浙江大学
ZHEJIANG UNIVERSITY

# KD3A: Unsupervised Multi-Source Decentralized Domain Adaptation via Knowledge Distillation

**Hao-zhe Feng**[1✉], Zhaoyang You[2], Minghao Chen[1], Tianye Zhang[1], Minfeng Zhu[1], Fei Wu[2], Chao Wu[3], Wei Chen[1*]

1 State Key Lab of CAD&CG, Zhejiang University, Hangzhou, China

2 College of Computer Science and Technology, Zhejiang University, Hangzhou, China
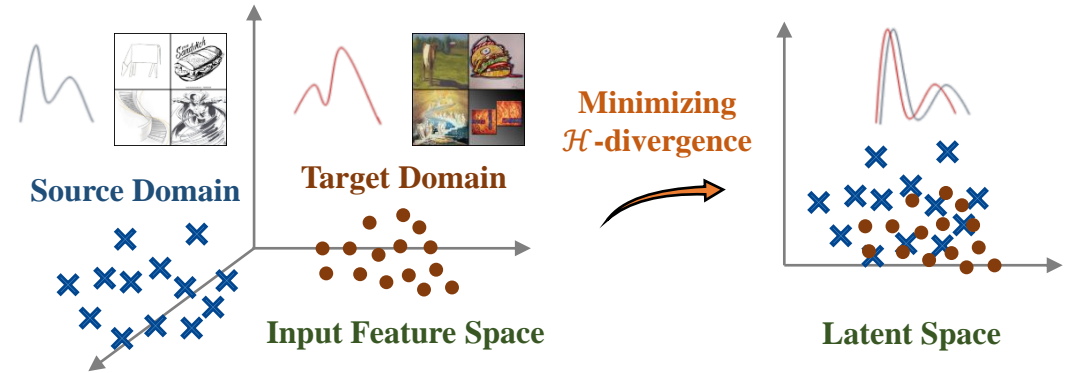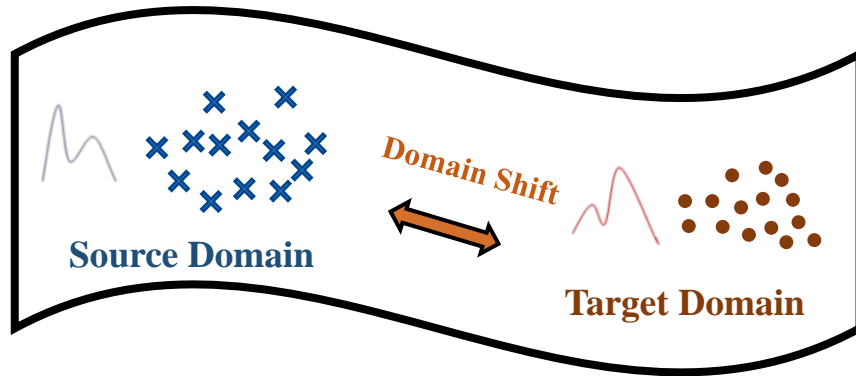
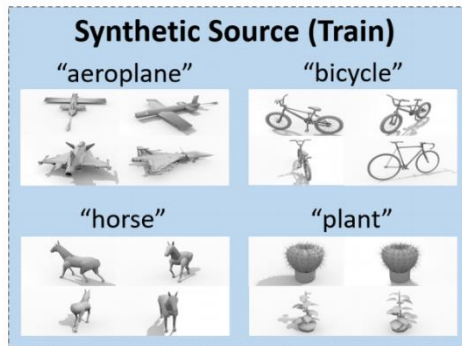3 School of Public Affairs, Zhejiang University, Hangzhou, China

✉fenghz@zju.edu.cn

*Corresponding author

Visual Analytics Group | State Key Lab of CAD&CG, Zhejiang University

# Unsupervised Multi-Source Domain Adaptation (UMDA)



Domain Shift

Source Domain

Target Domain



Source Domain

Target Domain

Minimizing $\mathcal{H}$-divergence

Input Feature Space

Latent Space

*[The prevailing unsupervised multi-source domain adaptation (UMDA) paradigm.]*



**Synthetic Source (Train)**

"aeroplane"       "bicycle"

"horse"       "plant"

**Real Target (Test)**

**Synthesized Data: Cheap and Abundant**

**Real Data: Expensive and Scarce**

*[Peng X et al. VisDA: The Visual Domain Adaptation Challenge, 2017.]*

**Theorem 1** *Let $\mathcal{H}$ be the model space, $\{\epsilon_{\mathbb{D}_S^k}(h)\}_{k=1}^K$ and $\epsilon_{\mathbb{D}_T}(h)$ be the task risks of source domains $\{\mathbb{D}_S^k\}_{k=1}^K$ and the target domain $\mathbb{D}_T$, and $\boldsymbol{\alpha} \in \mathcal{R}_+^K, \sum_{k=1}^K \boldsymbol{\alpha}_k = 1$ be the domain weights. Then for all $h \in \mathcal{H}$ we have:*

Task loss     $\mathcal{H}$-divergence

$$\epsilon_{\mathbb{D}_T}(h) \leq \sum_{k=1}^K \boldsymbol{\alpha}_k \left( \epsilon_{\mathbb{D}_S^k}(h) + d_{\mathcal{H}}(\mathbb{D}_S^k, \mathbb{D}_T) \right) + \lambda_0 \quad (3)$$

where $\lambda_0 = \min_{h \in \mathcal{H}} \sum_{k=1}^K \alpha_k \epsilon_{\mathbb{D}_S^k}(h) + \epsilon_{\mathbb{D}_T}(h)$ is a constant according to the task risk of the optimal model on the source domains and target domain.
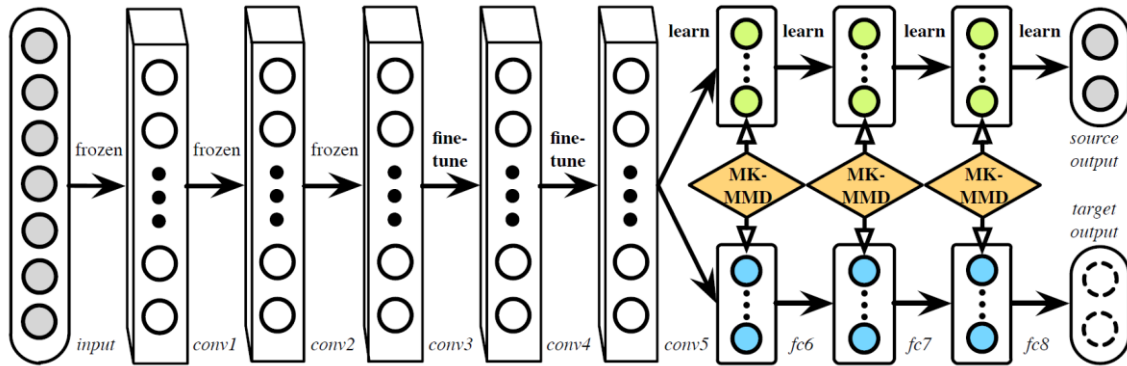
*[The Generalization Bound of UMDA.]*

# Unsupervised Multi-Source Domain Adaptation (UMDA)



*[Ganin Y et al. Unsupervised Domain Adaptation by Backpropagation (DANN), ICML 2015.]*



*[Long M et al. Learning Transferable Features with Deep Adaptation Networks, ICML 2015.]*

## How to optimize the $\mathcal{H}$-divergence?

- **Adversarial Learning**

  <span style="color:red">Making features from different domains undistinguished</span>

