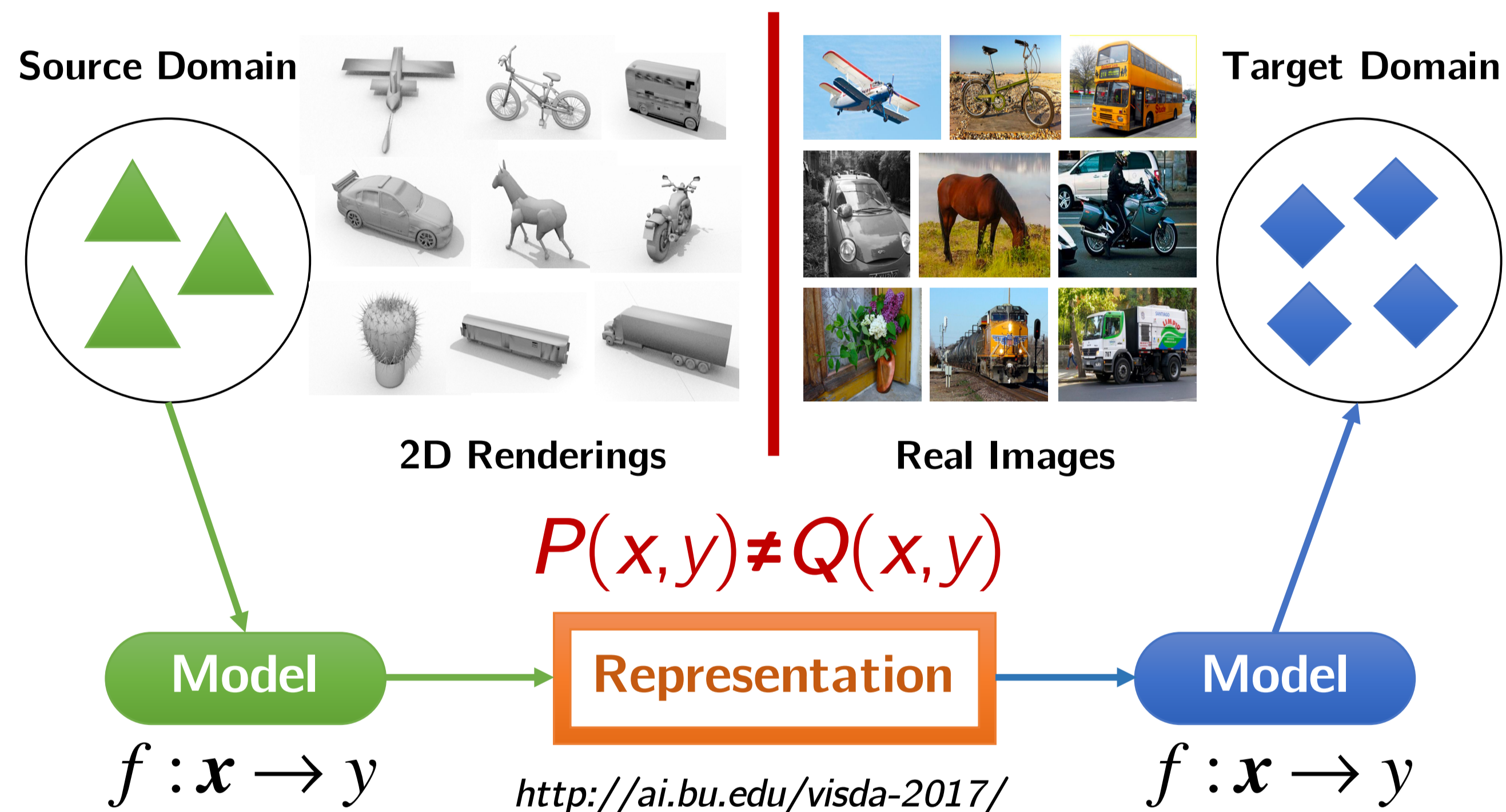


## Summary

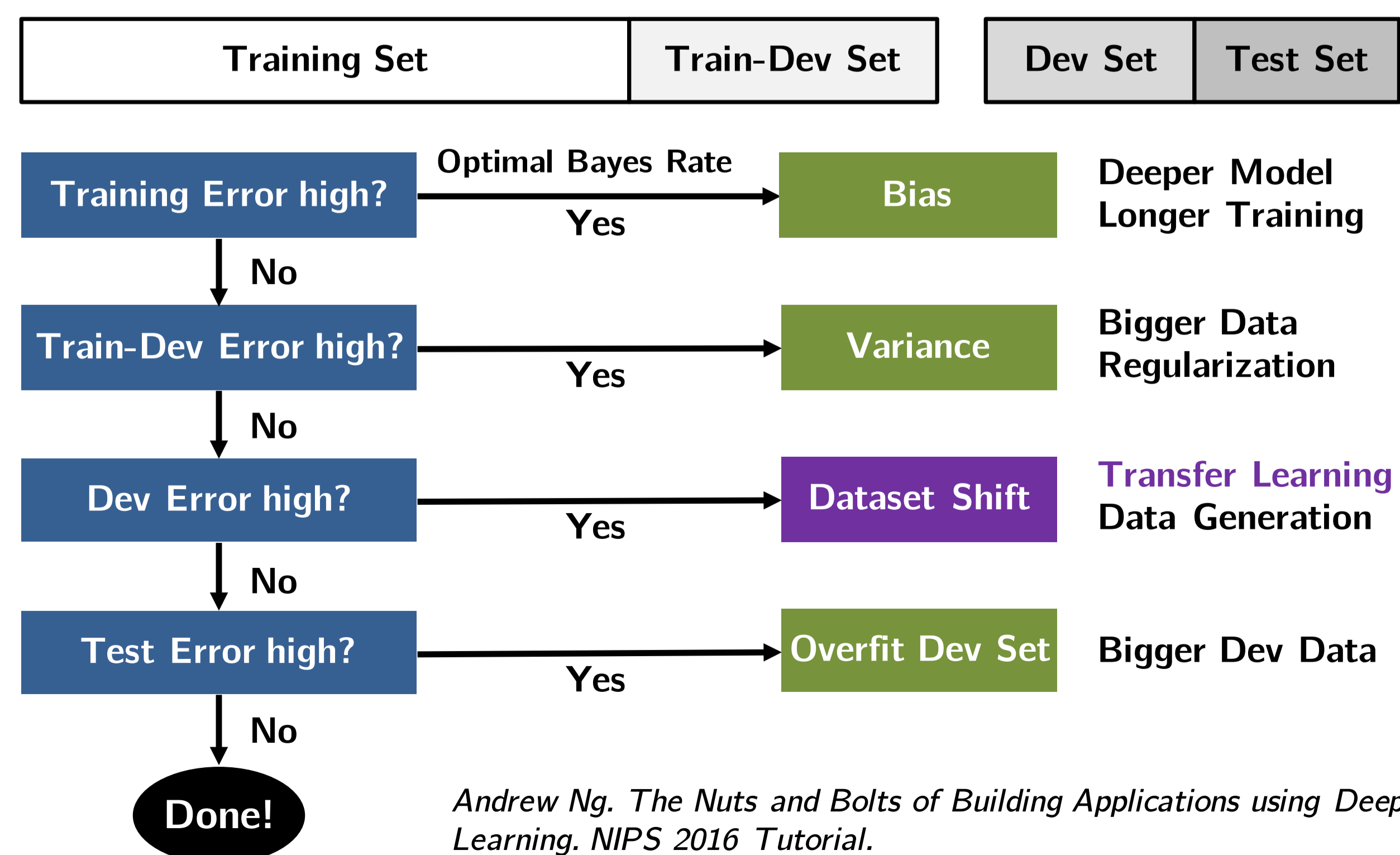
- ▶ A joint adaptation network framework for deep transfer learning
- ▶ Two main contributions:
  - ▶ **Joint** adaptation of multilayer features and classifier predictions
  - ▶ **Adversarial** adaptation with semi-parametric domain discriminator
- ▶ State-of-the-art results on visual & simulation-to-real datasets
- ▶ Open Problems
  - ▶ Randomized method for the multilinear operation across feature maps
  - ▶ Kernel approximation of the universal kernel for distribution matching
- ▶ Code@: <https://github.com/thuml/transfer-caffe>

