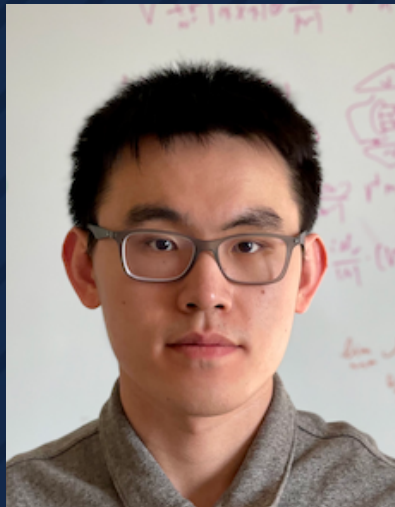




UNIVERSITY OF
ILLINOIS
URBANA-CHAMPAIGN

Understanding Gradual Domain Adaptation: *Improved Analysis, Optimal Path and Beyond* ICML 2022



Haoxiang Wang
PhD Candidate
ECE, UIUC

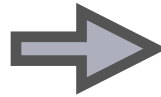


Bo Li
Assistant Professor
Computer Science, UIUC



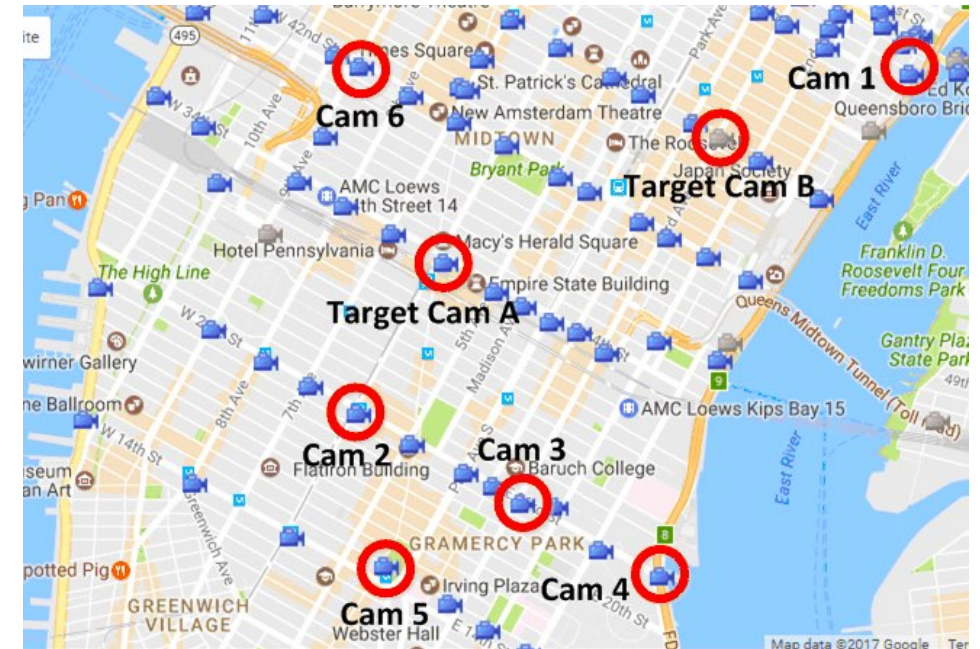
Han Zhao
Assistant Professor
Computer Science, UIUC