  1. Build domain classifier $h \in \mathcal{H}$ with:

     $$d_{\mathcal{H}}(\mathbb{D}_S, \mathbb{D}_T) = \sup_{h \in \mathcal{H}} |\mathbb{E}_{\mathbf{X}_i \sim \mathbb{D}_S} \mathbf{I}[h(\mathbf{X}_i) = 1] - \mathbb{E}_{\mathbf{X}_i \sim \mathbb{D}_T} \mathbf{I}[h(\mathbf{X}_i) = 1]|$$

  2. Seek the classifier $G_d(\cdot, \boldsymbol{\theta_d})$ that minimizing the loss of domain classifier: $\lambda \frac{\partial L_d}{\partial \theta_d}$

  3. Seek the feature mapping $G_f(\cdot, \boldsymbol{\theta_f})$ that maximizing the loss of domain classifier: $-\lambda \frac{\partial L_d}{\partial \theta_f}$

- **Maximum Mean Discrepancy (MMD)**

  1. Build reproducing kernel Hilbert space (RKHS) $\mathcal{H}$ with kernel $\kappa$:
     $h(\boldsymbol{X}) = \langle h(\cdot), \kappa(\cdot, \boldsymbol{X}) \rangle$

  2. Write the equivalence $\mathcal{H}$-divergence as

     $$d_{\mathcal{H}}(\mathbb{D}_S, \mathbb{D}_T) = \sup_{h \in \mathcal{H}} \langle h, \int_{\mathbf{X}} \kappa(\cdot, \mathbf{X}) d\mathbb{P}_S(\mathbf{X}) \rangle - \langle h, \int_{\mathbf{X}} \kappa(\cdot, \mathbf{X}) d\mathbb{P}_T(\mathbf{X}) \rangle$$

  3. Minimize kernel MMD

     $$d_{\mathrm{MMD}}^{\kappa}(\mathbb{D}_S, \mathbb{D}_T) = -2\mathbb{E}_{\mathbf{X}_S, \mathbf{X}_T \sim \mathbb{D}_S, \mathbb{D}_T} \kappa(\mathbf{X}_S, \mathbf{X}_T)$$
     $$+ \mathbb{E}_{\mathbf{X}_S, \mathbf{X}_S' \sim \mathbb{D}_S} \kappa(\mathbf{X}_S, \mathbf{X}_S') + \mathbb{E}_{\mathbf{X}_T, \mathbf{X}_T' \sim \mathbb{D}_T} \kappa(\mathbf{X}_T, \mathbf{X}_T')$$

# Domain Adaptation with Privacy-preserving Policy

## General Data Protection Regulation (GDPR)

- **Individual Rights**

Simply combining source data is forbidden.

1. **Right to be informed and accessible.** Individuals have the right to be informed about the collection and use of their personal data, such as the purposes of data analysis, with whom the data is shared, and the data collecting methods. **If the data is obtained from other sources, the privacy information must be provided within one month**. Individuals also have the right to access and receive a copy of their personal data and other supplementary information

2. **Right to object.** Individuals have the right to object **the directly usage in market**.

3. **Right to rectification and erasure.** Individuals have the right to rectify inaccurate personal data. They also have rights to 'forget' the influence of their data, such as to be total deleted in databases and to **control the impact in data-mining models**.

*[Individual rights in general data protection regulation,*
*https://ico.org.uk/for-organisations/guide-to-data-*
*protection/guide-to-the-general-data-protection-*
*regulation-gdpr/individual-rights.]*

4. **Right to restrict processing and automated decision making.** Individuals have the right to restrict the processing their personal data. **The automated profiling about evaluating certain things of individuals** or making automated decisions about individuals without human involvement **must be authorized**.

# Domain Adaptation with Privacy-preserving Policy
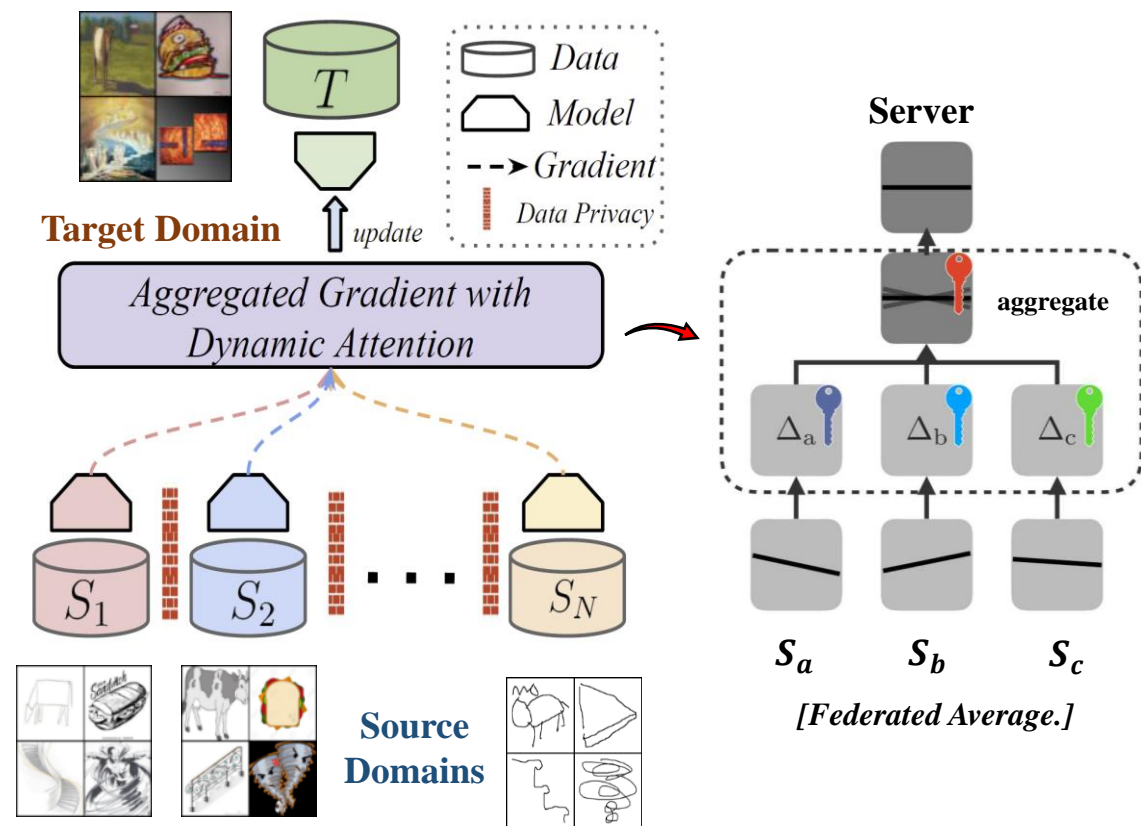
## Data Residency Laws and Federated Learning

- ## Corporate Responsibilities

Data localization or residency refers to the data location's physical storage within a country's national boundaries. For European Union**, Companies have to keep the data inside the EU**. Data can only be conditionally transferred to countries and organizations that have signed up to equivalent privacy protection as EU.

Regulated data types: **Profile; Finance; Employment; Health, Payment.**

- ## Solutions: Federated Learning

Federated learning (FL) is a machine learning setting where many clients (e.g. mobile devices or whole organizations) **collaboratively train a model** under the orchestration of a central server (e.g. service provider), **while keeping the training data decentralized**.



