

A High-Dimensional Statistical Method for Optimizing Transfer Quantities in Multi-Source Transfer Learning

Qingyue Zhang*, Haohao Fu*, Guanbo Huang*, Yaoyuan Liang, Chang Chu,
Tianren Peng, Yanru Wu, Qi Li, Yang Li, Shao-Lun Huang

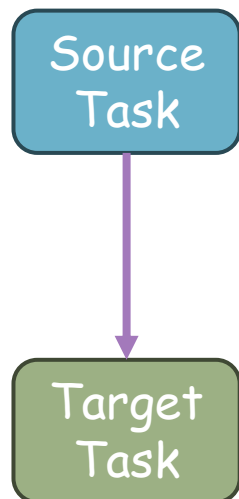
Tsinghua University

zhangqy23@mails.tsinghua.edu.cn

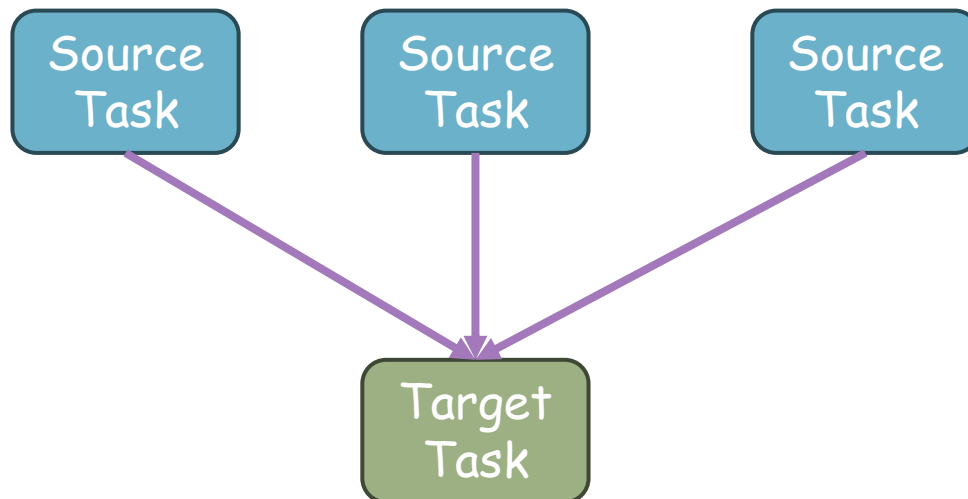
- **Background**
- Preliminaries
- Main Results
- Conclusion

