

# **Learning Semantic Representations for Unsupervised Domain Adaptation**

Moving Semantic Transfer Network (MSTN)

$$f : X_T \rightarrow Y_T$$

**“Ultimate goal is to develop a deep neural network that is able to predict labels for the samples from target domain”**

$X_t$  target sample     $Y_t$  target label  
 $X_s$  source sample     $Y_s$  source label

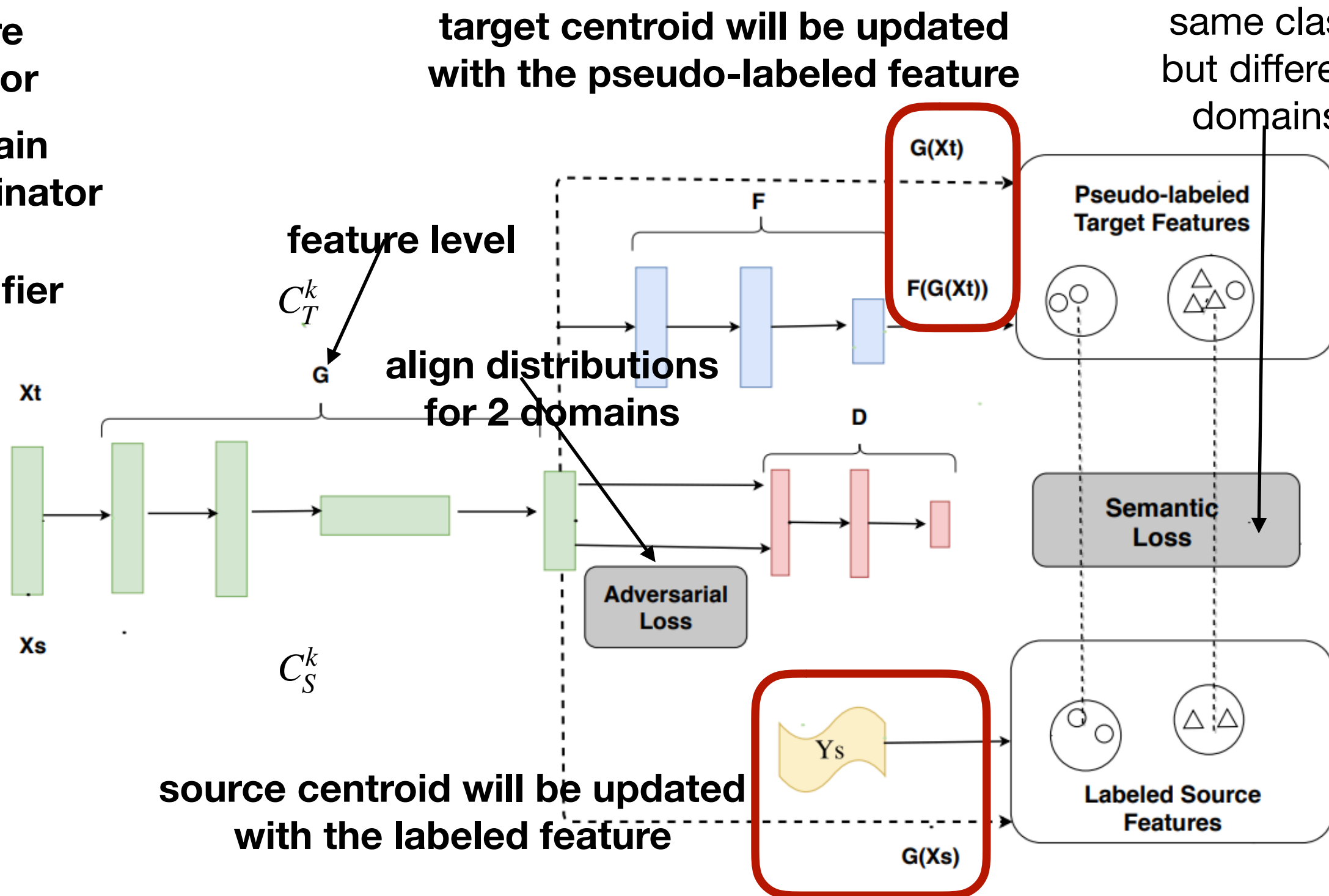
1. Standard source classification loss
2. Domain adversarial loss
3. Semantic Loss