*Data*
*Model*
*Gradient*
*Data Privacy*

**Target Domain**

*Aggregated Gradient with Dynamic Attention*

**Source Domains**

*[Peng X, et al. Federated adversarial domain adaptation (FADA), ICLR 2020.]*

**Server**

aggregate

$S_a$   $S_b$   $S_c$

*[Federated Average.]*

## Problem Formulation for Decentralized Scenario

1. **All the data and computations** on source domains must be kept **localized.**

2. **Available information** in each communication round**:**

- The size of the training sets $\{N_S^k\}_{k=1}^K$ on source domains.
- The parameters of K models $\{h_S^k\}_{k=1}^K$ trained on source domains.
- The target domain data containing $N_T$ unlabeled examples $\mathbb{D}_T \coloneqq \{X_i^T\}_{i=1}^{N_T}$.

## Challenges for Decentralized UMDA

1. **Minimizing the $\mathcal{H}$-divergence requires pairwise calculation of data.**

2. **The communication cost and privacy security.**

3. **The negative transfer problem.**



*[DomainNet Dataset.]*



*[Peng X, et al. Federated adversarial domain adaptation (FADA), ICLR 2020.]*
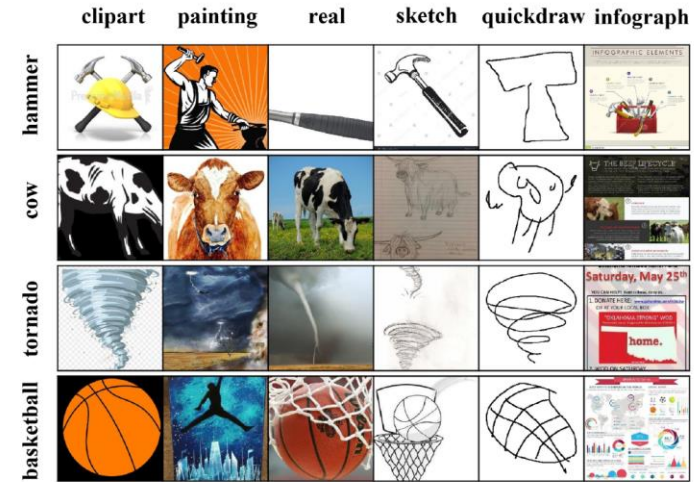
# Domain Adaptation with Privacy-preserving Policy

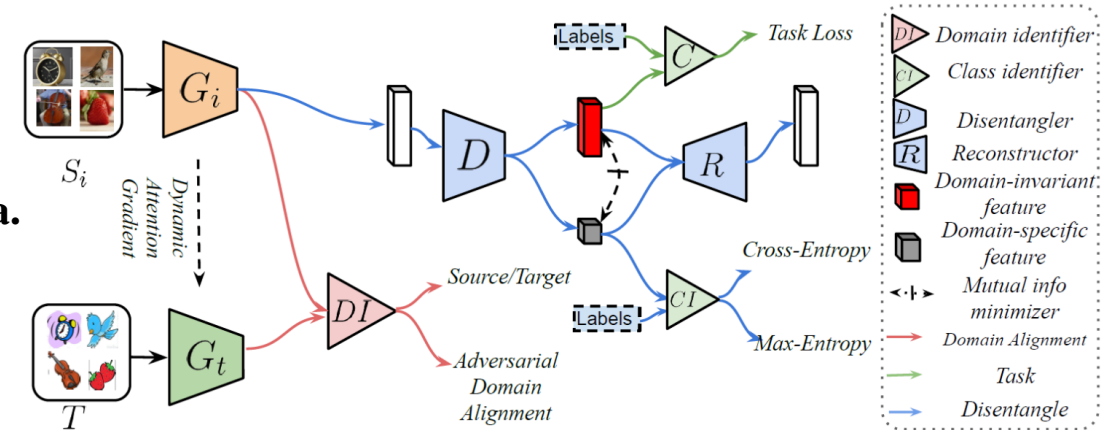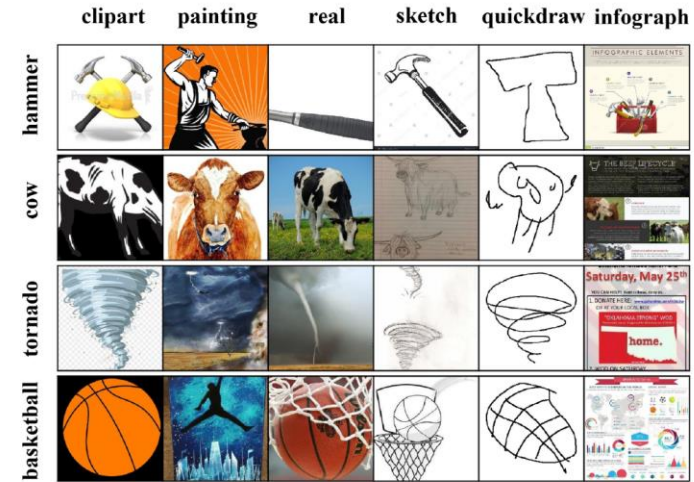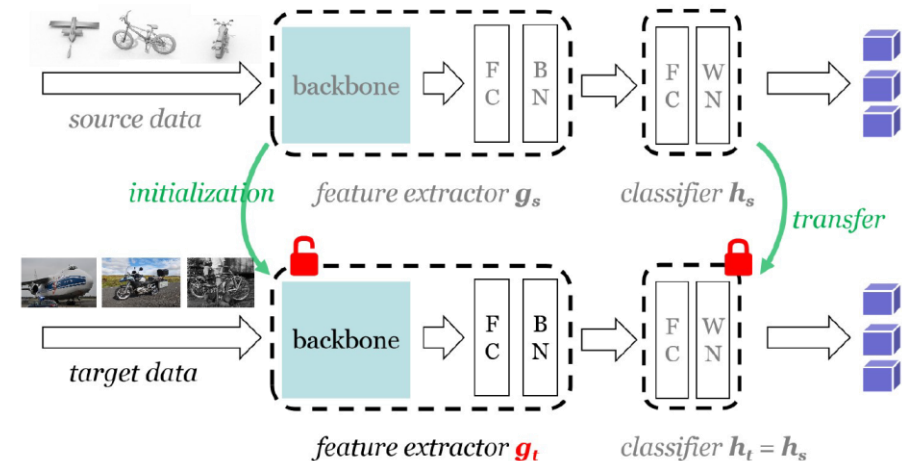## Problem Formulation for Decentralized Scenario

1. **All the data and computations** on source domains must be kept **localized.**

2. **Available information** in each communication round:

- The size of the training sets $\{N_S^k\}_{k=1}^K$ on source domains.
- The parameters of K models $\{h_S^k\}_{k=1}^K$ trained on source domains.
- The target domain data containing $N_T$ unlabeled examples $\mathbb{D}_T := \{X_i^T\}_{i=1}^{N_T}$.

