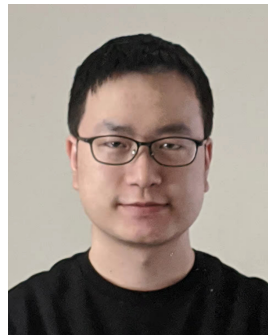


A Unified Meta-Learning Framework for Dynamic Transfer Learning



Jun Wu
UIUC



Jingrui He
UIUC

Roadmap

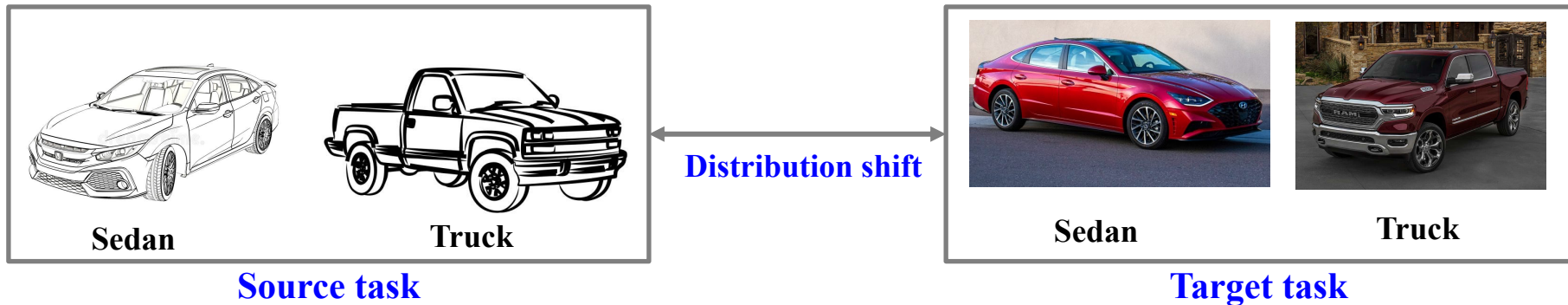


- Background
- Problem definition
- Theoretical analysis
- Proposed framework
- Experiments
- Conclusion

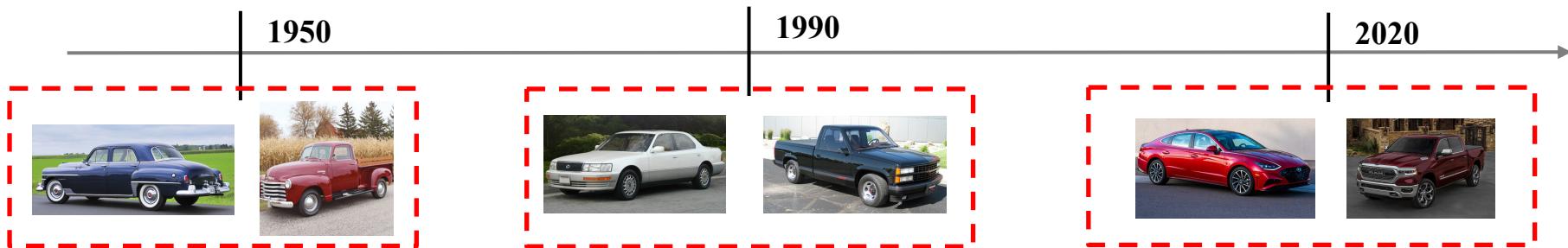
Background



- Distribution shift of transfer learning
 - E.g., sedans vs trucks



- Time-evolving distribution
 - New data are collected at different time stamps



Time-evolving data distribution

- Hoffman, Judy, Trevor Darrell, and Kate Saenko. "Continuous manifold based adaptation for evolving visual domains." CVPR. 2014

Problem Definition

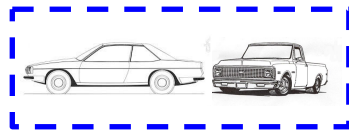


- Dynamic transfer learning

- Given: Labeled dynamic source task $\{\mathcal{D}_j^s\}_{j=1}^N$ (with data $D_j^s = \{x_{ij}^s, y_{ij}^s\}$);

- unlabeled dynamic target task $\{\mathcal{D}_j^t\}_{j=1}^N$ (with data $D_j^t = \{x_{ij}^t\}$)

Task 1



...



1st time stamp

2nd time stamp

Nth time stamp

Task 2



...