Background: Transfer Learning

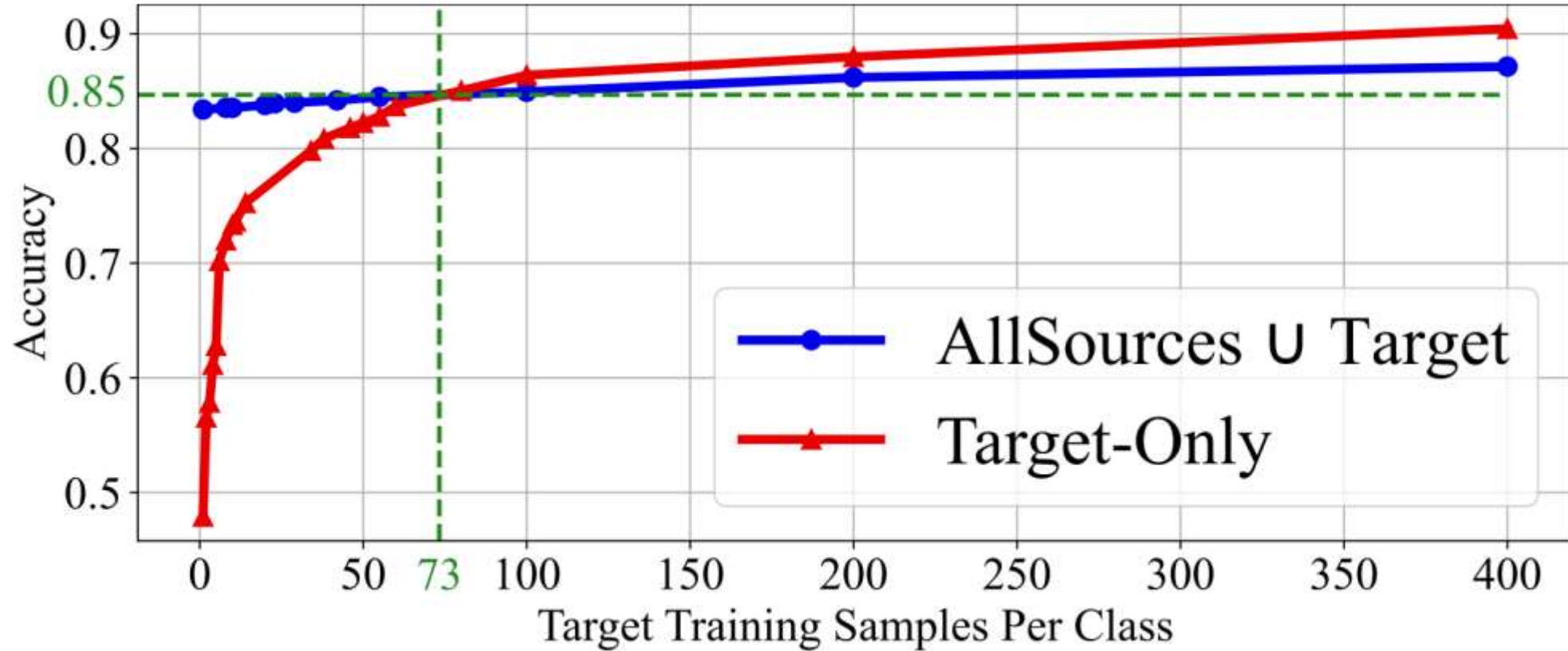
Single-source
Transfer Learning



Muti-source
Transfer Learning



Background: Transfer Learning



What is the **optimal transfer quantity** of each source task?

- Background
- **Preliminaries**
- Main Results
- Conclusion

- The maximum likelihood estimator (MLE) is defined as the maximizer of the empirical log-likelihood:

$$\hat{\underline{\theta}}_{\text{MLE}} = \arg \max_{\underline{\theta}} \frac{1}{n} \sum_{x \in \mathcal{D}} \log P_{X;\underline{\theta}}(x)$$

- **Asymptotic Normality:** As the sample size increases, the distribution of the normalized estimation error converges in law to a multivariate Gaussian distribution

$$\sqrt{n} \left(\hat{\underline{\theta}}_{\text{MLE}} - \underline{\theta}^* \right) \xrightarrow{d} \mathcal{N}(0, J(\underline{\theta}^*)^{-1})$$

- **The Fisher information matrix**, which characterizes the amount of information carried by the distribution about the parameter, is defined as

$$J(\underline{\theta})^{d \times d} = \mathbb{E} \left[\left(\frac{\partial}{\partial \underline{\theta}} \log P_{X;\underline{\theta}} \right) \left(\frac{\partial}{\partial \underline{\theta}} \log P_{X;\underline{\theta}} \right)^T \right].$$

- **Target task:** Training samples $X^{N_0} = \{x_j\}_{j=1}^{N_0} \stackrel{\text{i.i.d.}}{\sim} P_{X;\underline{\theta}_0}$.
- **Source tasks:** Each source task \mathcal{S}_i has $X_i^{N_i} = \{x_{i,j}\}_{j=1}^{N_i} \stackrel{\text{i.i.d.}}{\sim} P_{X;\underline{\theta}_i}$, $i \in [1, K]$.
- **Estimator:** The training process is formulated as a parameter estimation problem. The MLE $\hat{\underline{\theta}}$ is obtained using all target samples and a selected subset of source samples:

$$\hat{\underline{\theta}} = \arg \max_{\underline{\theta}} \left[\sum_{x \in X^{N_0}} \log P_{X;\underline{\theta}}(x) + \sum_{i=1}^K \sum_{x \in X^{n_i}} \log P_{X;\underline{\theta}}(x) \right].$$

- **Objective:** Find optimal transfer quantities n_1^*, \dots, n_K^* by minimizing the expected K-L divergence between the true target distribution $P_{X;\underline{\theta}_0}$ and the learned one $P_{X;\hat{\underline{\theta}}}$.

$$n_1^*, \dots, n_K^* = \arg \min_{n_1, \dots, n_K} \mathbb{E}[D(P_{X;\underline{\theta}_0} || P_{X;\hat{\underline{\theta}}})].$$

- Background
- Preliminaries
- **Main Results**
- Conclusion

Theorem (Theorem 1. Single-source transfer with 1-dimensional models)