## Challenges for Decentralized UMDA

1. **Minimizing the $\mathcal{H}$-divergence requires pairwise calculation of data.**

2. **The communication cost and privacy security.**
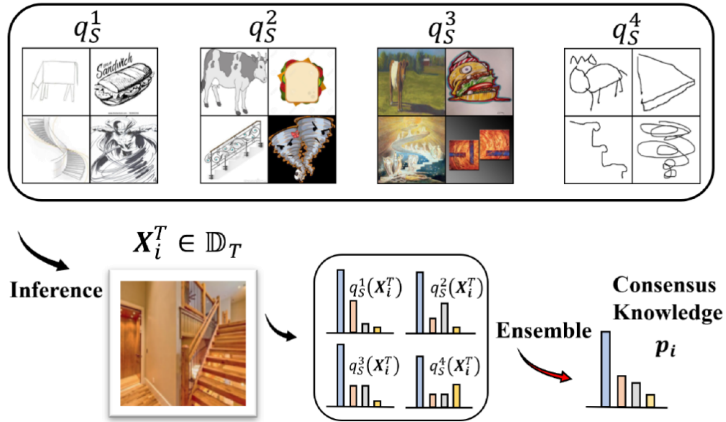
3. **The negative transfer problem.**
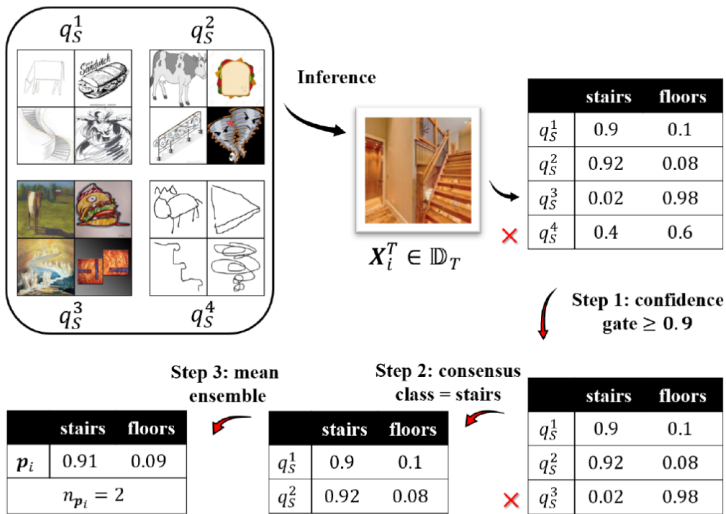


*[DomainNet Dataset.]*



*[Liang J, et al. Do we really need to access the source data? Source hypothesis transfer for unsupervised domain adaptation (SHOT), ICML 2020.]*

# Our Decentralized UMDA Paradigm: KD3A



(a) Knowledge distillation process in UMDA.



(b) Knowledge vote ensemble.

## KD3A: Using Three Components in Tandem to Solve the Decentralized UMDA Challenges

1. **Knowledge vote: producing high-quality domain consensus on $\mathbb{D}_T$**

2. **Consensus focus: against negative transfer**

3. **BatchNorm MMD: decentralized optimization of $\mathcal{H}$-divergence**

## Knowledge Distillation: Extending Source Domains with Consensus Knowledge

1. Get an extended source domain $\mathbb{D}_S^{K+1}$ as $\mathbb{D}_S^{K+1} = \{(\mathbf{X}_i^T, \mathbf{p}_i)\}_{i=1}^{N_T}$

2. Train the extended source model $h_S^{K+1}$ through knowledge distillation loss as $L^{\text{kd}}(\mathbf{X}_i^T, q_S^{K+1}) = D_{\text{KL}}(\mathbf{p}_i \| q_S^{K+1}(\mathbf{X}_i^T))$.

3. Target model is the aggregation of source models as $h_T := \sum_{k=1}^{K+1} \alpha_k \, h_S^k$

4. The related task risk for $\mathbb{D}_S^{K+1}$: $\epsilon_{\mathbb{D}_S^{K+1}}(h) = \Pr_{(\mathbf{X},\mathbf{p}) \sim \mathbb{D}_S^{K+1}}[h(\mathbf{X}) \neq \arg_c \max \mathbf{p}_c]$.

(a) Knowledge distillation process in UMDA.



(b) Knowledge vote ensemble.

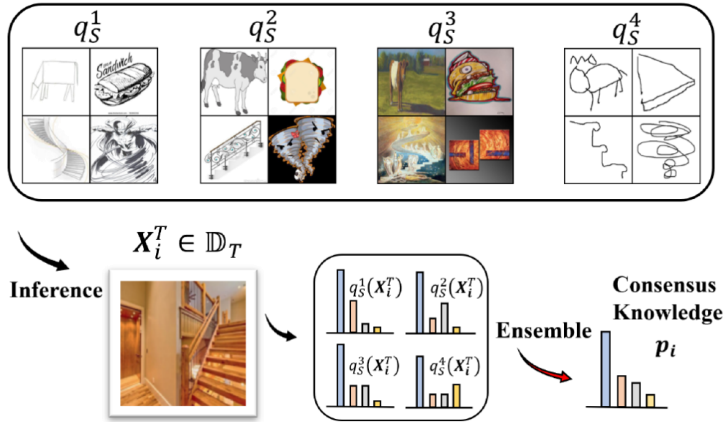## KD3A: Using Three Components in Tandem to Solve the Decentralized UMDA Challenges

1. **Knowledge vote: producing high-quality domain consensus on $\mathbb{D}_T$**

2. **Consensus focus: against negative transfer**

3. **BatchNorm MMD: decentralized optimization of $\mathcal{H}$-divergence**

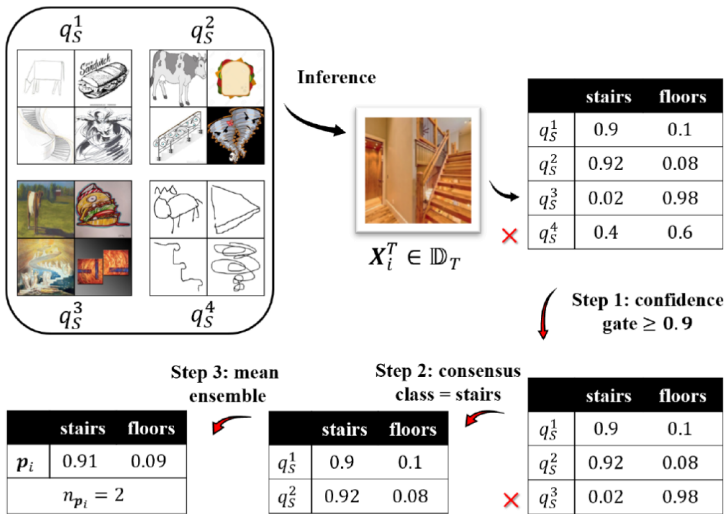## Knowledge Distillation: Extending Source Domains with Consensus Knowledge

**Proposition 1** (*The generalization bound for knowledge distillation*). *Let $\mathcal{H}$ be the model space and $\epsilon_{\mathbb{D}_S^{K+1}}(h)$ be the task risk of the new source domain $\mathbb{D}_S^{K+1}$ based on knowledge distillation. Then for all $h_T \in \mathcal{H}$, we have:*

$$\epsilon_{\mathbb{D}_T}(h_T) \leq \epsilon_{\mathbb{D}_S^{K+1}}(h_T) + \frac{1}{2}d_{\mathcal{H}\Delta\mathcal{H}}(\mathbb{D}_S^{K+1}, \mathbb{D}_T)$$
$$+ \min\{\lambda_1, \sup_{h\in\mathcal{H}} |\epsilon_{\mathbb{D}_S^{K+1}}(h) - \epsilon_{\mathbb{D}_T}(h)|\}$$

(4)

Notice: the new domain $\mathbb{D}_S^{K+1}$ can improve the original bound if the consensus knowledge is good enough to represent the ground-truth label, that is:

$$\sup_{h\in\mathcal{H}} |\epsilon_{\mathbb{D}_T}(h) - \epsilon_{\mathbb{D}_S^{K+1}}(h)| \leq \lambda_1.$$

(a) Knowledge distillation process in UMDA.



(b) Knowledge vote ensemble.