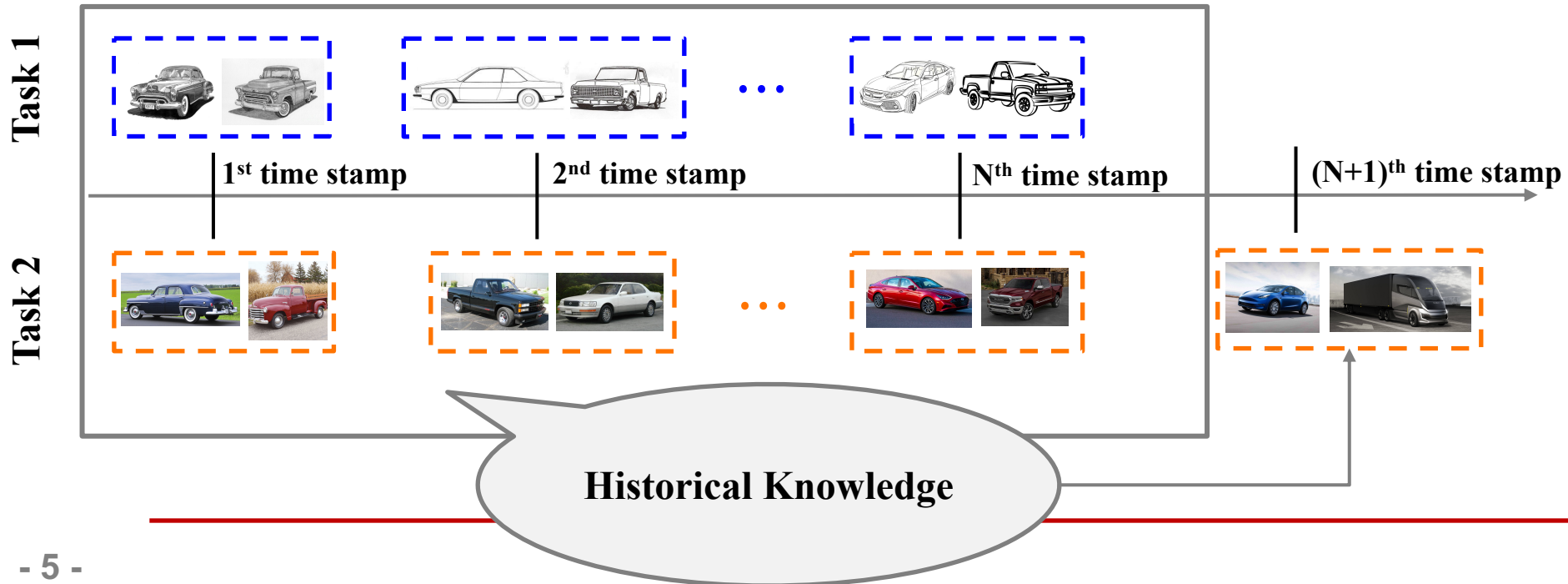


Problem Definition



■ Dynamic transfer learning

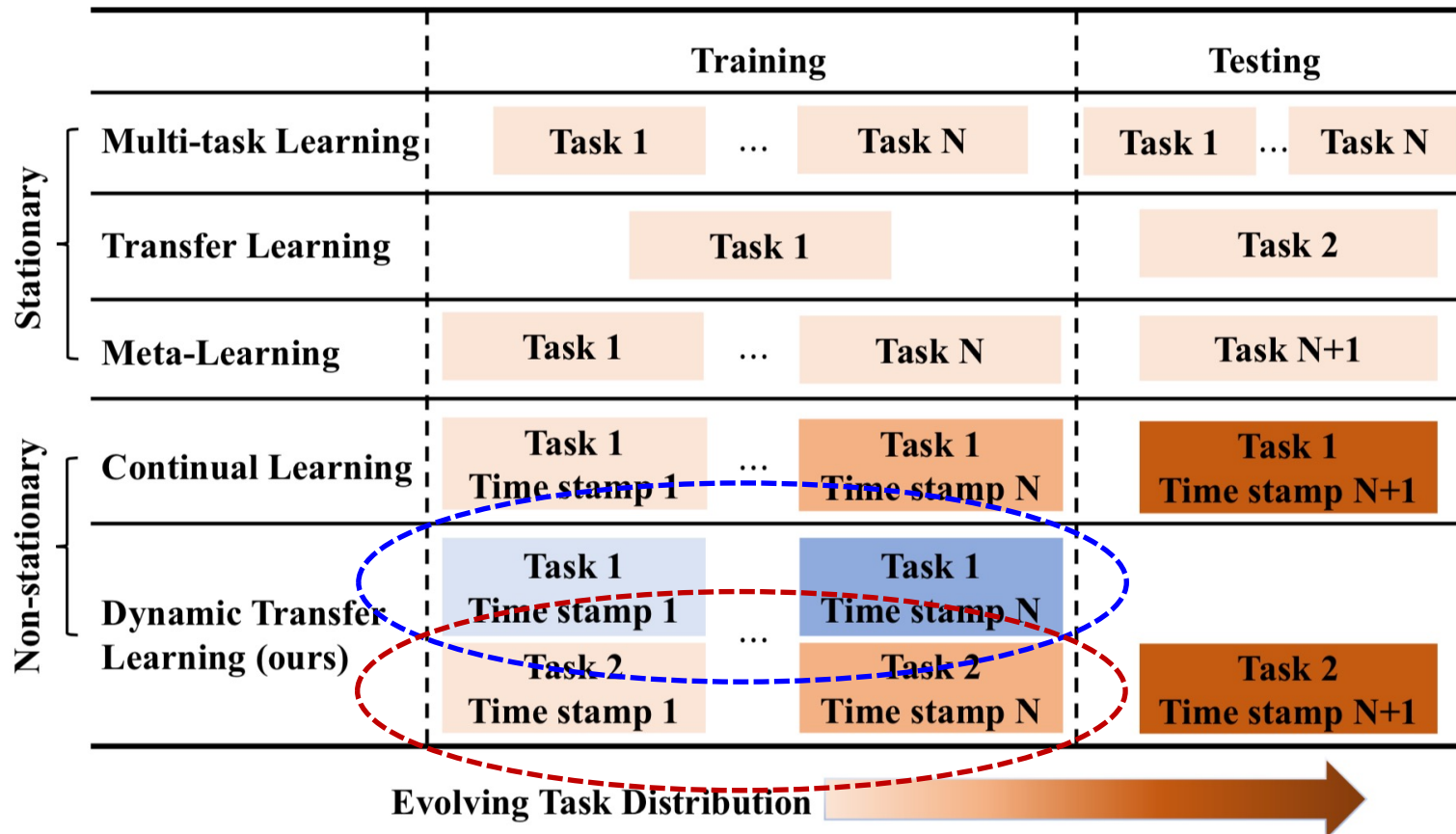
- Given: Labeled dynamic source task $\{\mathcal{D}_j^s\}_{j=1}^N$ (with data $D_j^s = \{x_{ij}^s, y_{ij}^s\}$);
unlabeled dynamic target task $\{\mathcal{D}_j^t\}_{j=1}^N$ (with data $D_j^t = \{x_{ij}^t\}$)
- Goal: Learn the prediction function on the newest target task \mathcal{D}_{N+1}^t



Problem Comparison



- Dynamic transfer learning



- Sener, Ozan, and Vladlen Koltun. "Multi-task learning as multi-objective optimization." NeurIPS, 2018.
- Finn, Chelsea, Pieter Abbeel, and Sergey Levine. "Model-agnostic meta-learning for fast adaptation of deep networks." ICML, 2017.
- Rolnick, David, Arun Ahuja, Jonathan Schwarz, Timothy Lillicrap, and Gregory Wayne. "Experience replay for continual learning." NeurIPS, 2019.

Assumptions



- A1: The class labels of the source task are available at any time stamp
- A2: The source and target tasks are related at the initial time stamp
 - $d(\mathcal{D}_1^s, \mathcal{D}_1^t) \leq \Delta$ at time stamp $j = 1$
- A3: The data distributions of both source and target tasks are continuously changing over time
 - $d(\mathcal{D}_j^s, \mathcal{D}_{j+1}^s) \leq \Delta$ and $d(\mathcal{D}_j^t, \mathcal{D}_{j+1}^t) \leq \Delta$ for time stamp $j \geq 1$

Theoretical Analysis



- Error bound on the newest target task \mathcal{D}_{N+1}^t
 - ① Empirical errors of historical source and target tasks;
 - ② Maximal distribution discrepancy across tasks and across time stamps;
 - ③ Maximal labeling difference across tasks and across time stamps;
 - ④ Average Rademacher complexity

Theorem 1: Assume that the loss function L is μ -admissible and obeys the triangle inequality, with probability at least $1 - \delta$, the expected error ϵ_{N+1}^t for the newest target task \mathcal{D}_{N+1}^t is bounded by

