



Boulder

# Computer Vision; Image Classification; Domain Adaptation

---

**Maziar Raissi**

**Assistant Professor**

Department of Applied Mathematics

University of Colorado Boulder

[maziar.raissi@colorado.edu](mailto:maziar.raissi@colorado.edu)



Boulder

# Domain-Adversarial Training of Neural Networks

## Domain Adaptation

Data at training and test time come from similar but different distributions!

$X \rightarrow$  input space

$Y = \{0, 1, \dots, L - 1\} \rightarrow$  set of  $L$  possible labels

$\mathcal{D}_S \rightarrow$  source domain

$\mathcal{D}_T \rightarrow$  target domain

$\mathcal{D}_S, \mathcal{D}_T \rightarrow$  distributions over  $X \times Y$

$S \rightarrow$  labeled source sample drawn i.i.d from  $\mathcal{D}_S$

$T \rightarrow$  unlabeled target sample drawn i.i.d from  $\mathcal{D}_T^X$

$\mathcal{D}_T^X \rightarrow$  marginal distribution of  $\mathcal{D}_T$  over  $X$

$S = \{(x_i, y_i)\}_{i=1}^n \sim (\mathcal{D}_S)^n$

$T = \{x_i\}_{i=n+1}^N \sim (\mathcal{D}_T^X)^{n'}$

$N = n + n' \rightarrow$  total number of samples

Build a classifier  $\eta : X \rightarrow Y$  with a low target risk:

$$\mathcal{R}_{\mathcal{D}_T}(\eta) = \Pr_{(x,y) \sim \mathcal{D}_T}(\eta(x) \neq y)$$

## Theorem

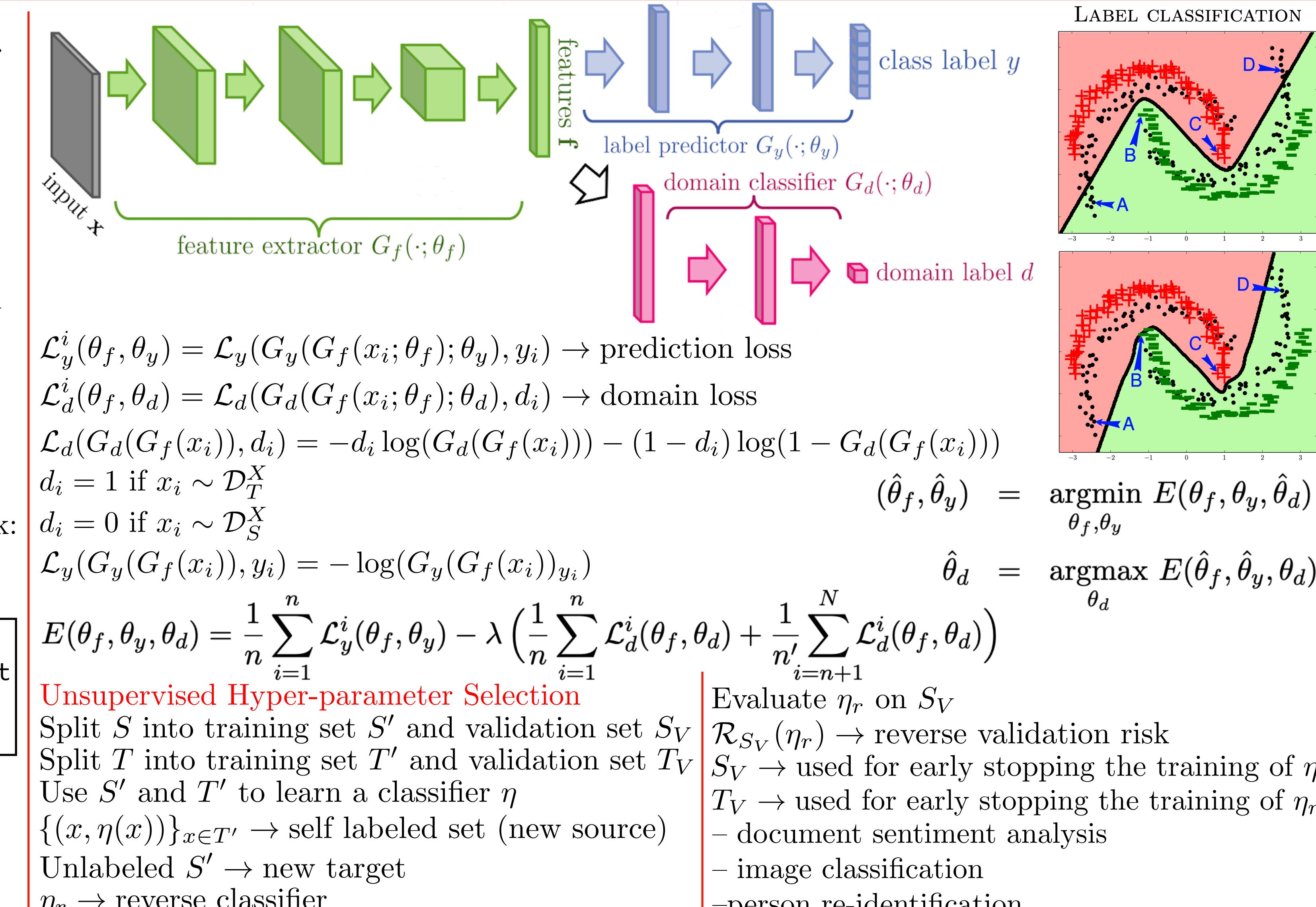
For effective domain transfer to be achieved, predictions must be made based on features that cannot discriminate between the training (source) and test (target) domains.