## Knowledge Vote: Producing Good Consensus

**Main motivation:** if a certain consensus knowledge is supported by more source domains with high confidence (e.g., > 0.9), it will be more likely to be the true label.

**Three steps:**

1. **Confidence gate.** For each $X_i^T \in \mathbb{D}_T$, we use a high-level confidence gate to filter the predictions $\{q_S^k(X_i^T)\}_{k=1}^K$ of teacher models and eliminate the unconfident models.

2. **Consensus class vote.** For the remained models, the predictions are added up to find the consensus class with the maximum accumulated confidence. Then the inconsistent models are dropped.

3. **Mean ensemble** is conducted after the class vote to get the consensus knowledge $p_i$. The number of domains that support $p_i$ is also recorded as $n_{p_i}$. For those $X^T$ with all teacher models eliminated by confidence gate, the naïve mean ensemble are conducted on all teachers to get $p$ and a low weight is assigned to $X^T$ as $n_p = 0.001$.

Knowledge vote builds $\mathbb{D}_S^{K+1} = \{(X_i^T, p_i, n_{p_i})\}_{i=1}^{N_T}$ with the **loss objective:**

$$L^{\text{kv}}(\mathbf{X}_i^T, q) = n_{\mathbf{p}_i} \cdot D_{\text{KL}}(\mathbf{p}_i \| q(\mathbf{X}_i^T))$$

# Our Decentralized UMDA Paradigm: KD3A

## Consensus Focus: Against Negative Transfer

**Domain weight $\alpha$** determine the contribution of each source domain. To get better UMDA performance, we should assign the low weights to bad source domains.

- **Previous methods:** utilizing $\mathcal{H}$-divergence to re-weight source domains and identify irrelevant sources.

$$\alpha_k = N_k e^{-d_{\mathcal{H}}(\mathbb{D}_S^k, \mathbb{D}_T)} / \sum_k N_k e^{-d_{\mathcal{H}}(\mathbb{D}_S^k, \mathbb{D}_T)}.$$

Notice: requires pairwise calculation. Can not identify malicious domain!

- **Our motivation:** measure the quality of consensus from knowledge vote, assign high weights to those domains which provide high-quality consensus and penalize those domains which provide bad consensus.



Aggregated Gradient with Dynamic Attention

$\alpha_1$   $\alpha_2$   $\alpha_K$

$S_1$   $S_2$   $\cdots$   $S_N$

**Source Domains**

*[Aggregate Source Models with Domain Weight.]*

1. **Consensus quality**: if one consensus class is supported by more source domains with higher confidence, then it will be more likely to represent the true label, which leads to higher consensus quality.

$$\mathcal{S} = \{\mathbb{D}_S^k\}_{k=1}^K, \forall \, \mathcal{S}' \subset \mathcal{S}, \quad Q(\mathcal{S}') = \sum_{\mathbf{X}_i^T \in \mathbb{D}_T} n_{\mathbf{p}_i}(\mathcal{S}') \cdot \max \mathbf{p}_i(\mathcal{S}')$$

Total consensus quality.   Marginal contribution for $\mathbb{D}_S^k$.

2. **Domain quality**: describing the marginal contribution to the total consensus quality of each source domain as $\mathrm{CF}(\mathbb{D}_S^k) = Q(\mathcal{S}) - Q(\mathcal{S} \setminus \{\mathbb{D}_S^k\})$

3. **Re-weighting strategy from consensus focus**: $\alpha_{K+1} = \frac{N_T}{\sum_{k=1}^K N_k + N_T}, \quad \alpha_k^{\mathrm{CF}} = (1 - \alpha_{K+1}) \cdot \frac{N_k \cdot \mathrm{CF}(\mathbb{D}_S^k)}{\sum_{k=1}^K N_k \cdot \mathrm{CF}(\mathbb{D}_S^k)}$

Does not need to access original data. Can identify malicious domain!

## BatchNorm MMD: Decentralized Optimization of $\mathcal{H}$-divergence

**Minimizing the $\mathcal{H}$-divergence** can optimize the UMDA upper bound.

- **Previous methods:** optimizing $\mathcal{H}$-divergence with the kernel-based MMD.

$$\min_{h \in \mathcal{H}} \sum_{k=1}^{K+1} \alpha_k d_{\text{MMD}}^\kappa (\mathbb{D}_S^k, \mathbb{D}_T)$$

<span style="color:red">Need to access original data.</span>

- **Our motivation:** utilizes the mean and variance parameters in each BatchNorm layer to optimize the $\mathcal{H}$-divergence without accessing data.

**Theorem 1** *Let $\mathcal{H}$ be the model space, $\{\epsilon_{\mathbb{D}_S^k}(h)\}_{k=1}^K$ and $\epsilon_{\mathbb{D}_T}(h)$ be the task risks of source domains $\{\mathbb{D}_S^k\}_{k=1}^K$ and the target domain $\mathbb{D}_T$, and $\alpha \in \mathcal{R}_+^K$, $\sum_{k=1}^K \alpha_k = 1$ be the domain weights. Then for all $h \in \mathcal{H}$ we have:*

<span style="color:red">Task loss</span>   <span style="color:red">$\mathcal{H}$-divergence</span>

$$\epsilon_{\mathbb{D}_T}(h) \leq \sum_{k=1}^K \alpha_k \left( \epsilon_{\mathbb{D}_S^k}(h) + d_{\mathcal{H}}(\mathbb{D}_S^k, \mathbb{D}_T) \right) + \lambda_0 \quad (3)$$

*[The Generalization Bound of UMDA.]*

1. **Build Kernel-MMD** with batch-normalized feature $\pi_l$ and quartic kernel $\kappa(X^S, X^T) = \left( \langle \pi_l^S, \pi_l^T \rangle + \frac{1}{2} \right)^2$:

$$d_{\text{MMD}}(\mathbb{D}_S^k, \mathbb{D}_T) = \|\mathbb{E}(\pi_l^k) - \mathbb{E}(\pi_l^T)\|_2^2 + \|\mathbb{E}[\pi_l^k]^2 - \mathbb{E}[\pi_l^T]^2\|_2^2$$

2. **Obtaining** $\{\mathbb{E}(\pi_l^k), \text{Var}(\pi_l^k)\}_{l=1}^L$ from $L$ BatchNorm layers of the $K+1$ source domain models : $\mathbb{E}[\pi]^2 = Var(\pi) + [\mathbb{E}(\pi)]^2$

3. **Training** target model $h_T$:

$$\sum_{l=1}^L \sum_{k=1}^{K+1} \alpha_k \left( \|\mu(\pi_l^T) - \mathbb{E}(\pi_l^k)\|_2^2 + \|\mu[\pi_l^T]^2 - \mathbb{E}[\pi_l^k]^2\|_2^2 \right)$$

<span style="color:red">**Tips:**</span> directly optimizing this loss requires to traverse all BatchNorm layers, which is <span style="color:red">time-consuming</span>. Instead, we use a EM-liked method. First, getting the <span style="color:red">global optimal solution</span> as $\mu_{op}(\pi_l^T) = \Sigma_{k=1}^{K+1} \alpha_k \mathbb{E}(\pi_l^k)$, $\mu_{op}[\pi_l^T]^2 = \Sigma_{k=1}^{K+1} \alpha_k \mathbb{E}[\pi_l^k]^2$. Then, <span style="color:red">directly substituting</span> the solution into the target model. This heuristic method works well in practice.