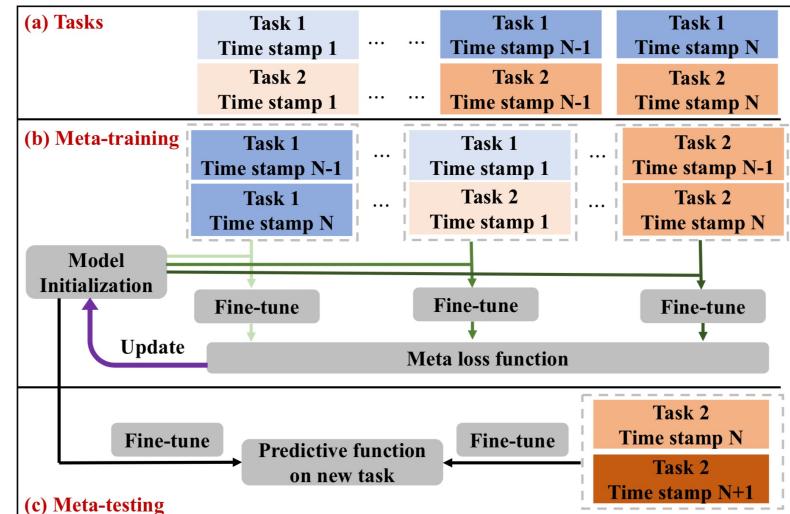
$$\epsilon_{N+1}^t(h) \leq \underbrace{\frac{1}{2N} \sum_{j=1}^N (\hat{\epsilon}_j^s(h) + \hat{\epsilon}_j^t(h))}_{\text{①}} + \underbrace{\frac{N+2}{2} d_{max}}_{\text{②}} + \underbrace{\frac{N+2}{2} \lambda_{max}}_{\text{③}} + \underbrace{\mathfrak{R}(H_L)}_{\text{④}} + \frac{\mu}{N} \sqrt{\frac{\log 1/\delta}{m}}$$

Proposed Framework: L2E



- A unified meta-learning framework
 - Learning to evolve (L2E)

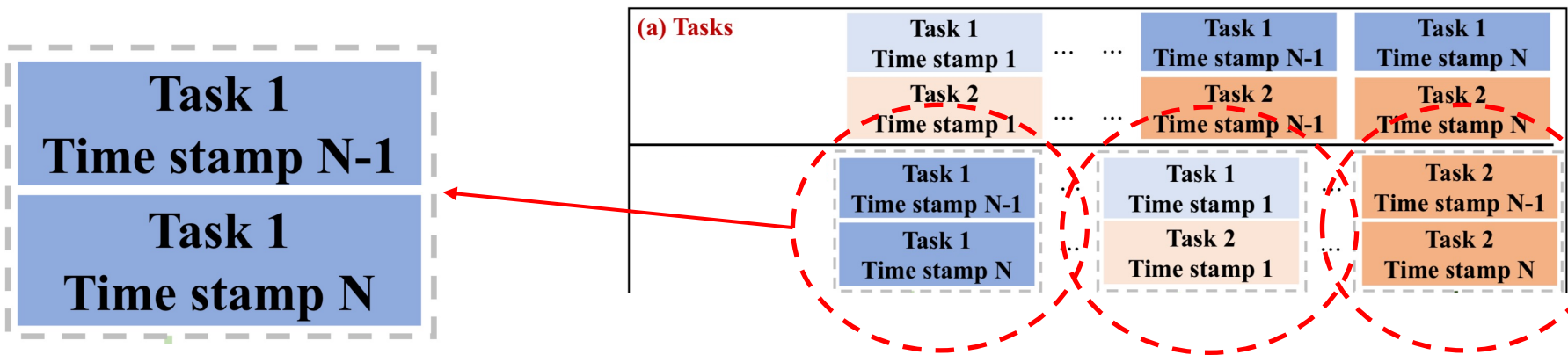
- Three key components
 - Meta-pairs of tasks
 - Meta-training
 - Meta-testing



Step 1: Meta-pair of Tasks



- Construction of meta-pair of tasks
 - Consecutive source/target task
 - Initial source and target tasks



② Maximal distribution discrepancy across tasks and across time stamps;

$$d_{max} = \max \left\{ \max_{1 \leq j \leq N-1} d(\mathcal{D}_j^s, \mathcal{D}_{j+1}^s), d(\mathcal{D}_1^s, \mathcal{D}_1^t), \max_{1 \leq j \leq N} d(\mathcal{D}_j^t, \mathcal{D}_{j+1}^t) \right\}$$