Domain Adaptation: Source Domain \longrightarrow Target Domain

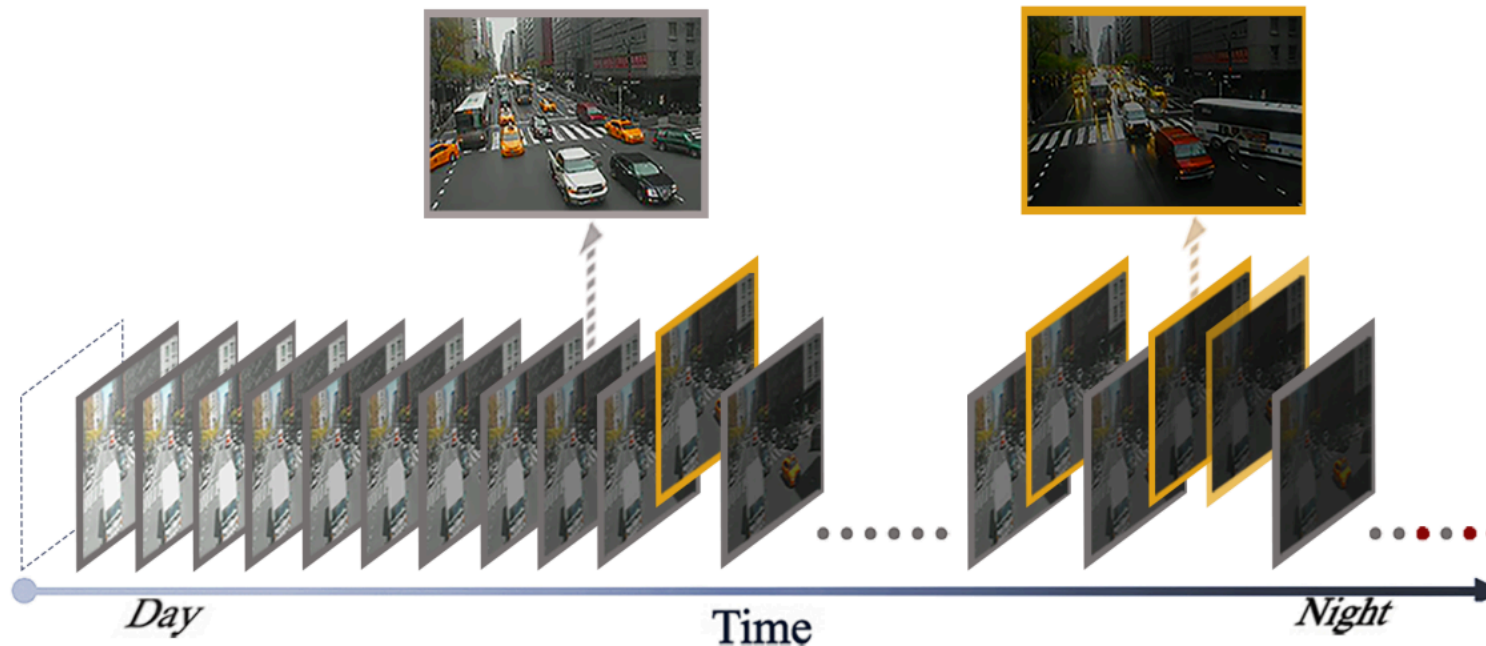


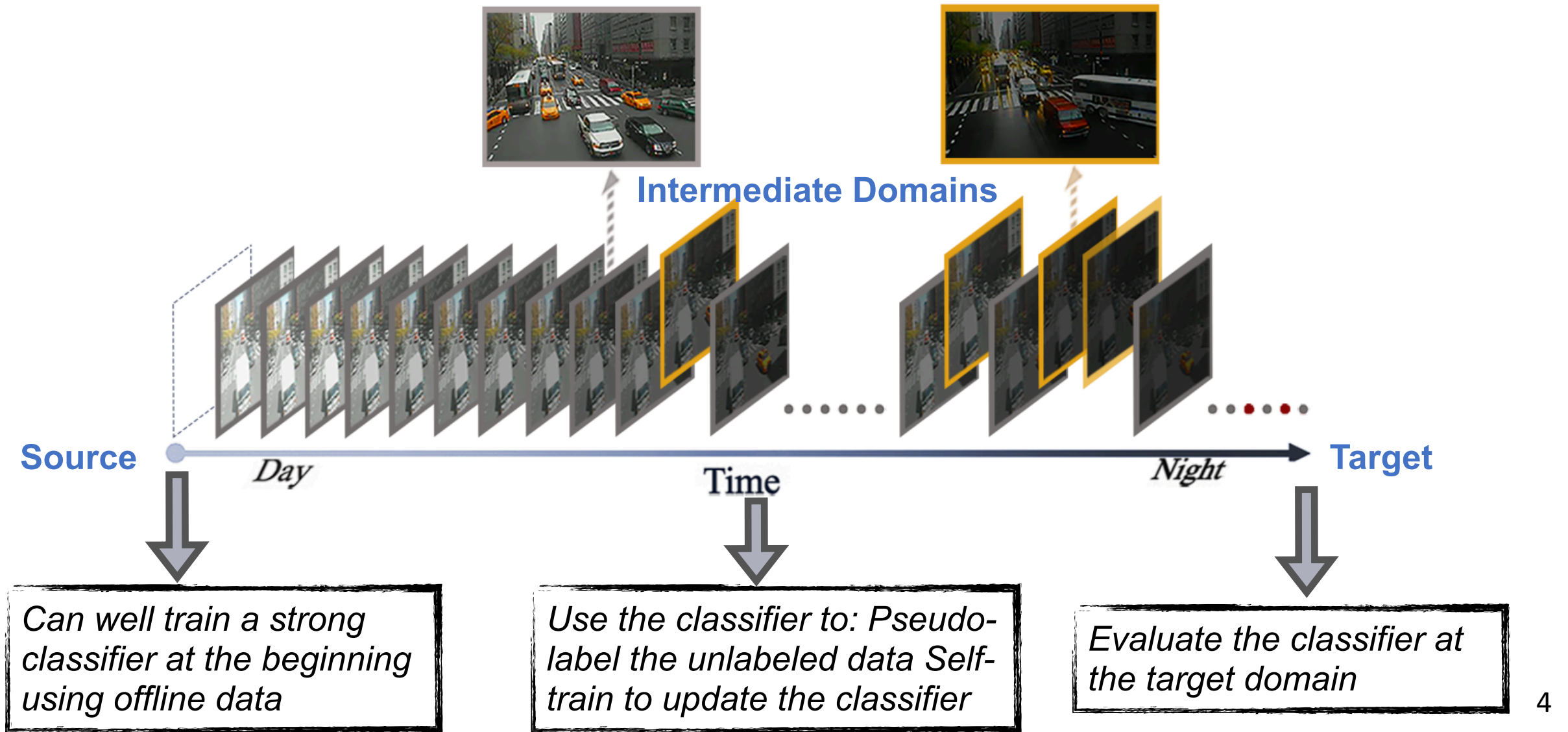
Source (with Labels)

Target (No Labels)



- In many real-world applications (e.g., self-driving cars, recommendation systems), distributions do not change abruptly
- Instead, they change **smoothly** in time / space
- Could we do better in this case?





With high probability,

$$\varepsilon_T(h_T) \leq e^{\mathcal{O}(T)} \left(\varepsilon_0(h_0) + \mathcal{O}\left(\frac{1}{\sqrt{n}} + \sqrt{\frac{\log T}{n}}\right) \right)$$

- h_0 : initial classifier; h_T : classifier evaluated at the target domain
- ε_0 : error in the source domain; ε_T : error in the target domain
- T : # of intermediate domains (time steps) - 1
- n : # of unlabeled data points in each intermediate/target domain

With high probability,

$$\varepsilon_T(h_T) \leq \boxed{e^{\mathcal{O}(T)}} \left(\varepsilon_0(h_0) + \mathcal{O}\left(\frac{1}{\sqrt{n}} + \sqrt{\frac{\log T}{n}}\right) \right)$$

The dependency on T is an **exponential** (pessimistic)

Is it possible to have a more optimistic generalization bound of the gradual self-training algorithm (hopefully with some insights for algorithm design)?

Theorem 1 (Generalization Bound for Gradual Self-Training). *For any $\delta \in (0, 1)$, the population loss of gradually self-trained classifier h_T in the target domain is upper bounded with probability at least $1 - \delta$ as*

$$\varepsilon_T(h_T) \leq \varepsilon_0(h_0) + \mathcal{O}\left(T\Delta + \frac{T}{\sqrt{n}} + T\sqrt{\frac{\log 1/\delta}{n}} + \frac{1}{\sqrt{nT}} + \sqrt{\frac{(\log nT)^{3L-2}}{nT}} + \sqrt{\frac{\log 1/\delta}{nT}}\right)$$

- h_0 : initial classifier; h_T : classifier evaluated at the target domain
- T : # of intermediate domains (time steps) - 1
- n : # of unlabeled data points in each domain
- Δ : average distributional distance between consecutive domains