In transfer between $P_{X;\theta_0}$ and $P_{X;\theta_1}$, where $\theta_0, \theta_1 \in \mathbb{R}$ and $|\theta_0 - \theta_1| = O(\frac{1}{\sqrt{N_0}})$. Then, the K-L measure $\mathbb{E}[D(P_{X;\theta_0} || P_{X;\hat{\theta}})]$ can be expressed as:

$$\frac{1}{2} \left(\underbrace{\frac{1}{N_0 + n_1}}_{\text{variance term}} + \underbrace{\frac{n_1^2}{(N_0 + n_1)^2} t}_{\text{bias term}} \right) + o\left(\frac{1}{N_0}\right), \text{ where } t \triangleq J(\theta_0)(\theta_1 - \theta_0)^2.$$

Optimal Transfer Quantity: The optimal transfer quantity n_1^* is

$$n_1^* = \begin{cases} N_1, & \text{if } N_0 \cdot t \leq 0.5 \\ \min\left(N_1, \frac{N_0}{2N_0 t - 1}\right), & \text{if } N_0 \cdot t > 0.5 \end{cases}. \quad (1)$$

Theorem (Theorem 2. Single-source transfer with high-dimensional models)

In transfer between $P_{X;\theta_0}$ and $P_{X;\theta_1}$, where $\theta_0, \theta_1 \in \mathbb{R}^d$ and $\|\theta_0 - \theta_1\| = O(\frac{1}{\sqrt{N_0}})$. Then, the K-L measure $\mathbb{E}[D(P_{X;\theta_0}||P_{X;\hat{\theta}})]$ can be expressed as:

$$\frac{d}{2} \left(\frac{1}{N_0 + n_1} + \frac{n_1^2}{(N_0 + n_1)^2} t \right) + o \left(\frac{1}{N_0} \right), \text{ where } t \triangleq \frac{(\theta_1 - \theta_0)^\top J(\theta_0)(\theta_1 - \theta_0)}{d}.$$

Optimal Transfer Quantity: The optimal transfer quantity n_1^* is

$$n_1^* = \begin{cases} N_1, & \text{if } N_0 \cdot t \leq 0.5 \\ \min \left(N_1, \frac{N_0}{2N_0 t - 1} \right), & \text{if } N_0 \cdot t > 0.5 \end{cases}. \quad (2)$$

Theorem (Theorem 3. Multi-source transfer with high-dimensional models)

In the multi-source setting, we denote $s = \sum_{i=1}^K n_i$, and $\alpha_i = \frac{n_i}{s}$. Then, the K-L measure can be expressed as:

$$\frac{d}{2} \left(\frac{N_0}{(N_0 + s)^2} + \frac{s^2}{(N_0 + s)^2} t \right) + o\left(\frac{1}{N_0}\right), t = \frac{\underline{\alpha}^T \Theta^T J(\underline{\theta}_0) \Theta \underline{\alpha}}{d},$$

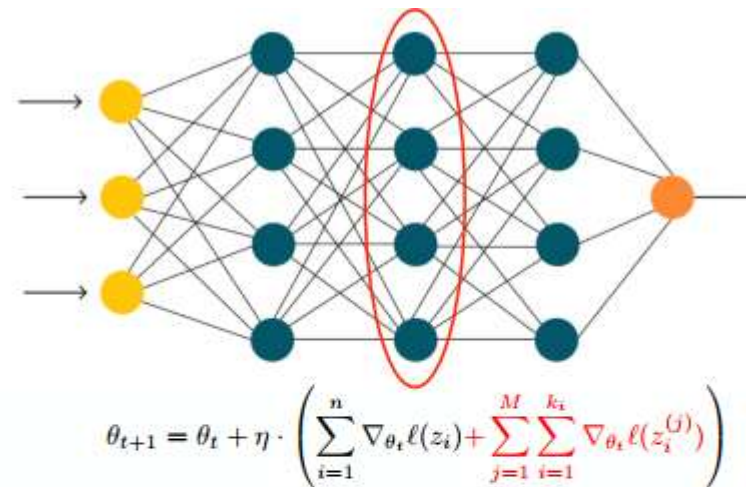
where $\underline{\alpha} = [\alpha_1, \dots, \alpha_K]^T$, and $\Theta^{d \times K} = [\underline{\theta}_1 - \underline{\theta}_0, \dots, \underline{\theta}_K - \underline{\theta}_0]$.