Step 2: Meta-training



- Objective function
 - Learn an optimal model initialization

$$\theta_N^* = \arg \min_{\theta} \sum_{k=1}^{N-1} \zeta_k(M_k(\theta); D_k^{val}) \quad \leftarrow \text{Meta-pair of tasks}$$

$$M_k(\theta) \leftarrow \theta - \alpha \cdot \nabla_{\theta} \zeta_k(\theta; D_k^{train}) \quad \leftarrow \text{One-step GD}$$

$$\zeta_k(\theta; D_k) = \left\{ \begin{array}{ll} \hat{\epsilon}_{-k+1}^s(\theta) + \gamma \cdot d(\mathcal{D}_{-k}^s, \mathcal{D}_{-k+1}^s); & k < 0 \\ \hat{\epsilon}_1^s(\theta) + \gamma \cdot d(\mathcal{D}_1^s, \mathcal{D}_1^t); & k = 0 \\ \hat{\epsilon}_k^t(\theta) + \gamma \cdot d(\mathcal{D}_k^t, \mathcal{D}_{k+1}^t); & k > 0 \end{array} \right\} \text{ Meta-pairs of tasks}$$

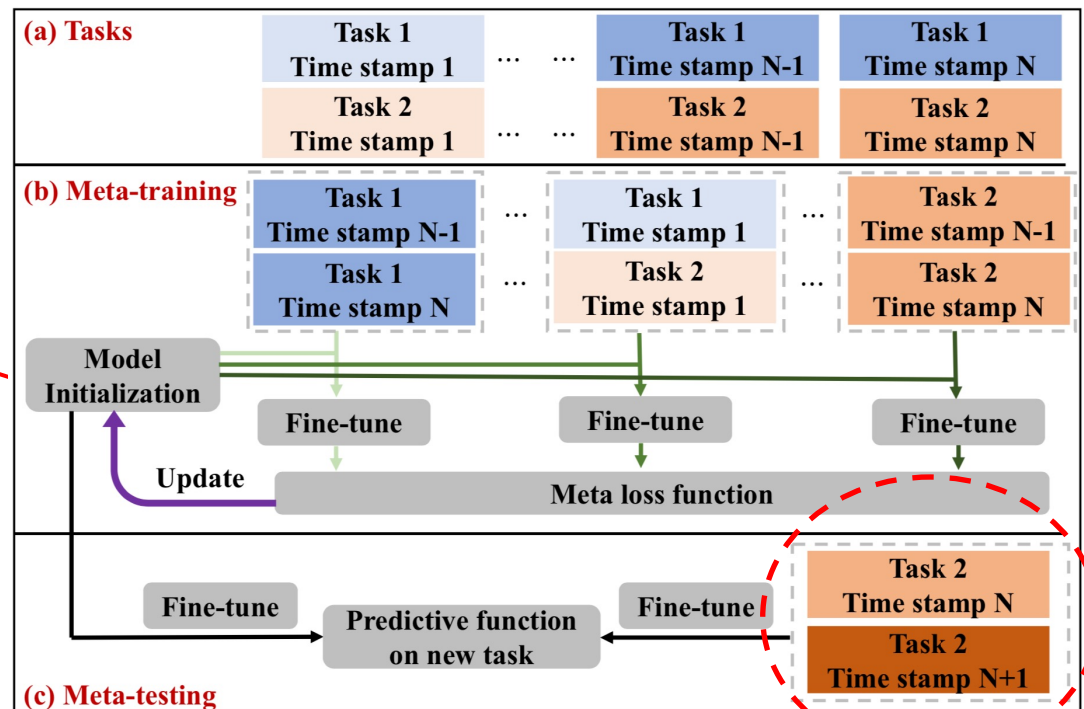
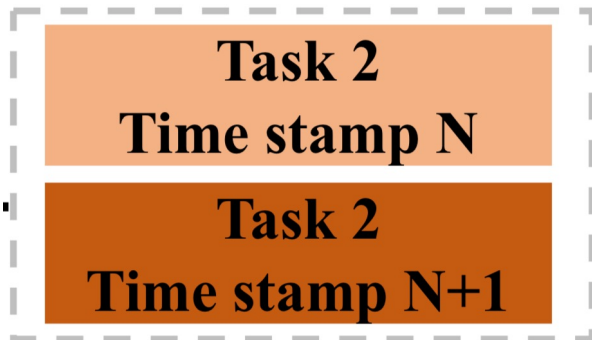
It is equivalent to standard transfer learning for each meta-pair of tasks

Step 3: Meta-testing



- Fine-tune θ_N^* for the newest target task

$$\theta_{N+1} = M_N(\theta_N^*) \leftarrow \theta_N^* - \alpha \cdot \nabla_{\theta} \zeta_N(\theta; D_N^{train})$$



• Finn, Chelsea, Pieter Abbeel, and Sergey Levine. "Model-agnostic meta-learning for fast adaptation of deep networks." ICML, 2017.