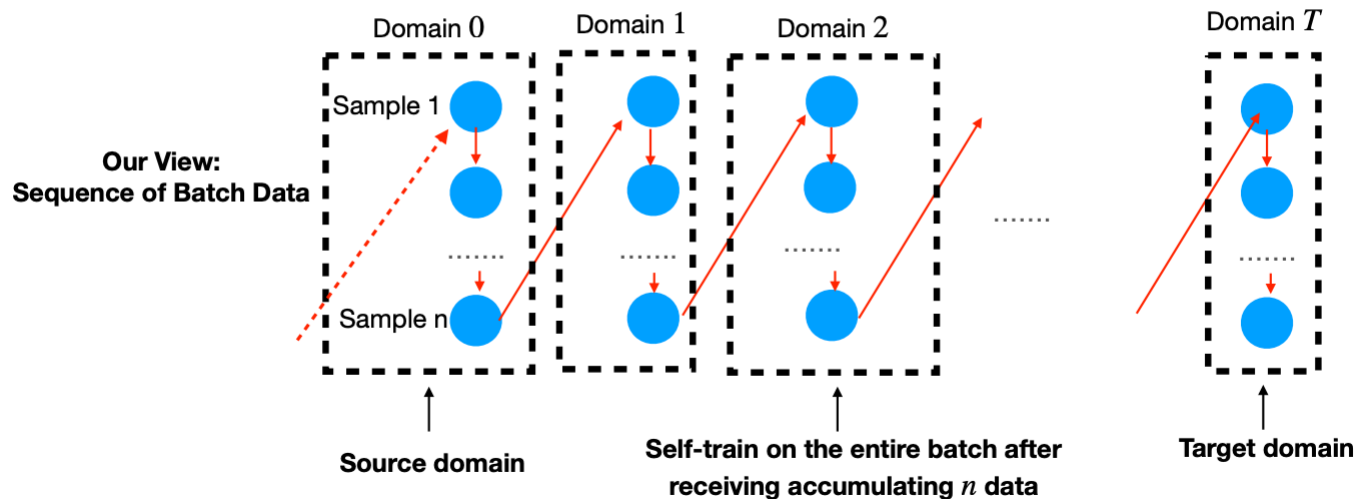
1. If the distribution distance is small, then for any smooth model, the error difference is also small.

$$|\varepsilon_\mu(h) - \varepsilon_\nu(h)| \leq \overbrace{\rho}^{\text{Lipschitz constants}} \sqrt{\overbrace{R^2}^{\text{p-Wasserstein Distance}} + 1} \overbrace{W_p(\mu, \nu)}^{\text{p-Wasserstein Distance}}$$

2. The self-training algorithm is stable:


$$\hat{h} = \arg \min_{f \in \mathcal{H}} \sum_{x \in S} \ell(f(x), h(x)) \quad \text{s.t.} \quad |\varepsilon_\mu(\hat{h}) - \varepsilon_\nu(h)| \leq \mathcal{O} \left(W_p(\mu, \nu) + \frac{\rho B + \sqrt{\log \frac{1}{\delta}}}{\sqrt{n}} \right)$$

3. A reduction to online learning (treating the pseudo-label as if they were ground-truth)




The dependency on T is an **exponential** (pessimistic)

[Kumar, Ma, Liang. ICML'20]: $\varepsilon_T(h_T) \leq \boxed{e^{\mathcal{O}(T)}} \left(\varepsilon_0(h_0) + \mathcal{O}\left(\frac{1}{\sqrt{n}} + \sqrt{\frac{\log T}{n}}\right) \right)$



Ours: $\varepsilon_T(h_T) \leq \varepsilon_0(h_0) + \boxed{\mathcal{O}\left(T\Delta + \frac{T}{\sqrt{n}}\right)} + \tilde{\mathcal{O}}\left(\frac{1}{\sqrt{nT}}\right)$

