## Transferable Calibration: Bias Reduction

- Bias between the estimated ECE and the ground-truth one

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

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

$$\mathbb{E}_{x \sim p} [(w(x) - \widehat{w}(x))^2] \leq (M+1)^4 \mathbb{E}_{x \sim p} \left[ \left( P(d=1|x) - \widehat{P}(d=1|x) \right)^2 \right]. \quad (7)$$

- Use  $\lambda$  ( $0 \leq \lambda \leq 1$ ) to control the bound  $M$  of  $\widehat{w}(x)$

$$T^* = \arg \min_{T, \lambda} \mathbb{E}_{x_v \sim p} [\widetilde{w}(x_v) \mathcal{L}_{\text{ECE}}(\sigma(\phi(x_v) / T), y)], \quad \widetilde{w}(x_v) = [\widehat{w}(x_v)]^\lambda. \quad (8)$$

## Transferable Calibration: Variance Reduction

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

$$\begin{aligned} u^* &= u + \eta_1(t_1 - \tau_1) \\ u^{**} &= u^* + \eta_2(t_2 - \tau_2) \end{aligned} \quad (9)$$

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

$$\mathbb{E}_q^*(\hat{y}, y) = \widetilde{\mathbb{E}}_q(\hat{y}, y) - \frac{1}{n_s} \frac{\text{Cov}(\mathcal{L}_{\text{ECE}}, \tilde{w}(x))}{\text{Var}[\tilde{w}(x)]} \sum_{i=1}^{n_s} [\tilde{w}(x_s^i) - 1]. \quad (10)$$

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

$$\mathbb{E}_q^{**}(\hat{y}, y) = \mathbb{E}_q^*(\hat{y}, y) - \frac{1}{n_s} \frac{\text{Cov}(\mathcal{L}_{\text{ECE}}^{\tilde{w}^*}, r(x))}{\text{Var}[r(x)]} \sum_{i=1}^{n_s} [r(x_s^i) - c], \quad (11)$$

## 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>	13.2	8.2	26.0	8.8	18.1	16.8
	CPCS [38]	35.0	29.4	8.3	21.3	29.0	<b>5.6</b>	<b>19.9</b>	9.1	20.3	19.8
TransCal (w/o Bias)	<b>21.7</b>	<b>10.8</b>	<b>5.8</b>	27.6	<b>9.2</b>	6.0	27.4	5.2	16.9	<b>14.5</b>	
	TransCal (w/o Variance)	31.2	16.4	6.5	31.1	14.7	16.1	27.5	<b>4.1</b>	20.0	18.6
	TransCal (ours)	<b>22.9</b>	<b>9.3</b>	<b>5.1</b>	21.7	14.0	6.4	21.6	<b>4.5</b>	<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	6.2

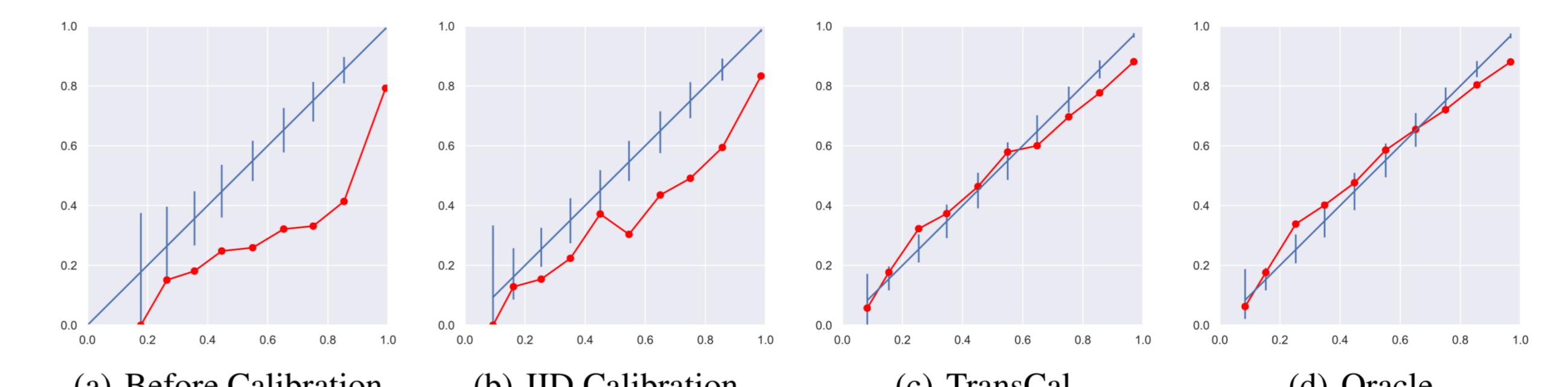


Figure 2: Reliability diagrams from *Clipart* to *Product* with CDAN [30] before and after calibration.

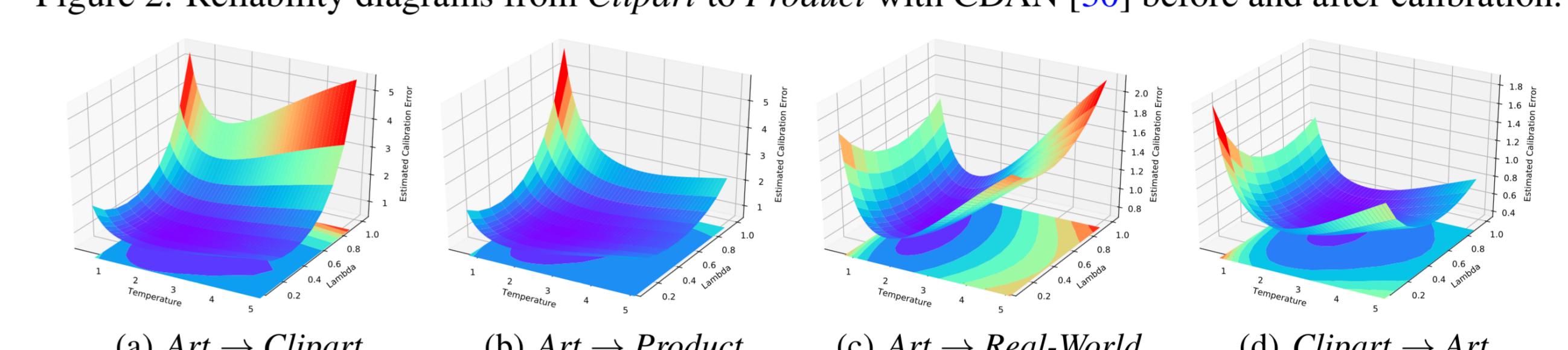


Figure 3: The estimated calibration error with respect to different values of temperature  $T$  and meta parameter  $\lambda$  (both are *learnable*), showing that different models achieve optimal values at different  $\lambda$ .