■ Data sets

- Office-31
 - Image-CLEF
 - Caltran
- } Generate the dynamic task by adding the random noise and rotation to the original images

■ Baselines

- Static adaptation: SourceOnly, DANN, MDD
- Multi-source adaptation: MDAN, M3SDA, DARN
- Dynamic adaptation: CUA, TransLATE, GST

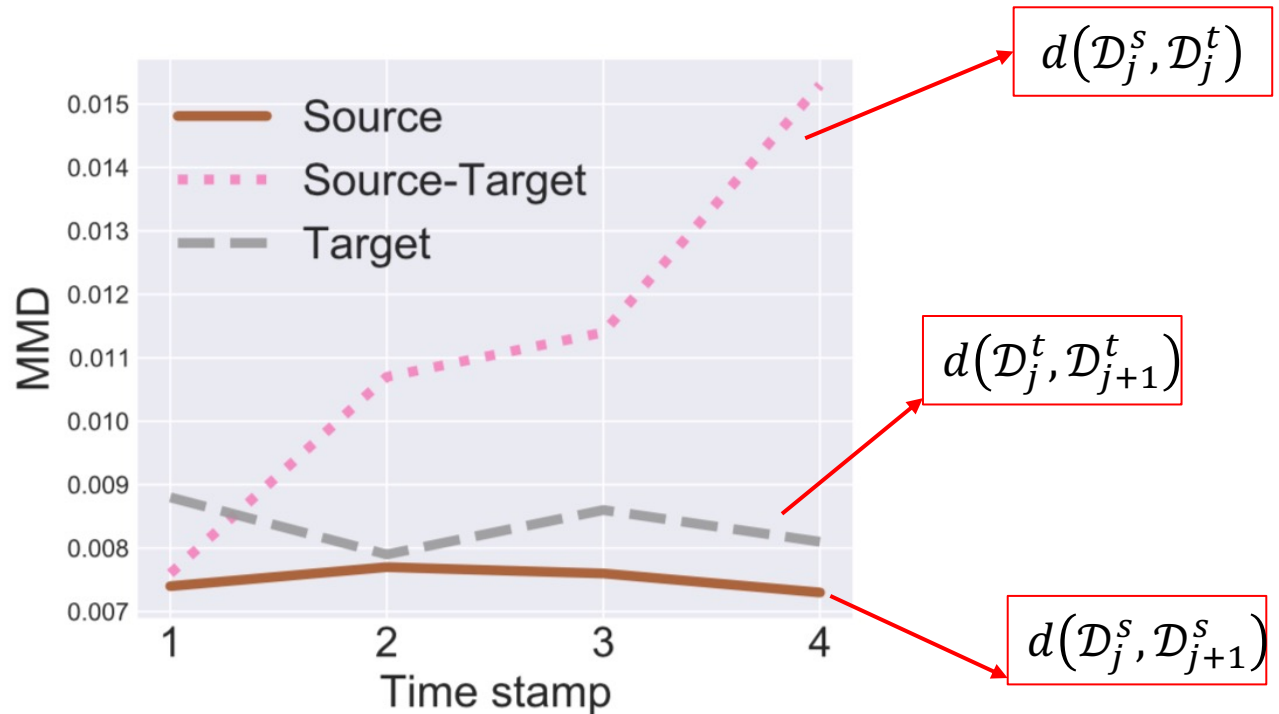
■ Metric

- Acc: Classification accuracy on the newest target task
- H-Acc: Average classification accuracy on all the historical target tasks

Evolution of Tasks



- Visualizing the distribution discrepancy on Image-CLEF (B→P)
 - MMD: Maximum mean discrepancy



- (1) The source and target tasks are changing smoothly
- (2) The relatedness between source and target tasks is decreasing over time