# Our Decentralized UMDA Paradigm: KD3A

## Algorithm of KD3A

**Algorithm 1** KD3A training process with epoch t.

**Input:**

    Source domains $\mathcal{S} = \{\mathbb{D}_S^k\}_{k=1}^K$. Target domain $\mathbb{D}_T$;

    Target model $h_T^{(t-1)}$ with parameters $\Theta^{(t-1)}$;

    Confidence gate $g^{(t)}$;

**Output:**

    Target model $h_T^{(t)}$ with parameters $\Theta^{(t)}$.

1: // Locally training on source domains:
2: **for** $\mathbb{D}_S^k$ in $\mathcal{S}$ **do**
3:    Model initialize: $(h_S^k, \Theta_S^k) \leftarrow (h^{(t-1)}, \Theta^{(t-1)})$.
4:    Train $h_S^k$ with classification loss on $\mathbb{D}_S^k$.
5: **end for**
6: Upload $\{(h_S^k, \Theta_S^k)\}_{k=1}^K$ to the target domain.
7: // Knowledge Vote:
8: $\mathbb{D}_S^{K+1} \leftarrow$ KnowledgeVote$(\mathbb{D}_T, g^{(t)}, \{h_S^k\}_{k=1}^K)$.
9: Train $h_S^{K+1}$ with $L^{kv}$ loss (5) on $\mathbb{D}_S^{K+1}$.
10: // Consensus Focus:
11: $\boldsymbol{\alpha}^{CF} \leftarrow$ ConsensusFocus$(\mathbb{D}_T, \{h_S^k\}_{k=1}^K, \{N_k\}_{k=1}^K)$.
12: // Model Aggregation:
13: $\Theta^{(t)} \leftarrow \sum_{k=1}^{K+1} \boldsymbol{\alpha}_k^{CF} \cdot \Theta_S^k$.
14: // BatchNorm MMD:
15: Obtain $\{\mathbb{E}[\boldsymbol{\pi}_l^k]^i, i = 1, 2\}_{l,k=1}^{L,K+1}$ from $\{(h_S^k, \Theta_S^k)\}_{k=1}^{K+1}$.
16: Train $h_T^{(t)}$ with BatchNorm MMD on $\mathbb{D}_T$.
17: Return $(h_T^{(t)}, \Theta^{(t)})$.

**Step 1** (lines 2-6)
**Step 2** (lines 8-13)
**Step 3** (lines 15-16)

## Generalization Bound for KD3A

**Theorem 2** (*The decentralized generalization bound for KD3A*). *Let $h_T$ be the target model of KD3A, $\{\mathbb{D}_S^k\}_{k=1}^{K+1}$ be the extended source domains through Knowledge Vote and $\boldsymbol{\alpha}^{CF} \in \mathcal{R}_+^{K+1}, \sum_{k=1}^{K+1} \alpha_k^{CF} = 1$ be the domain weights through Consensus Focus. Then we have:*

$$\epsilon_{\mathbb{D}_T}(h_T) \leq \sum_{k=1}^{K+1} \alpha_k^{CF}\left(\epsilon_{\mathbb{D}_S^k}(h_T) + \frac{1}{2}d_{\mathcal{H}\Delta\mathcal{H}}(\mathbb{D}_S^k, \mathbb{D}_T)\right) + \lambda_2$$

$$(13)$$

The generalization performance of KD3A bound (13) depends on the quality of the consensus knowledge, as the following proposition shows (see Appendix C for proof):

**Proposition 2** *The KD3A bound (13) is a tighter bound than the original bound (2), if the task risk gap between the knowledge distillation domain $\mathbb{D}_S^{K+1}$ and the target domain $\mathbb{D}_T$ is smaller than the following upper-bound for all source domain $k \in \{1, \cdots, K\}$, that is, $\epsilon_{\mathbb{D}_S^{K+1}}(h)$ should satisfy:*

$$\sup_{h \in \mathcal{H}} |\epsilon_{\mathbb{D}_S^{K+1}}(h) - \epsilon_{\mathbb{D}_T}(h)| \leq \inf_{h \in \mathcal{H}} |\epsilon_{\mathbb{D}_S^{K+1}}(h) - \epsilon_{\mathbb{D}_S^k}(h)|$$
$$+ \frac{1}{2}d_{\mathcal{H}\Delta\mathcal{H}}(\mathbb{D}_S^k, \mathbb{D}_T) + \lambda_S^k$$

**For good domains:**

KD3A provides better consensus knowledge with *Knowledge Vote*, making $\epsilon_{\mathbb{D}_S^{K+1}}$ closer to $\epsilon_{\mathbb{D}^T}$

**For bad domains:**

KD3A filters out their knowledge with *Consensus Focus*, making $\epsilon_{\mathbb{D}_S^{K+1}}$ stay away from bad domains.
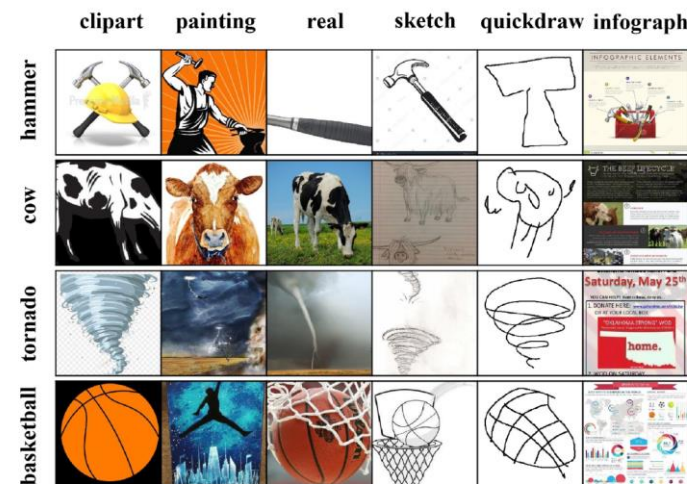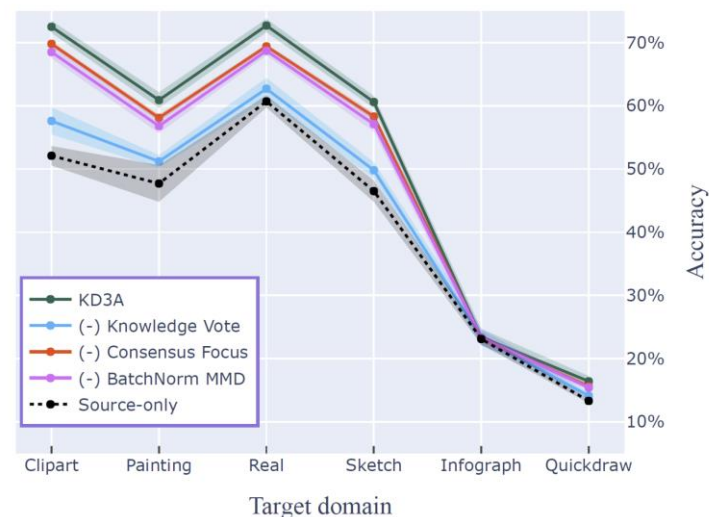
## Domain Adaptation Performance on DomainNet