Optimal Transfer Quantities: *The optimal quantities n_i^* are obtained by minimizing the K-L measure. Specifically, we perform a grid search over the feasible range of s , and for each candidate s' , solve a quadratic programming problem to obtain the optimal $\underline{\alpha}'$ under the constraint set $\mathcal{A}(s')$. The final optimal pair $(s^*, \underline{\alpha}^*)$ is selected as the one yielding the minimum value of the objective function among all candidate pairs. Finally, we use $n_i^* = s^* \alpha_i^*$.*

Algorithm 1 OTQMS: Training

- 1: **Input:** Target data $D_{\mathcal{T}} = \{(z_{\mathcal{T}}^i, y_{\mathcal{T}}^i)\}_{i=1}^{N_0}$, source data $\{D_{S_k} = \{(z_{S_k}^i, y_{S_k}^i)\}_{i=1}^{N_k}\}_{k=1}^K$, model type f_{θ} and its parameters θ_0 for target task and $\{\theta_k\}_{k=1}^K$ for source tasks, parameter dimension d .
// z represents the feature and y represents the label
- 2: **Parameter:** Learning rate η .
- 3: **Initialize:** randomly initialize θ_0 , use parameters of pretrained source models to initialize $\{\theta_k\}_{k=1}^K$.
- 4: **Output:** a well-trained θ_0 for target task model f_{θ_0} .
- 5: $D_{train} \leftarrow D_{\mathcal{T}}$ // Initialize the training dataset by target task samples
- 6: **repeat** // Use dynamic strategy to train the target task
- 7: $\mathcal{L}_{train} \leftarrow \frac{1}{|D_{train}|} \sum_{(y^i, z^i) \in D_{train}} \ell(y^i, f_{\theta_0}(z^i))$
- 8: $\theta_0 \leftarrow \theta_0 - \eta \nabla_{\theta_0} \mathcal{L}_{train}$
- 9: $\Theta \leftarrow [\theta_1 - \theta_0, \dots, \theta_K - \theta_0]^T$
- 10: $J(\theta_0) \leftarrow (\nabla_{\theta_0} \mathcal{L}_{train})(\nabla_{\theta_0} \mathcal{L}_{train})^T$
- 11: $(s^*, \alpha^*) \leftarrow \arg \min_{(s, \alpha)} \frac{d}{2} \left(\frac{1}{N_0 + s} + \frac{s^2}{(N_0 + s)^2} \frac{\alpha^T \Theta^T J(\theta_0) \Theta \alpha}{d} \right)$
- 12: $D_{source} \leftarrow \bigcup_{k=1}^K \left\{ D_{S_k}^* \mid D_{S_k}^* \stackrel{\text{rand}}{\subseteq} D_{S_k}, |D_{S_k}^*| = s^* \alpha_k^* \right\}$
- 13: $D_{train} \leftarrow D_{source} \cup D_{\mathcal{T}}$ // Update the training dataset
- 14: **until** θ_0 converges;

Dynamic Strategy



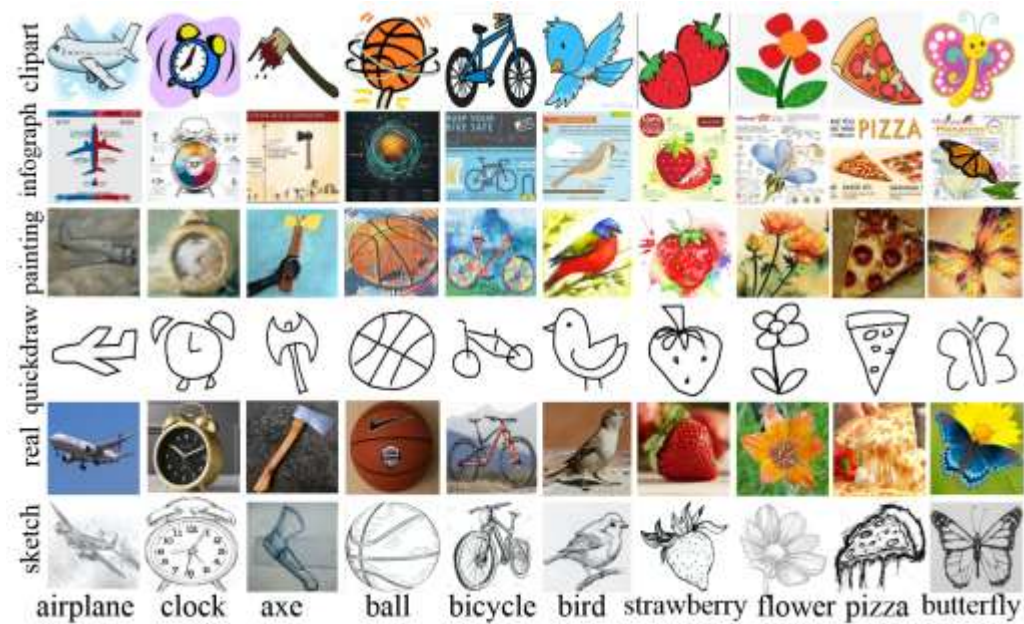
- The quantities in computing the optimal number of samples can be estimated from current parameters;
- Update the parameters in SGD algorithm with both target samples and k_i samples from source task i ;
- Iteratively update the optimal sample numbers and network parameters.