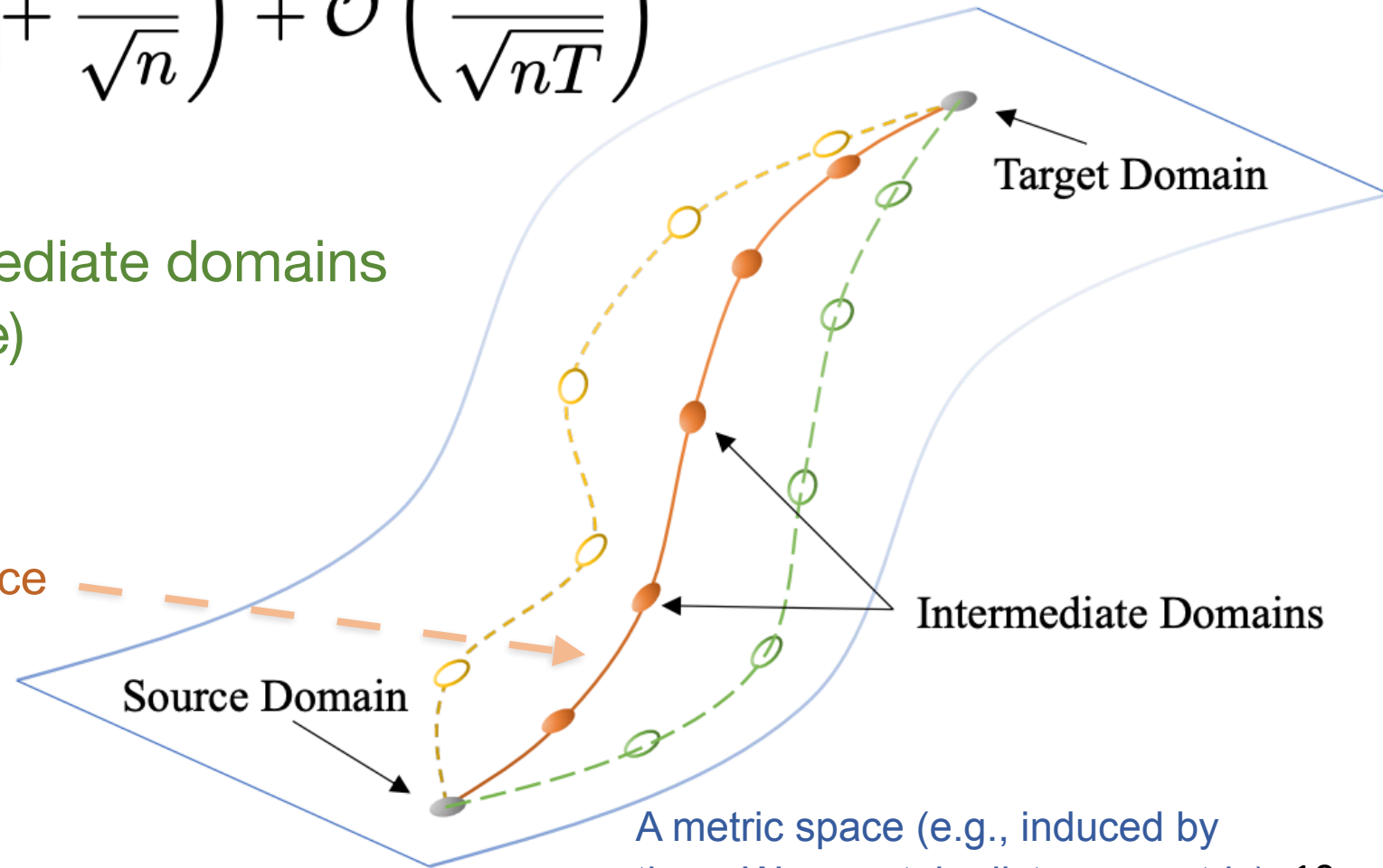


The dependency on T becomes **linear** (optimistic)

$$\varepsilon_T(h_T) \leq \varepsilon_0(h_0) + \mathcal{O} \left(\boxed{T\Delta} + \frac{T}{\sqrt{n}} \right) + \tilde{\mathcal{O}} \left(\frac{1}{\sqrt{nT}} \right)$$

$T\Delta$: length of the path of intermediate domains
(measured in some metric space)