## Domain-Adversarial Neural Network (DANN)

$G_f(\cdot; \theta_f) \rightarrow$  feature extractor

$G_y(\cdot; \theta_y) \rightarrow$  label predictor

$G_d(\cdot; \theta_d) \rightarrow$  domain classifier



## Unsupervised Hyper-parameter Selection

Split  $S$  into training set  $S'$  and validation set  $S_V$

Split  $T$  into training set  $T'$  and validation set  $T_V$

Use  $S'$  and  $T'$  to learn a classifier  $\eta$

$\{(x, \eta(x))\}_{x \in T'} \rightarrow$  self labeled set (new source)

Unlabeled  $S' \rightarrow$  new target

$\eta_r \rightarrow$  reverse classifier

Evaluate  $\eta_r$  on  $S_V$

$\mathcal{R}_{S_V}(\eta_r) \rightarrow$  reverse validation risk

$S_V \rightarrow$  used for early stopping the training of  $\eta$

$T_V \rightarrow$  used for early stopping the training of  $\eta_r$

- document sentiment analysis

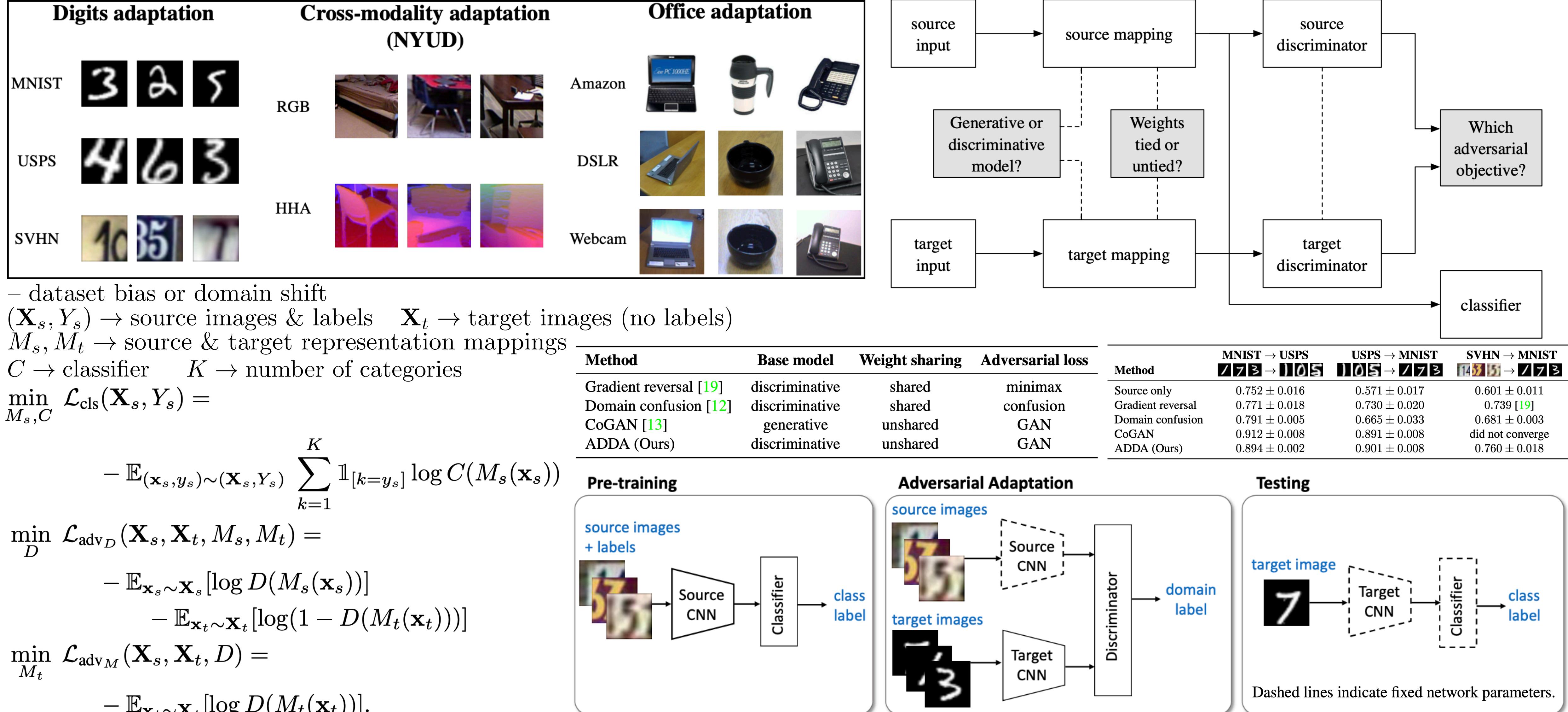
- image classification

- person re-identification



Boulder

# Adversarial Discriminative Domain Adaptation





Boulder

# CyCADA: Cycle-Consistent Adversarial Domain Adaptation

Adversarial Adaptation Models:

- discovering domain invariant representations (feature space methods)
- mapping between unpaired image domains (image space methods)

## Unsupervised Adaptation

$X_S \rightarrow$  source data       $Y_S \rightarrow$  source labels  
 $X_T \rightarrow$  target data      no target labels!

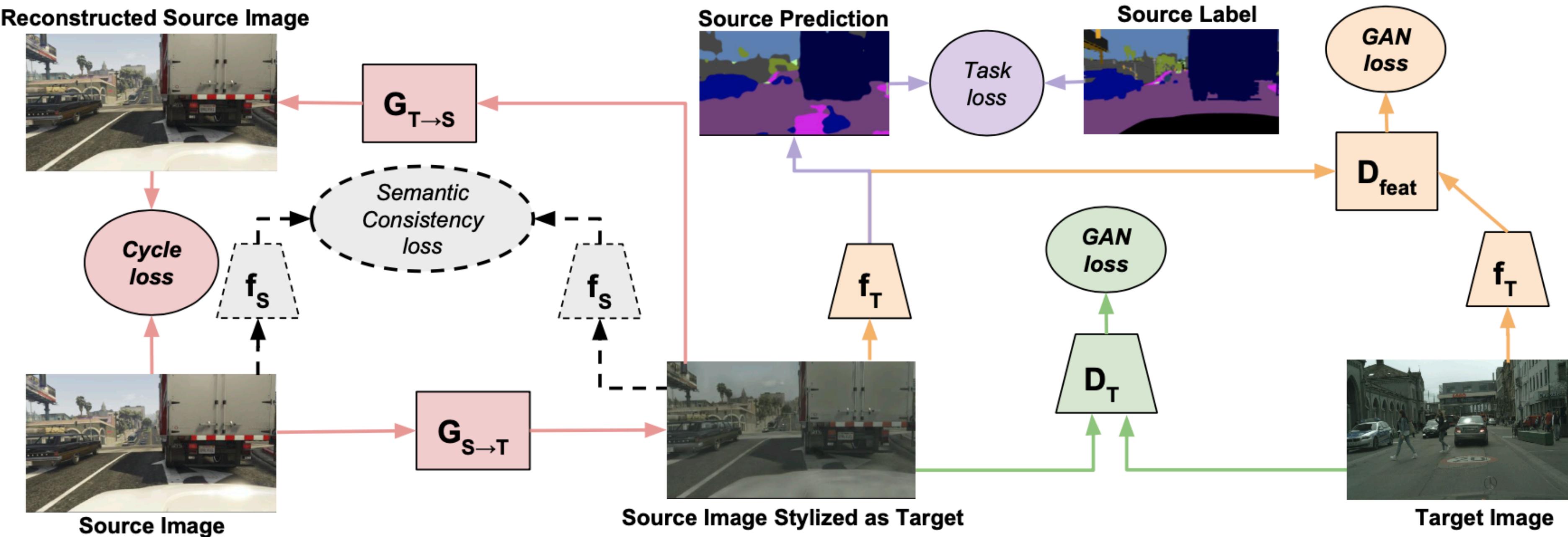
Learn a model  $f_T$  that correctly predicts the label for the target data  $X_T$