$G$  Feature Extractor

$D$  Domain Discriminator

$F$  Classifier

align embedding by restricting the distance between centroid in same class but different domains



$$f = F \circ G$$

a visual classifier is trained by minimizing the source classification error and the discrepancy between source domain and target domain

cross entropy loss  
(source classification error)

domain adversarial similarity loss  
(discrepancy between source/target domain)

$$\mathcal{L}(\mathcal{X}_S, \mathcal{Y}_S, \mathcal{X}_T) = \mathcal{L}_C(\mathcal{X}_S, \mathcal{Y}_S) + \lambda \mathcal{L}_{DC}(\mathcal{X}_S, \mathcal{X}_T) + \gamma \mathcal{L}_{SM}(\mathcal{X}_S, \mathcal{Y}_S, \mathcal{X}_T),$$

making alignment semantic  
(centroid alignment)

$$\mathcal{L} = \underbrace{\mathbb{E}_{(x,y) \sim D_S} [J(f(x), y)]}_{\mathcal{L}_C(\mathcal{X}_S, \mathcal{Y}_S)} + \lambda \underbrace{d(\mathcal{X}_S, \mathcal{X}_T)}_{\mathcal{L}_{DC}(\mathcal{X}_S, \mathcal{X}_T)} \quad (1)$$

1. **Standard source classification loss (cross entropy loss) + divergence between two domains (domain adversarial similarity loss)**

$$d(\mathcal{X}_S, \mathcal{X}_T) = \frac{\mathbb{E}_{x \sim D_S} [\log(1 - D \circ G(x))] + \mathbb{E}_{x \sim D_T} [\log(D \circ G(x))]}{2} \quad (2)$$

**domain adversarial similarity loss (Domain classifier D to tell whether features from G arise from source or target domain)**

**discriminability**

$$\mathcal{L}_{SM}^{UDA}(\mathcal{X}_S, \mathcal{Y}_S, \mathcal{X}_T) = \underbrace{\sum_{k=1}^K \Phi(C_S^k, C_T^k)}_{\mathcal{L}_{SM}(\mathcal{X}_S, \mathcal{Y}_S, \mathcal{X}_T)}, \quad (4)$$

3. **making alignment semantic with pseudo labeled target domain (centroid alignment)**

$X_t$  target sample     $Y_t$  target label  
 $X_s$  source sample     $Y_s$  source label