Results



Method	I \rightarrow C		B \rightarrow P		
	Acc	H-Acc	Acc	H-Acc	
Static adaptation	SourceOnly	0.26 \pm 0.01	0.51 \pm 0.01	0.24 \pm 0.00	0.43 \pm 0.01
	DANN	0.36 \pm 0.00	0.58 \pm 0.01	0.27 \pm 0.01	0.43 \pm 0.00
	MDD	0.41 \pm 0.01	0.62 \pm 0.01	0.28 \pm 0.03	0.42 \pm 0.01
Multi-source adaptation	MDAN	0.62 \pm 0.03	0.77 \pm 0.00	0.37 \pm 0.05	0.51 \pm 0.02
	M3SDA	0.56 \pm 0.03	0.74 \pm 0.01	0.39 \pm 0.02	0.52 \pm 0.02
	DARN	0.55 \pm 0.02	0.76 \pm 0.02	0.39 \pm 0.02	0.52 \pm 0.01
Dynamic adaptation	CUA	0.58 \pm 0.01	0.74 \pm 0.01	0.36 \pm 0.03	0.51 \pm 0.00
	TransLATE	0.64 \pm 0.01	0.76 \pm 0.00	0.40 \pm 0.03	0.55 \pm 0.01
	GST	0.39 \pm 0.01	0.54 \pm 0.03	0.32 \pm 0.01	0.31 \pm 0.02
L2E (ours)	0.66\pm0.02	0.80\pm0.01	0.44\pm0.04	0.57\pm0.02	

Method	I \rightarrow C		B \rightarrow P		
	Acc	H-Acc	Acc	H-Acc	
Static adaptation	SourceOnly	0.26 \pm 0.01	0.51 \pm 0.01	0.24 \pm 0.00	0.43 \pm 0.01
	DANN	0.36 \pm 0.00	0.58 \pm 0.01	0.27 \pm 0.01	0.43 \pm 0.00
	MDD	0.41 \pm 0.01	0.62 \pm 0.01	0.28 \pm 0.03	0.42 \pm 0.01
Multi-source adaptation	MDAN	0.62 \pm 0.03	0.77 \pm 0.00	0.37 \pm 0.05	0.51 \pm 0.02
	M3SDA	0.56 \pm 0.03	0.74 \pm 0.01	0.39 \pm 0.02	0.52 \pm 0.02
	DARN	0.55 \pm 0.02	0.76 \pm 0.02	0.39 \pm 0.02	0.52 \pm 0.01
Dynamic adaptation	CUA	0.58 \pm 0.01	0.74 \pm 0.01	0.36 \pm 0.03	0.51 \pm 0.00
	TransLATE	0.64 \pm 0.01	0.76 \pm 0.00	0.40 \pm 0.03	0.55 \pm 0.01
	GST	0.39 \pm 0.01	0.54 \pm 0.03	0.32 \pm 0.01	0.31 \pm 0.02
L2E (ours)	0.66\pm0.02	0.80\pm0.01	0.44\pm0.04	0.57\pm0.02	