Geodesic in the metric space



A metric space (e.g., induced by the p -Wasserstein distance metric)

$$\varepsilon_T(h_T) \leq \varepsilon_0(h_0) + \mathcal{O}\left(T\Delta + \frac{T}{\sqrt{n}}\right) + \tilde{\mathcal{O}}\left(\frac{1}{\sqrt{nT}}\right)$$

Minimize our error bound w.r.t. T

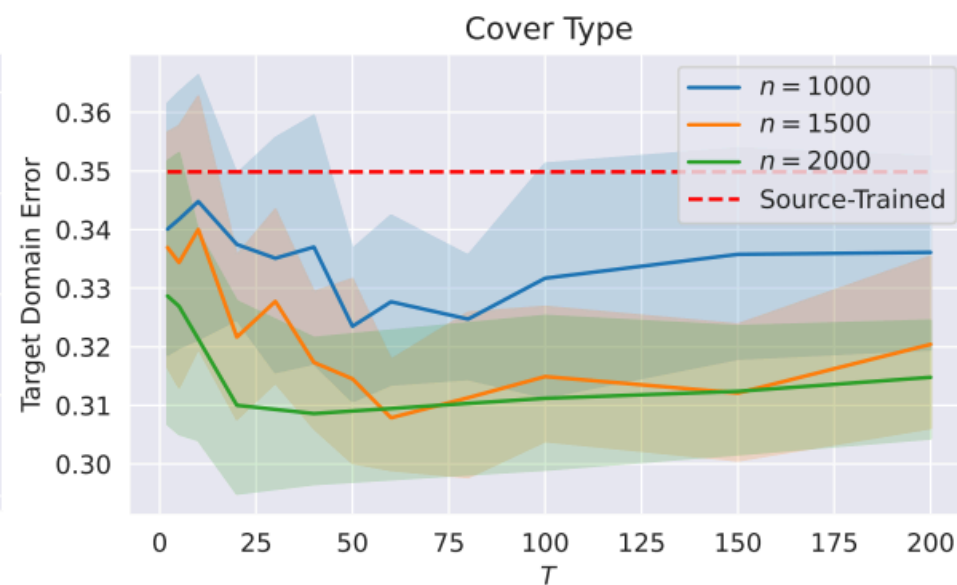
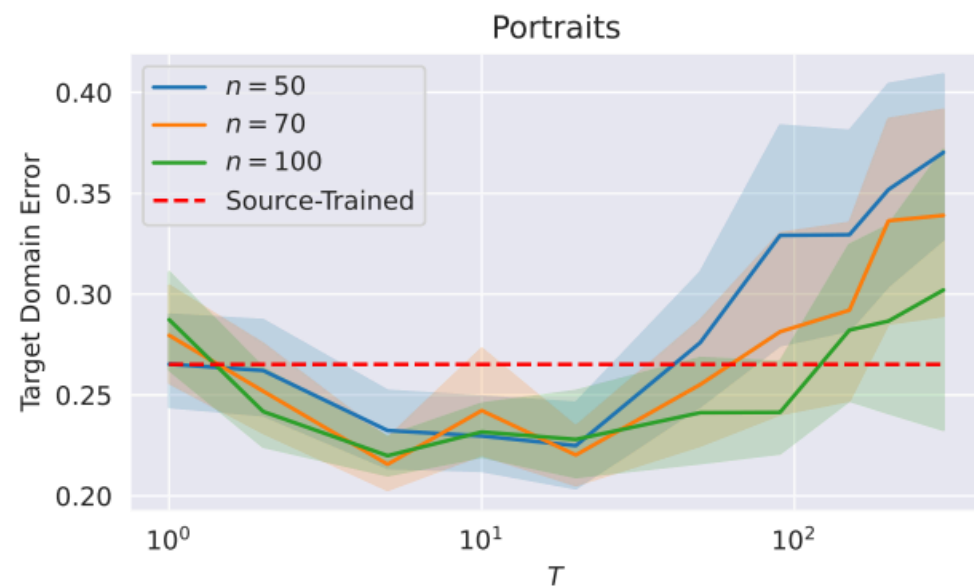
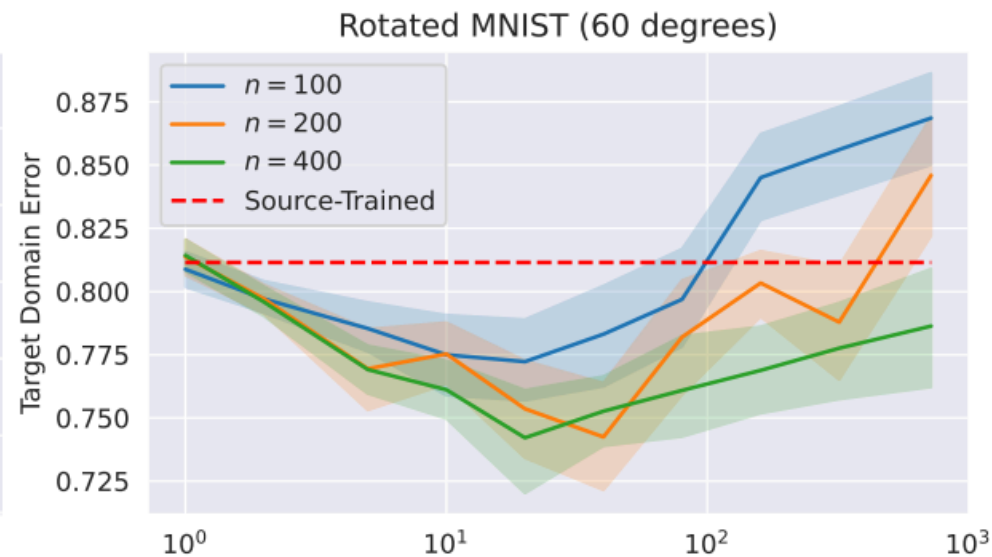
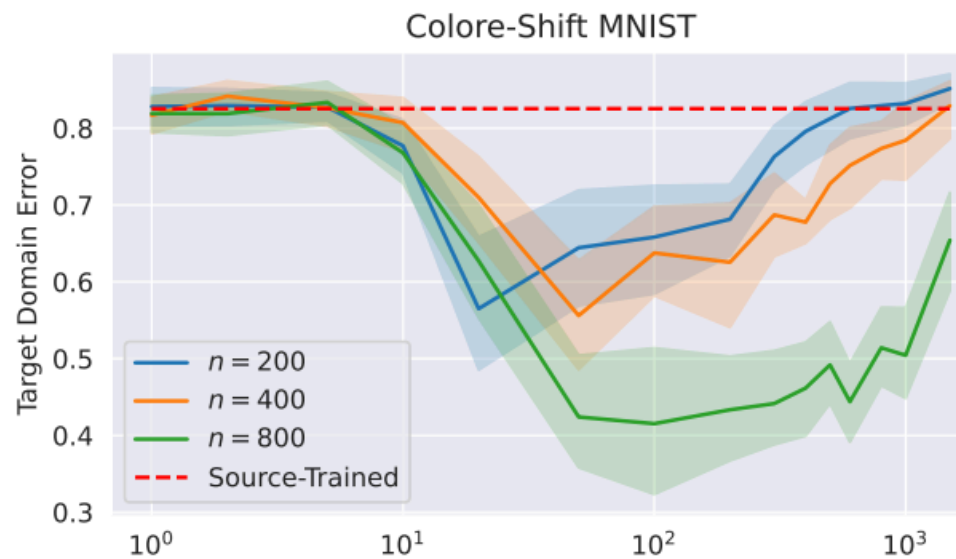


$$T^* = \max \left\{ \frac{L}{\Delta_{\max}}, \tilde{\mathcal{O}} \left(\left(\frac{1}{1 + \Delta_{\max} \sqrt{n}} \right)^{2/3} \right) \right\}$$

- Δ_{\max} : average distributional distance between consecutive domains

Is there an optimal T^* in practice?

$$T^* = \max \left\{ \frac{L}{\Delta_{\max}}, \tilde{\mathcal{O}} \left(\left(\frac{1}{1 + \Delta_{\max} \sqrt{n}} \right)^{2/3} \right) \right\}$$



Code: github.com/Haoxiang-Wang/gradual-domain-adaptation

Contact Information:

- Haoxiang Wang: hwang264@illinois.edu
- Bo Li: lbo@illinois.edu
- Han Zhao: hanzhao@illinois.edu