Main Results: Dataset of Experiment



Office-Home

Four domains, each formulated as a 65-class classification task.



Domainnet

Six domains, each formulated as a 345-class classification task.

Table 2: **Multi-Source Transfer Performance on DomainNet and Office-Home.** The arrows indicate transferring from the rest tasks. The highest/second-highest accuracy is marked in Bold/Underscore form respectively.

Method	Backbone	DomainNet							Office-Home				
		→C	→I	→P	→Q	→R	→S	Avg	→Ar	→Cl	→Pr	→Rw	Avg
Unsupervised-all-shots													
MSFDA [19]	ResNet50	66.5	21.6	56.7	20.4	70.5	54.4	48.4	75.6	62.8	84.8	85.3	77.1
DATE [9]	ResNet50	-	-	-	-	-	-	-	75.2	60.9	85.2	84.0	76.3
M3SDA [17]	ResNet101	57.2	24.2	51.6	5.2	61.6	49.6	41.5	-	-	-	-	-
Supervised-10-shots													
Few-Shot Methods:													
H-ensemble [29]	ViT-S	53.4	21.3	54.4	19.0	70.4	44.0	43.8	71.8	47.5	77.6	79.1	69.0
MADA [30]	ViT-S	51.0	12.8	60.3	15.0	81.4	22.7	40.5	78.4	58.3	82.3	85.2	76.1
MADA [30]	ResNet50	66.1	23.9	60.4	31.9	75.4	52.5	51.7	72.2	64.4	82.9	81.9	75.4
MCW [12]	ViT-S	54.9	21.0	53.6	20.4	70.8	42.4	43.9	68.9	48.0	77.4	86.0	70.1
WADN [20]	ViT-S	68.0	29.7	59.1	16.8	74.2	55.1	50.5	60.3	39.7	66.2	68.7	58.7
Source-Ablation Methods:													
Target-Only	ViT-S	14.2	3.3	23.2	7.2	41.4	10.6	16.7	40.0	33.3	54.9	52.6	45.2
Single-Source-Avg	ViT-S	50.4	22.1	44.9	24.7	58.8	42.5	40.6	65.2	53.3	74.4	72.7	66.4
Single-Source-Best	ViT-S	60.2	28.0	55.4	28.4	66.0	49.7	48.0	72.9	60.9	80.7	74.8	72.3
AllSources \cup Target	ViT-S	71.7	32.4	60.0	31.4	71.7	58.5	54.3	77.0	62.3	84.9	84.5	77.2
OTQMS (Ours)	ViT-S	72.8	33.8	61.2	33.8	73.2	59.8	55.8	78.1	64.5	85.2	84.9	78.2