## Pre-train Source Task Model

$$\mathcal{L}_{\text{task}}(f_S, X_S, Y_S) = -\mathbb{E}_{(x_s, y_s) \sim (X_S, Y_S)} \sum_{k=1}^K \mathbb{1}_{[k=y_s]} \log \left( \sigma(f_S^{(k)}(x_s)) \right)$$

softmax function

domain shift leads to reduced performance when evaluating on target data



## Pixel-level Adaptation

$G_{S \rightarrow T} \rightarrow$  mapping from source to target

$\cup$  trained to produce target samples that fool adversarial discriminator  $D_T$

$$\mathcal{L}_{\text{GAN}}(G_{S \rightarrow T}, D_T, X_T, X_S) = \mathbb{E}_{x_t \sim X_T} [\log D_T(x_t)] + \mathbb{E}_{x_s \sim X_S} [\log(1 - D_T(G_{S \rightarrow T}(x_s)))]$$

$$\mathcal{L}_{\text{task}}(f_T, G_{S \rightarrow T}(X_S), Y_S) \rightarrow \text{learn a target model } f_T$$

No way to guarantee that  $G_{S \rightarrow T}(x_s)$  preserves the structure and content of  $x_s$ !

$$\mathcal{L}_{\text{GAN}}(G_{T \rightarrow S}, D_S, X_S, X_T)$$

$$\begin{aligned} \mathcal{L}_{\text{cyc}}(G_{S \rightarrow T}, G_{T \rightarrow S}, X_S, X_T) = & \mathbb{E}_{x_s \sim X_S} [||G_{T \rightarrow S}(G_{S \rightarrow T}(x_s)) - x_s||_1] \\ & + \mathbb{E}_{x_t \sim X_T} [||G_{S \rightarrow T}(G_{T \rightarrow S}(x_t)) - x_t||_1] \end{aligned}$$

## Prevent Label Flipping

semantic consistency  $\leftarrow \mathcal{L}_{\text{sem}}(G_{S \rightarrow T}, G_{T \rightarrow S}, X_S, X_T, f_S) =$

$$\mathcal{L}_{\text{task}}(f_S, G_{T \rightarrow S}(X_T), p(f_S, X_T))$$

$$+ \mathcal{L}_{\text{task}}(f_S, G_{S \rightarrow T}(X_S), p(f_S, X_S))$$

$p(f, X) = \arg \max f(x)$

$\cup$  predicted label

$f_S \rightarrow$  pretrained and fixed

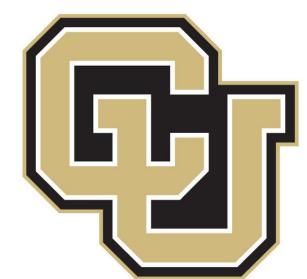
## Feature-level Adaptation

$$\mathcal{L}_{\text{GAN}}(f_T, D_{\text{feat}}, f_S(G_{S \rightarrow T}(X_S)), X_T)$$

## Complete Objective

$$f_T^* = \arg \min_{f_T} \min_{G_{S \rightarrow T}} \max_{D_S, D_T} \mathcal{L}_{\text{CyCADA}}$$

Model	Accuracy (%)
Source only	67.1
CyCADA - no feat adapt, no semantic loss	70.3
CyCADA - no feat adapt	71.2
CyCADA - no cycle consistency	75.7
CyCADA - no pixel adapt	83.8
CyCADA (Full)	<b>90.4</b>
Target Fully Supervised	99.2



Boulder

# Questions?

---