| Standards | Methods | Clipart | Infograph | Painting | Quickdraw | Real | Sketch | Avg |
|---|---|---|---|---|---|---|---|---|
| W/o DA | Oracle | $69.3_{\pm0.37}$ | $34.5_{\pm0.42}$ | $66.3_{\pm0.67}$ | $66.8_{\pm0.51}$ | $80.1_{\pm0.59}$ | $60.7_{\pm0.48}$ | 63.0 |
| | Source-only | $52.1_{\pm0.51}$ | $23.1_{\pm0.28}$ | $47.7_{\pm0.96}$ | $13.3_{\pm0.72}$ | $60.7_{\pm0.32}$ | $46.5_{\pm0.56}$ | 40.6 |
| $\mathcal{H}-$div. | MDAN | $60.3_{\pm0.41}$ | $25.0_{\pm0.43}$ | $50.3_{\pm0.36}$ | $8.2_{\pm1.92}$ | $61.5_{\pm0.46}$ | $51.3_{\pm0.58}$ | 42.8 |
| | M$^3$SDA | $58.6_{\pm0.53}$ | $\mathbf{26.0_{\pm0.89}}$ | $52.3_{\pm0.55}$ | $6.3_{\pm0.58}$ | $62.7_{\pm0.51}$ | $49.5_{\pm0.76}$ | 42.6 |
| Knowledge Ensemble | DAEL | $70.8_{\pm0.14}$ | $26.5_{\pm0.13}$ | $57.4_{\pm0.28}$ | $12.2_{\pm0.7}$ | $65.0_{\pm0.23}$ | $60.6_{\pm0.25}$ | 48.7 |
| Source Selection | CMSS | $64.2_{\pm0.18}$ | $28.0_{\pm0.2}$ | $53.6_{\pm0.39}$ | $16.0_{\pm0.12}$ | $63.4_{\pm0.21}$ | $53.8_{\pm0.35}$ | 46.5 |
| Others | DSBN* | 60.3 | 22.6 | 52.3 | 9.1 | 62.7 | 47.6 | 42.4 |
| Decentralized UMDA | SHOT* | 61.7 | 22.2 | 52.6 | 12.2 | 67.7 | 48.6 | 44.2 |
| | FADA* | 59.1 | 21.7 | 47.9 | 8.8 | 60.8 | 50.4 | 41.5 |
| | FADA | $45.3_{\pm0.7}$ | $16.3_{\pm0.8}$ | $38.9_{\pm0.7}$ | $7.9_{\pm0.4}$ | $46.7_{\pm0.4}$ | $26.8_{\pm0.4}$ | 30.3 |
| | KD3A | $\mathbf{72.5_{\pm0.62}}$ | $23.4_{\pm0.43}$ | $\mathbf{60.9_{\pm0.71}}$ | $\mathbf{16.4_{\pm0.28}}$ | $\mathbf{72.7_{\pm0.55}}$ | $\mathbf{60.6_{\pm0.32}}$ | $\mathbf{51.1}$ |

Table 1. UMDA accuracy (%) on the DomainNet dataset. Our model KD3A achieves 51.1% accuracy, significantly outperforming all other baselines. Moreover, KD3A achieves the oracle performance on two domains: clipart and sketch. *: The best results recorded in our re-implementation.



[DomainNet Dataset.]

# Experiments

## Domain Adaptation Performance on DigitFive

| Methods | mt | mm | sv | syn | usps | Avg |
|---|---|---|---|---|---|---|
| Oracle | $99.5_{\pm0.08}$ | $95.4_{\pm0.15}$ | $92.3_{\pm0.14}$ | $98.7_{\pm0.04}$ | $99.2_{\pm0.09}$ | 97.0 |
| Source-only | $92.3_{\pm0.91}$ | $63.7_{\pm0.83}$ | $71.5_{\pm0.75}$ | $83.4_{\pm0.79}$ | $90.71_{\pm0.54}$ | 80.3 |
| MDAN | $97.2_{\pm0.98}$ | $75.7_{\pm0.83}$ | $82.2_{\pm0.82}$ | $85.2_{\pm0.58}$ | $93.3_{\pm0.48}$ | 86.7 |
| M$^3$SDA | $98.4_{\pm0.68}$ | $72.8_{\pm1.13}$ | $81.3_{\pm0.86}$ | $89.6_{\pm0.56}$ | $96.2_{\pm0.81}$ | 87.7 |
| CMSS | $99.0_{\pm0.08}$ | $75.3_{\pm0.57}$ | $88.4_{\pm0.54}$ | $93.7_{\pm0.21}$ | $97.7_{\pm0.13}$ | 90.8 |
| DSBN* | 97.2 | 71.6 | 77.9 | 88.7 | 96.1 | 86.3 |
| FADA | $91.4_{\pm0.7}$ | $62.5_{\pm0.7}$ | $50.5_{\pm0.3}$ | $71.8_{\pm0.5}$ | $91.7_{\pm1}$ | 73.6 |
| FADA* | 92.5 | 64.5 | 72.1 | 82.8 | 91.7 | 80.8 |
| SHOT | $98.2_{\pm0.37}$ | $80.2_{\pm0.41}$ | $84.5_{\pm0.32}$ | $\mathbf{91.1}_{\pm0.23}$ | $97.1_{\pm0.28}$ | 90.2 |
| KD3A$^\dagger$ | $99.1_{\pm0.15}$ | $86.9_{\pm0.11}$ | $82.2_{\pm0.26}$ | $89.2_{\pm0.19}$ | $98.4_{\pm0.11}$ | 91.2 |
| KD3A | $\mathbf{99.2}_{\pm0.12}$ | $\mathbf{87.3}_{\pm0.23}$ | $\mathbf{85.6}_{\pm0.17}$ | $89.4_{\pm0.28}$ | $\mathbf{98.5}_{\pm0.25}$ | **92.0** |



[DigitFive Dataset.]

Table 6. UMDA accuracy (%) on the **Digit-5**. *: The best results recorded in our re-implementation. †: Methods trained without data-augmentation. Our model KD3A achieves 92.0% accuracy and outperforms all other baselines.

# Experiments

## Domain Adaptation Performance on Office-Caltech10

| Methods | A | C | D | W | Avg |
|---|---|---|---|---|---|
| Oracle | 99.7 | 98.4 | 99.8 | 99.7 | 99.4 |
| Source-only | 86.1 | 87.8 | 98.3 | 99.0 | 92.8 |
| MDAN | 98.9 | 98.6 | 91.8 | 95.4 | 96.1 |
| $M^3SDA$ | 94.5 | 92.2 | **99.2** | 99.5 | 96.4 |
| CMSS | 96.0 | 93.7 | 99.3 | 99.6 | 97.2 |
| DSBN* | 93.2 | 91.6 | 98.9 | 99.3 | 95.8 |
| FADA | $84.2_{\pm 0.5}$ | $88.7_{\pm 0.5}$ | $87.1_{\pm 0.6}$ | $88.1_{\pm 0.4}$ | 87.1 |
| SHOT | 96.4 | 96.2 | 98.5 | **99.7** | 97.7 |
| KD3A[†] | $96.0_{\pm 0.07}$ | $95.2_{\pm 0.08}$ | $97.9_{\pm 0.11}$ | $99.6_{\pm 0.03}$ | 97.2 |
| KD3A | $\mathbf{97.4}_{\pm 0.08}$ | $\mathbf{96.4}_{\pm 0.11}$ | $98.4_{\pm 0.08}$ | $\mathbf{99.7}_{\pm 0.02}$ | **97.9** |

Table 9. UMDA accuracy (%) on the Office-Caltech10. *: The best results recorded in our re-implementation. †: Methods trained without data-augmentation.