(1) Higher performance on the newest target task

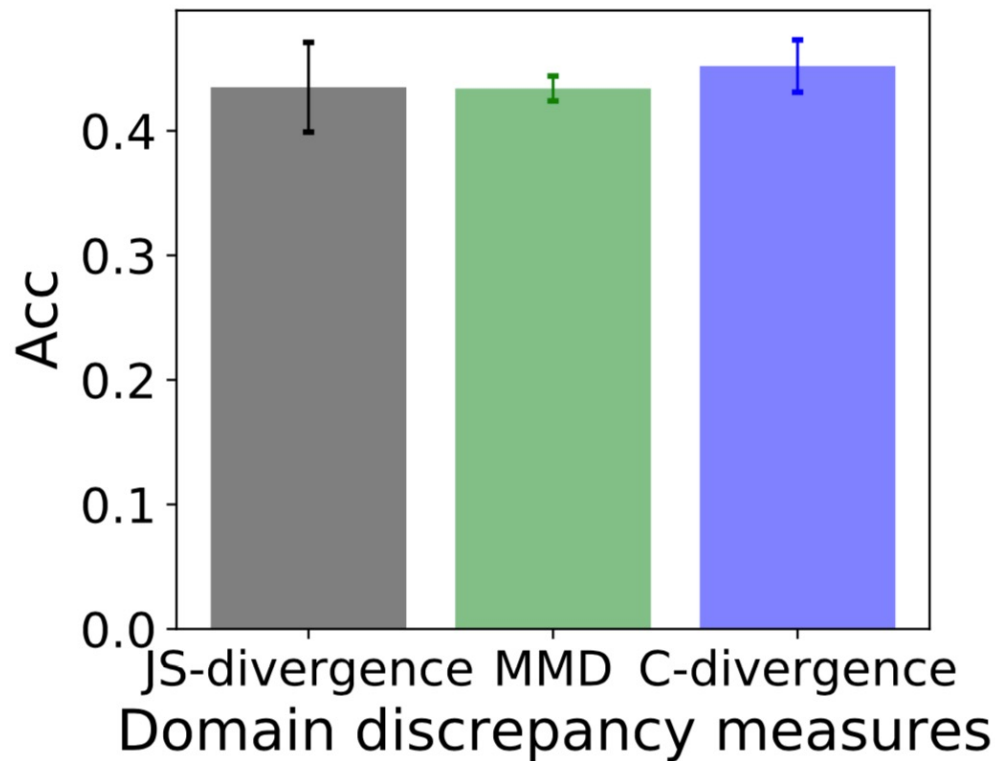
Results



Method	I \rightarrow C		B \rightarrow P		
	Acc	H-Acc	Acc	H-Acc	
Static adaptation	SourceOnly	0.26 \pm 0.01	0.51 \pm 0.01	0.24 \pm 0.00	0.43 \pm 0.01
	DANN	0.36 \pm 0.00	0.58 \pm 0.01	0.27 \pm 0.01	0.43 \pm 0.00
	MDD	0.41 \pm 0.01	0.62 \pm 0.01	0.28 \pm 0.03	0.42 \pm 0.01
Multi-source adaptation	MDAN	0.62 \pm 0.03	0.77 \pm 0.00	0.37 \pm 0.05	0.51 \pm 0.02
	M3SDA	0.56 \pm 0.03	0.74 \pm 0.01	0.39 \pm 0.02	0.52 \pm 0.02
	DARN	0.55 \pm 0.02	0.76 \pm 0.02	0.39 \pm 0.02	0.52 \pm 0.01
Dynamic adaptation	CUA	0.58 \pm 0.01	0.74 \pm 0.01	0.36 \pm 0.03	0.51 \pm 0.00
	TransLATE	0.64 \pm 0.01	0.76 \pm 0.00	0.40 \pm 0.03	0.55 \pm 0.01
	GST	0.39 \pm 0.01	0.54 \pm 0.03	0.32 \pm 0.01	0.31 \pm 0.02
L2E (ours)	0.66\pm0.02	0.80\pm0.01	0.44\pm0.04	0.57\pm0.02	

- (1) Higher performance on the newest target task
- (2) Higher performance on the historical target task

- Flexibility: L2E with different distribution discrepancy measures



- Theoretical results

- Derive the **generalization error bounds** of dynamic transfer learning

Theorem 1: Assume that the loss function L is μ -admissible and obeys the triangle inequality, with probability at least $1 - \delta$, the expected error ϵ_{N+1}^t for the newest target task \mathcal{D}_{N+1}^t is bounded by

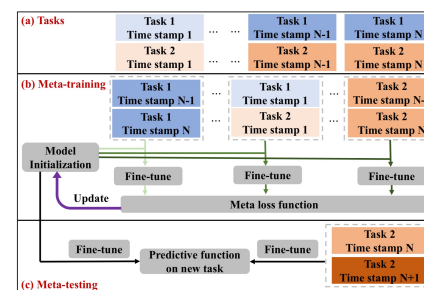
$$\epsilon_{N+1}^t(h) \leq \frac{1}{2N} \sum_{j=1}^N (\hat{\epsilon}_j^s(h) + \hat{\epsilon}_j^t(h)) + \frac{N+2}{2} (d_{max} + \lambda_{max}) + \mathfrak{R}(H_L) + \frac{\mu}{N} \sqrt{\frac{\log 1/\delta}{m}}$$

- Generic framework

- Propose a **meta-learning framework** (L2E) by reformulating the meta-pairs of tasks

- Empirical evaluation

- Demonstrate the **effectiveness** of our L2E framework on dynamic tasks



Method	I → C		B → P	
	Acc	H-Acc	Acc	H-Acc
SourceOnly	0.26±0.01	0.51±0.01	0.24±0.00	0.43±0.01
DANN	0.36±0.00	0.58±0.01	0.27±0.01	0.43±0.00
MDD	0.41±0.01	0.62±0.01	0.28±0.03	0.42±0.01
MDAN	0.62±0.03	0.77±0.00	0.37±0.05	0.51±0.02
M3SDA	0.56±0.03	0.74±0.01	0.39±0.02	0.52±0.02
DARN	0.55±0.02	0.76±0.02	0.39±0.02	0.52±0.01
CUA	0.58±0.01	0.74±0.01	0.36±0.03	0.51±0.00
TransLATE	0.64±0.01	0.76±0.00	0.40±0.03	0.55±0.01
GST	0.39±0.01	0.54±0.03	0.32±0.01	0.31±0.02
L2E (ours)	0.66±0.02	0.80±0.01	0.44±0.04	0.57±0.02



Thank You!

Please email me via junwu3@illinois.edu if you have any question.