## Deep Transfer Learning

- ▶ Deep learning across domains of different distributions  $P \neq Q$

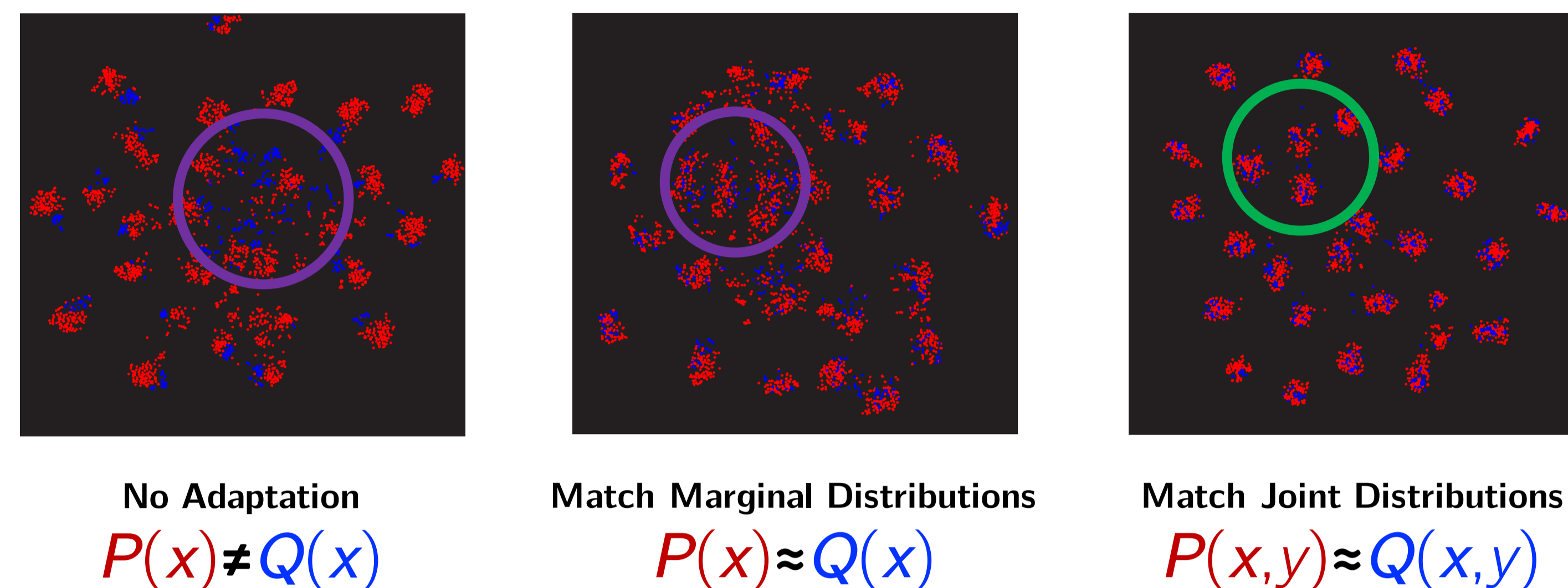


## Deep Transfer Learning: Why?

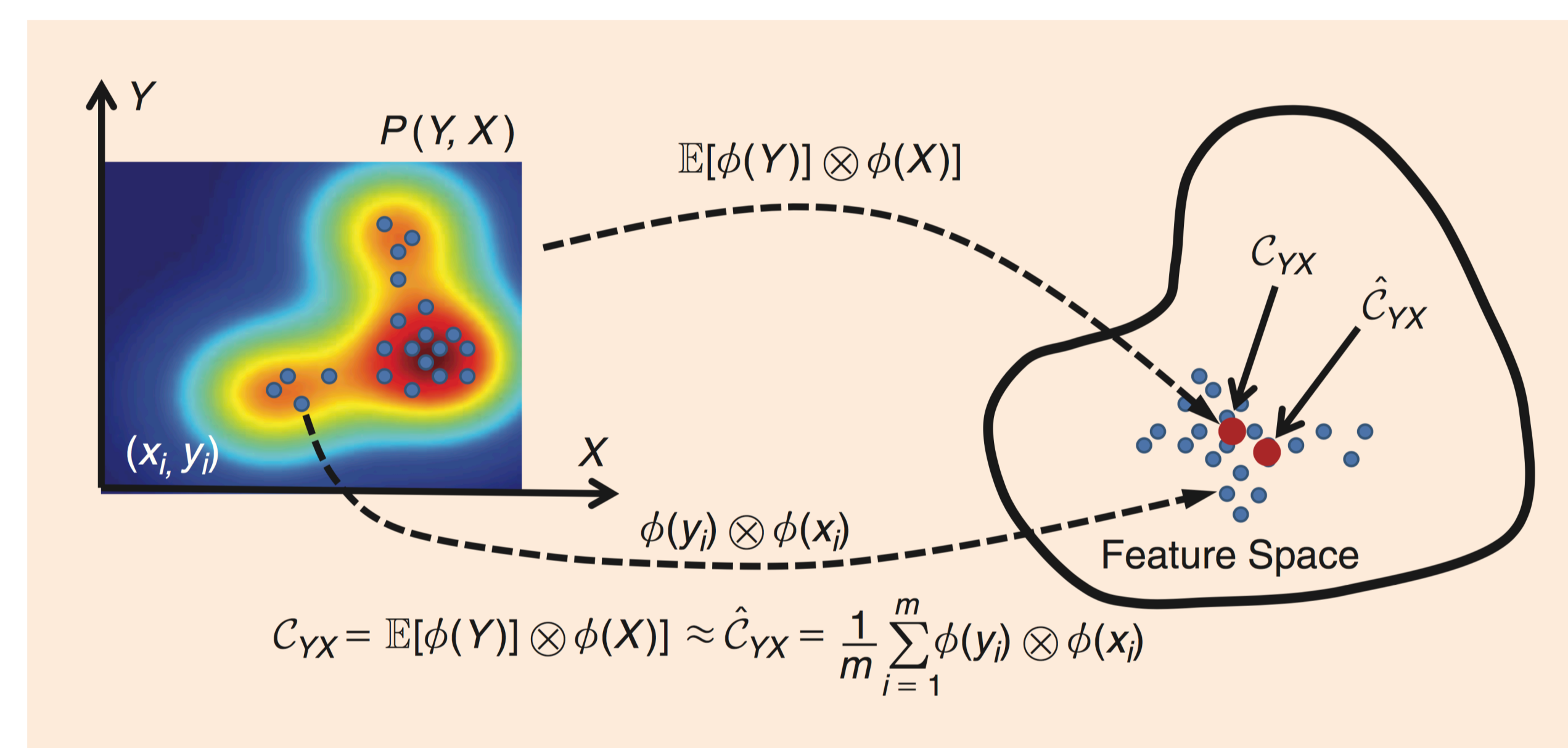


## Main Idea of This Work

- ▶ Directly model and match joint distributions  $P(\mathbf{x}, y) & Q(\mathbf{x}, y)$



## Kernel Embedding of Joint Distributions



$$C_{\mathbf{X}^{1:m}}(P) \triangleq \mathbb{E}_{\mathbf{X}^{1:m}} \left[ \otimes_{\ell=1}^m \phi^\ell(\mathbf{X}^\ell) \right] \approx \hat{C}_{\mathbf{X}^{1:m}} = \frac{1}{n} \sum_{i=1}^n \otimes_{\ell=1}^m \phi^\ell(\mathbf{x}_i^l) \quad (1)$$

Le Song et al. *Kernel Embeddings of Conditional Distributions*. IEEE, 2013.