Main Results: Experiment

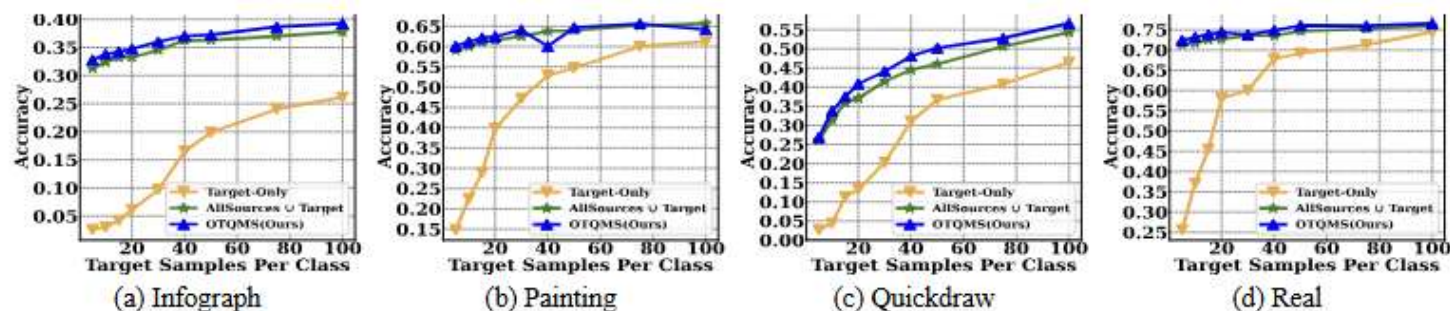


Figure 3: Performance comparison with increasing target shots up to 100 per class on DomainNet dataset (I, P, Q and R domains). OTQMS (blue) outperforms other methods.

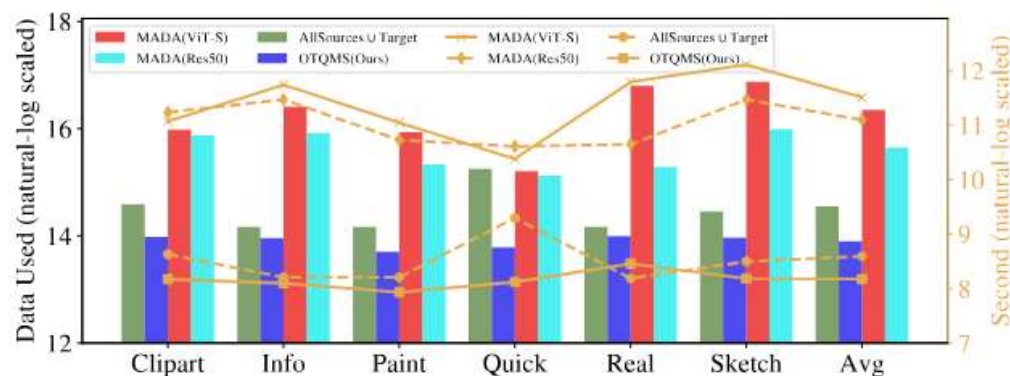


Figure 4: Data efficiency comparison of average sample usage and training time on DomainNet dataset, the left vertical axis represents the amount of sample usage, with green bars indicating AllSources \cup Target data counts, blue bars about OTQMS, red bars about MADA(ViT-S) and cyan bars about MADA(Res50), while the right orange vertical axis and lines represent training time.

Table 6: Multi-Source Transfer with LoRA on Office-Home. We apply LoRA on ViT-B backbone for PEFT.

Method	Backbone	Office-Home				
		→Ar	→Cl	→Pr	→Rw	Avg
<i>Supervised-10-shots Source-Ablation:</i>						
Target-Only	ViT-B	59.8	42.2	69.5	72.0	60.9
Single-Source-avg	ViT-B	72.2	59.9	82.6	81.0	73.9
Single-Source-best	ViT-B	74.4	61.8	84.9	81.9	75.8
AllSources \cup Target	ViT-B	<u>81.1</u>	<u>66.0</u>	<u>88.0</u>	<u>89.2</u>	<u>81.1</u>
OTQMS (Ours)	ViT-B	81.5	68.0	89.2	90.3	82.3

Table 4: Multi-task performance on four tasks of Office-Home.

Method	Backbone	Office-Home				
		Ar	Cl	Pr	Rw	Avg
Single-task	ViT-S	66.7	62.3	87.8	68.6	71.4
OTQMS	ViT-S	81.7	76.0	88.6	87.5	83.5

- Background
- Preliminaries
- Main Results
- **Conclusion**

- **Theoretical Framework:** We formulate multi-source transfer learning as a parameter estimation problem and derive solutions for the **optimal transfer quantity** of each source by minimizing a K-L divergence–based generalization error in the asymptotic regime.
- **Algorithm Design:** We develop **OTQMS**, a dynamic and data-efficient algorithm that iteratively updates transfer quantities using empirical Fisher information, enabling adaptive resampling and improved target model training.
- **Experimental Results:** Experiments on **DomainNet** and **Office-Home** demonstrate that OTQMS achieves higher accuracy and better data efficiency than state-of-the-art methods, and remains robust under different shot settings and architectures.

Thanks for your attention !

<https://github.com/zqy0126/OTQMS>