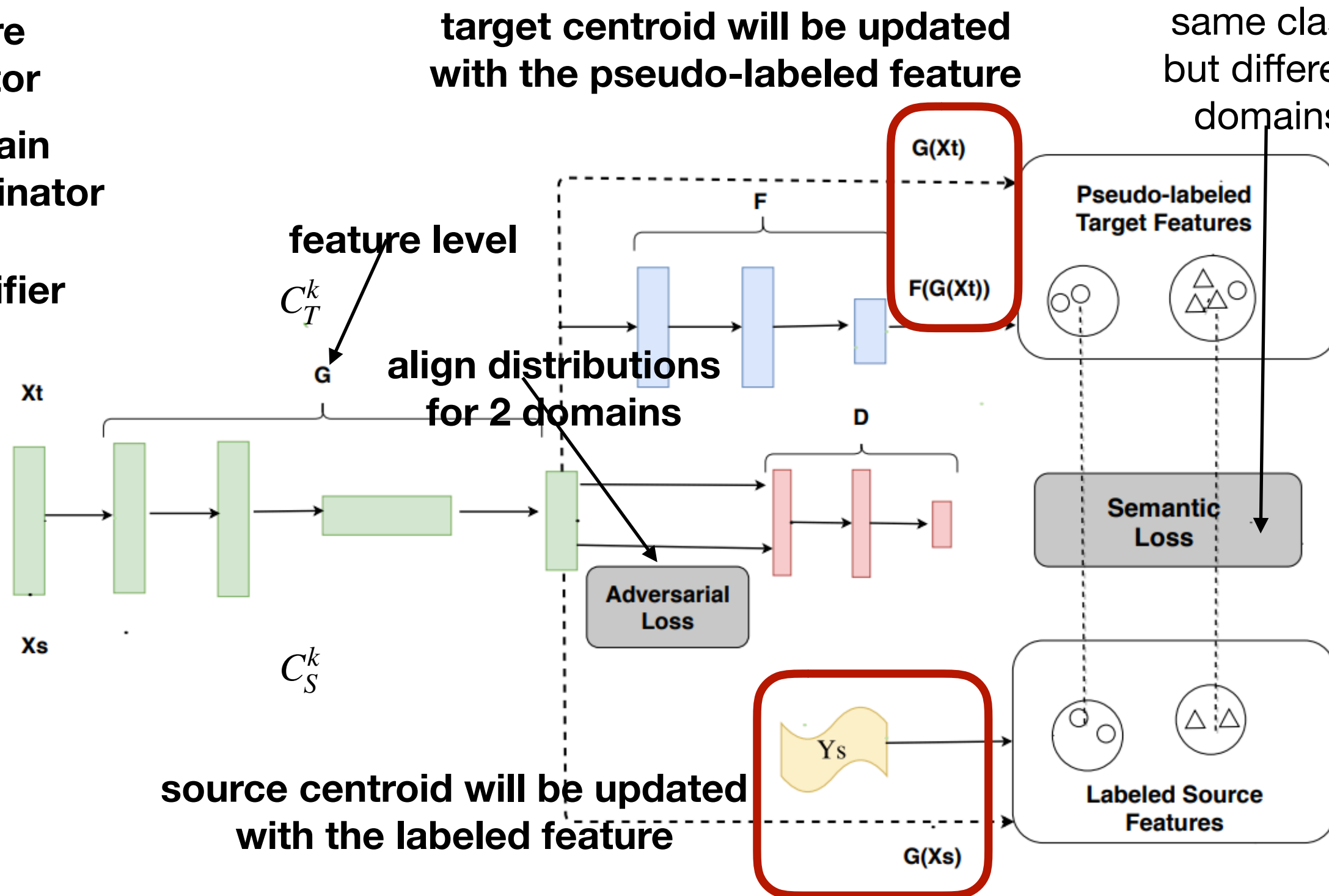
1. Standard source classification loss
2. Domain adversarial loss
3. Semantic Loss

$G$  Feature Extractor

$D$  Domain Discriminator

$F$  Classifier

align embedding by restricting the distance between centroid in same class but different domains



---

**Algorithm 1** Moving semantic transfer loss computation in iteration  $t$  in our model.  $K$  is the number of classes.

---

**Input:** Labeled set  $S$ , Unlabeled set  $T$ ,  $N$  is the batch size, Training classifier  $f$ , Global centroids for two domains:

$\{C_S^k\}_{k=1}^K$  and  $\{C_T^k\}_{k=1}^K$

- 1:  $S_t = \text{RANDOMSAMPLE}(S, N)$
  - 2:  $\underline{T}_t = \text{RANDOMSAMPLE}(T, N)$
  - 3:  $\widehat{T}_t = \text{Labeling}(G, f, T_t)$
  - 4:  $\mathcal{L}_{SM} = 0$
  - 5: **for**  $k = 1$  to  $K$  **do** feature extraction
  - 6:  $C_{S(t)}^k \leftarrow \frac{1}{|S_t^k|} \sum_{(x_i, y_i) \in S_t^k} G(x_i)$  (From Scratch)
  - 7:  $C_{\widehat{T}(t)}^k \leftarrow \frac{1}{|\widehat{T}_t^k|} \sum_{(x_i, y_i) \in \widehat{T}_t^k} G(x_i)$  (From Scratch)
  - 8:  $C_S^k \leftarrow \theta C_S^k + (1 - \theta) C_{S(t)}^k$  (Moving Average) align distributions
  - 9:  $C_T^k \leftarrow \theta C_T^k + (1 - \theta) C_{\widehat{T}(t)}^k$  (Moving Average)
  - 10:  $\mathcal{L}_{SM} \leftarrow \mathcal{L}_{SM} + \Phi(C_S^k, C_T^k)$  centroid alignment
  - 11: **end for**
  - 12: **return**  $\mathcal{L}_{SM}$
-