## Joint Maximum Mean Discrepancy (JMMD)

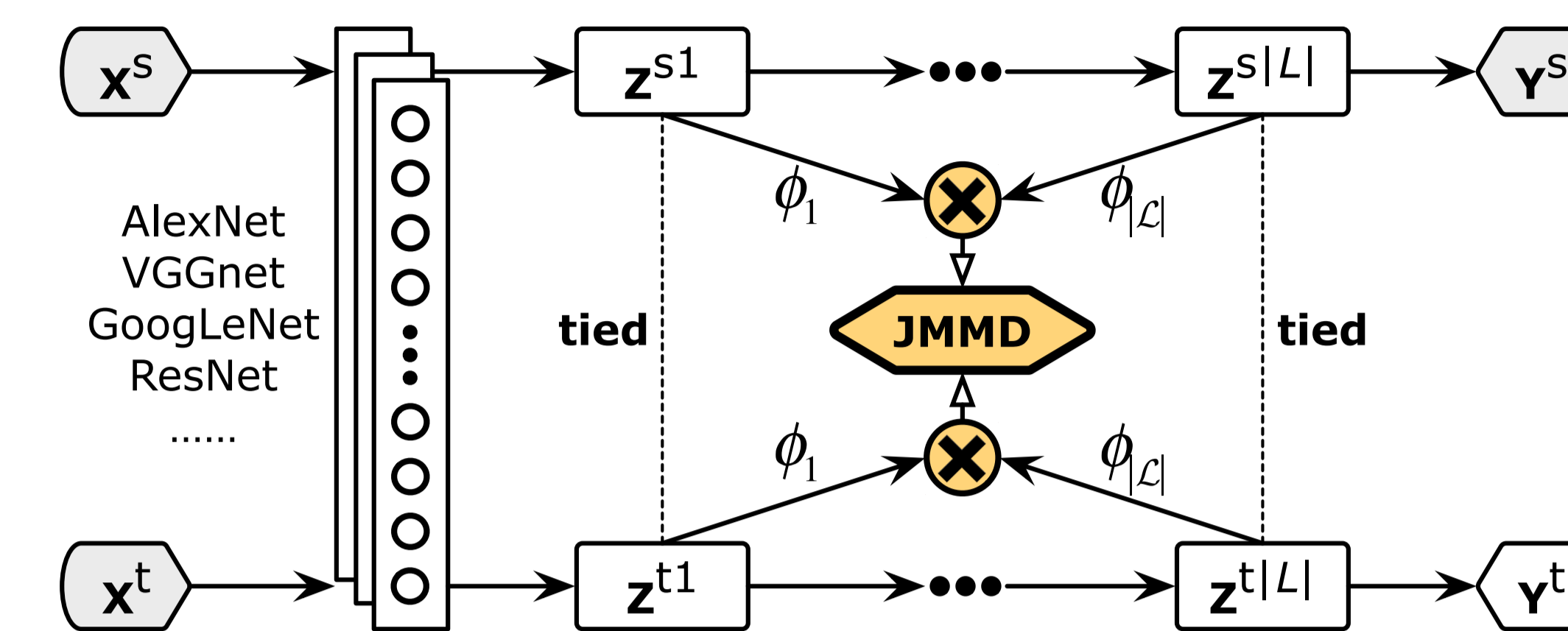
Distance between embeddings  $P(\mathbf{Z}^{s1}, \dots, \mathbf{Z}^{s|L|})$   $Q(\mathbf{Z}^{t1}, \dots, \mathbf{Z}^{t|L|})$

$$D_{\mathcal{L}}(P, Q) \triangleq \|\mathbf{C}_{\mathbf{Z}^{s,1:|L|}}(P) - \mathbf{C}_{\mathbf{Z}^{t,1:|L|}}(Q)\|_{\otimes_{\ell=1}^{|L|} \mathcal{H}^\ell} \quad (2)$$

$$\begin{aligned} \hat{D}_{\mathcal{L}}(P, Q) &= \frac{1}{n_s^2} \sum_{i=1}^{n_s} \sum_{j=1}^{n_s} \prod_{\ell \in \mathcal{L}} k^\ell(\mathbf{z}_i^{s\ell}, \mathbf{z}_j^{s\ell}) \\ &+ \frac{1}{n_t^2} \sum_{i=1}^{n_t} \sum_{j=1}^{n_t} \prod_{\ell \in \mathcal{L}} k^\ell(\mathbf{z}_i^{t\ell}, \mathbf{z}_j^{t\ell}) \\ &- \frac{2}{n_s n_t} \sum_{i=1}^{n_s} \sum_{j=1}^{n_t} \prod_{\ell \in \mathcal{L}} k^\ell(\mathbf{z}_i^{s\ell}, \mathbf{z}_j^{t\ell}). \end{aligned} \quad (3)$$

- ▶  $P = Q$  iff.  $\hat{D}_{\mathcal{L}}(P, Q) = 0$  (In practice,  $\hat{D}_{\mathcal{L}}(P, Q) < \epsilon$ )