*[Office-Caltech10 Dataset.]*

backpack monitor headphone bike mouse

# Experiments

## Domain Adaptation Performance on AmazonReview

| Methods | *Books* | *DVDs* | *Elec.* | *Kitchen* | Avg. |
|---|---|---|---|---|---|
| Source-only | 74.4 | 79.2 | 73.5 | 71.4 | 74.6 |
| MDAN | 78.6 | **80.7** | 85.4 | 86.3 | 82.8 |
| FADA | 78.1 | 82.7 | 77.4 | 77.5 | 78.9 |
| KD3A | **79.0** | 80.6 | **85.6** | **86.9** | **83.1** |

*Table 8.* The UMDA performance on Amazon Review dataset.

## Data Augmentation in Domain Adaptation

| | Clipart | Infograph | Painting | Avg |
|---|---|---|---|---|
| KD3A$^\dagger$ | $69.7_{\pm 0.67}$ | $21.2_{\pm 0.35}$ | $58.8_{\pm 0.66}$ | 48.8 |
| KD3A | $\mathbf{72.5}_{\pm \mathbf{0.62}}$ | $\mathbf{23.4}_{\pm \mathbf{0.43}}$ | $\mathbf{60.9}_{\pm \mathbf{0.71}}$ | **51.1** |
| | Quickdraw | Real | Sketch | |
| KD3A$^\dagger$ | $15.1_{\pm 0.21}$ | $70.4_{\pm 0.54}$ | $57.9_{\pm 0.41}$ | 48.8 |
| KD3A | $\mathbf{16.4}_{\pm \mathbf{0.28}}$ | $\mathbf{72.7}_{\pm \mathbf{0.55}}$ | $\mathbf{60.6}_{\pm \mathbf{0.32}}$ | **51.1** |

*Table 7.* The ablation study for data-augmentation strategies on DomainNet.†: Methods trained without data-augmentation.

| | |
|---|---|
| B | The most gorgeous artwork in comic books. It contains the most extraordinary and finest artwork of Alex Ross. |
| D | In my opinion it is the best American animated film ever released. It has a beautiful story with a ton of laughs, a lot of teachable moments. |
| E | My advice is if you need a CD rack that holds a lot of CD's? Save your money and invest in something nicer and more sturdy. |
| K | I absolutely love this product. my neighbor has four little yippers and my hepard/chow mix was antogonized on our side of the fence. |

*[AmazonReview Dataset.]*

**Augmentation Strategy:**

To reduce hyper-parameters, we use mixup as a unified augmentation strategy and simply set the mix-parameter $\alpha = 0.2$ in all experiments.

## Robustness To Negative Transfer

We construct irrelevant and malicious source domains on DomainNet and conduct synthesized experiments to show the robustness of KD3A to negative transfer.

Since *Quickdraw* is very different from others, and all models perform bad on it, we take *Quickdraw* as the irrelevant domain, denoted by IR-qdr. Moreover, we perform poisoning attack on the high-quality domain *Real* with m% wrong labels to construct malicious domain, denoted by MA-m.

quickdraw

| | $\mathcal{H}$-divergence | Info gain | Consensus focus | Domain drop |
|---|---|---|---|---|
| IR-qdr | 57.9 | 57.7 | 58.1 | **58.3** |
| MA-15 | 50.5 | 50.5 | **52.1** | |
| MA-30 | 49.8 | 48.9 | **51.1** | 50.7 |
| MA-50 | 47.6 | 46.3 | 50.6 | |

Table 2. Average UMDA accuracy (%) with irrelevant and malicious domains. IR-qdr means to use the *Quickdraw* as the irrelevant source domain, while MA-m means to construct a malicious source domain with $m\%$ mislabeled data. With consensus focus, our KD3A is robust to negative transfer.
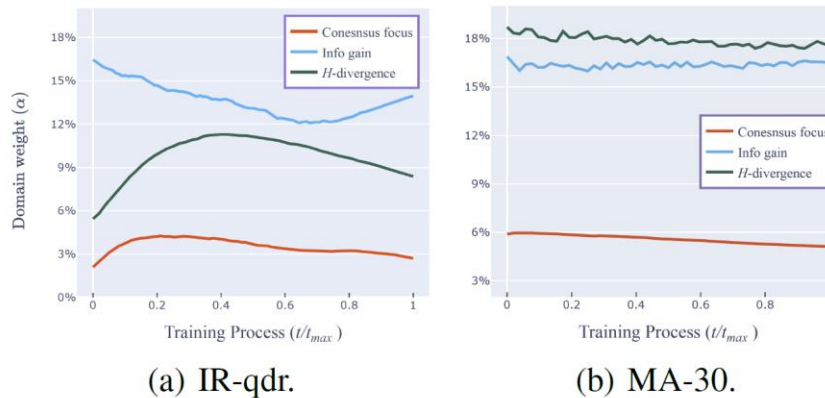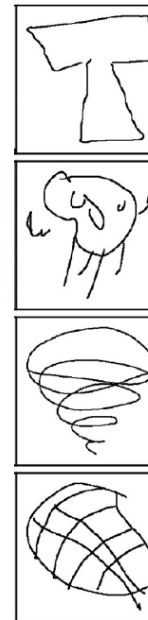
(a) IR-qdr.

(b) MA-30.

Figure 4. Weights assigned to the irrelevant and malicious domains in the training process. Our consensus focus can identify these bad domains with the low weights.

## Communication Efficiency And Privacy Security

| $r$ | 0.2 | 0.5 | 1 | 2 | 10 | 100 |
|------|------|------|------|------|------|------|
| FADA | 39.2 | 40.3 | 40.5 | 40.5 | 40.8 | 41.5 |
| KD3A | **50.5** | **50.9** | **51.1** | **51.3** | **51.3** | **52.0** |

*Table 3.* Average UMDA accuracy (%) with different communication rounds $r$ for our KD3A and FADA. KD3A achieves good performance with low communication cost (e.g., $r \leq 1$).
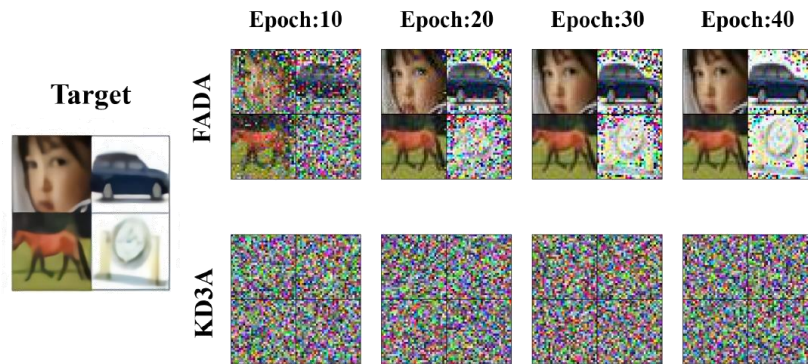


*Figure 5.* The gradient leakage attack (Zhu et al., 2019) on decentralized training strategy. KD3A is robust to this attack while FADA causes the privacy leakage.

**Conclusions:**

1. Due to the adversarial training strategy, FADA works under large communication rounds (i.e. $r = \textbf{100}$).

2. KD3A works under the low communication cost with $r = \textbf{1}$, leading to a $\textbf{100} \times$ communication reduction.

3. KD3A is robust to communication rounds. For example, the accuracy only drops **0.9%** when $r$ decreases from 100 to 1.

4. KD3A works under extremely low communication cost (e.g., $r=\textbf{0.2\&0.5}$).

5. Due to the low communication cost, our KD3A is robust to the advanced gradient leakage attack, which demonstrates high privacy security.

# Thanks

## Acknowledgement

Visual Analytics Group | State Key Lab of CAD&CG, Zhejiang University

浙江大学 ZHEJIANG UNIVERSITY