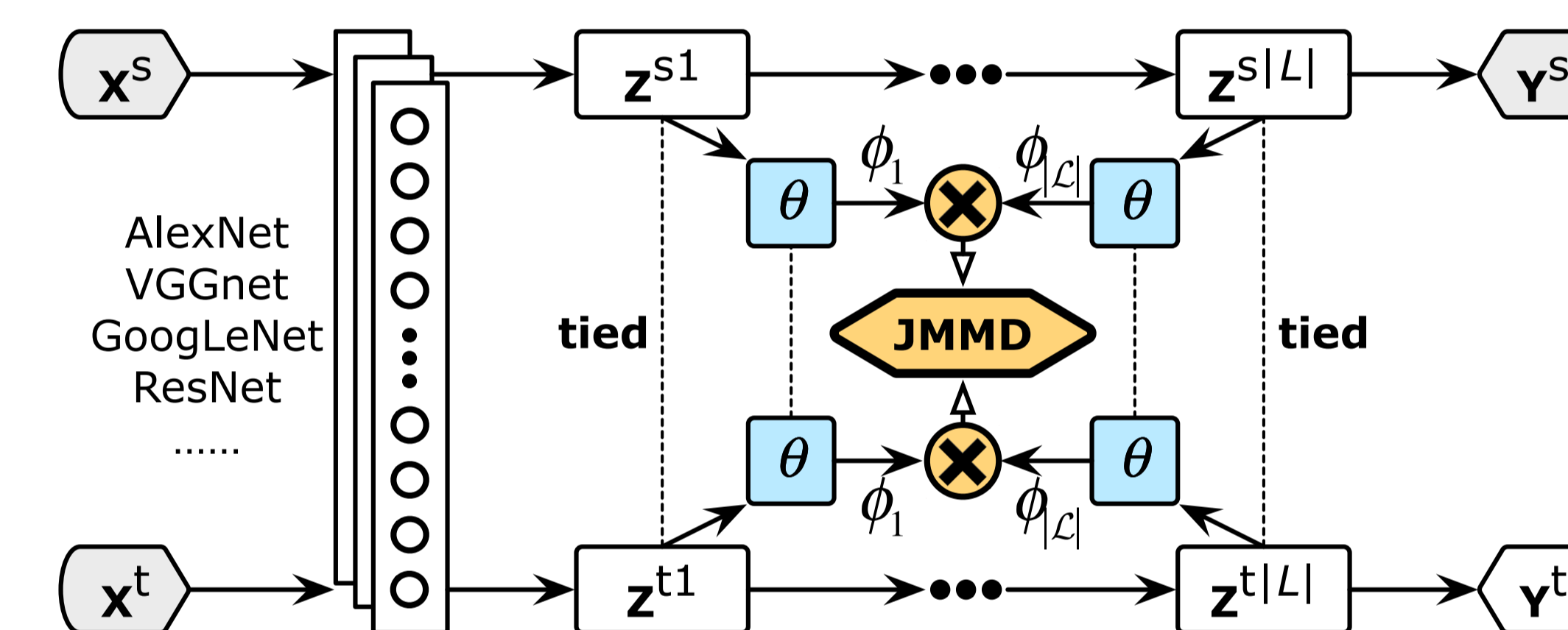
## Joint Adaptation Network (JAN)



Joint adaptation: match joint distributions of features/predictions

$$\min_f \frac{1}{n_s} \sum_{i=1}^{n_s} J(f(\mathbf{x}_i^s), \mathbf{y}_i^s) + \lambda \hat{D}_{\mathcal{L}}(P, Q) \quad (4)$$

## Adversarial Joint Adaptation Network (JAN-A)



Optimal matching: maximize JMMD as semi-parametric adversary

$$\min_f \max_{\theta} \frac{1}{n_s} \sum_{i=1}^{n_s} J(f(\mathbf{x}_i^s), \mathbf{y}_i^s) + \lambda \hat{D}_{\mathcal{L}}(P, Q; \theta) \quad (5)$$

## Experimental Results

Method	A → W	D → W	W → D	A → D	D → A	W → A	Avg
AlexNet	61.6±0.5	95.4±0.3	99.0±0.2	63.8±0.5	51.1±0.6	49.8±0.4	70.1
RevGrad	73.0±0.5	96.4±0.3	99.2±0.3	72.3±0.3	53.4±0.4	51.2±0.5	74.3
JAN	74.9±0.3	96.6±0.2	99.5±0.2	71.8±0.2	58.3±0.3	55.0±0.4	76.0
JAN-A	75.2±0.4	96.6±0.2	99.6±0.1	72.8±0.3	57.5±0.2	56.3±0.2	76.3
ResNet	68.4±0.2	96.7±0.1	99.3±0.1	68.9±0.2	62.5±0.3	60.7±0.3	76.1
RevGrad	82.0±0.4	96.9±0.2	99.1±0.1	79.7±0.4	68.2±0.4	67.4±0.5	82.2
JAN	85.4±0.3	97.4±0.2	99.8±0.2	84.7±0.3	68.6±0.3	70.0±0.4	84.3
JAN-A	86.0±0.4	96.7±0.3	99.7±0.1	85.1±0.4	69.2±0.4	70.7±0.5	84.6

