



---

# Recent Advances in Transfer Learning

---

Mingsheng Long  
Tsinghua University  
Jan 25, 2018

<https://github.com/thumtl>

# Machine Learning



Learner:  $f : x \rightarrow y$

Distribution:  $(x, y) \sim P(x, y)$



fish

bird

mammal

tree

flower

.....

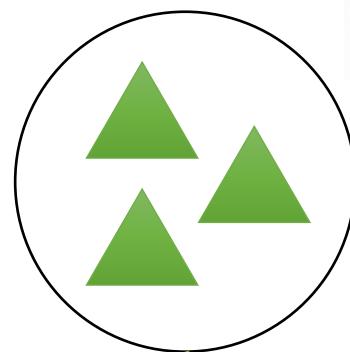
Error Bound:  $\epsilon_{\text{test}} \leq \hat{\epsilon}_{\text{train}} + \sqrt{\frac{\text{complexity}}{n}}$

# Transfer Learning



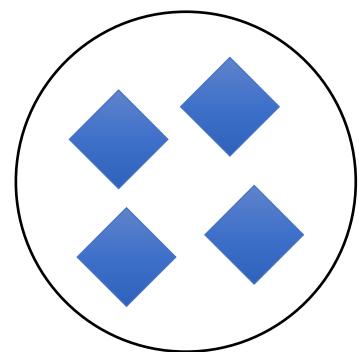
Learning across domains with **non-IID** distributions  $P \neq Q$

Source Domain



2D Renderings

Target Domain



Real Images

$$P(x,y) \neq Q(x,y)$$

Representation

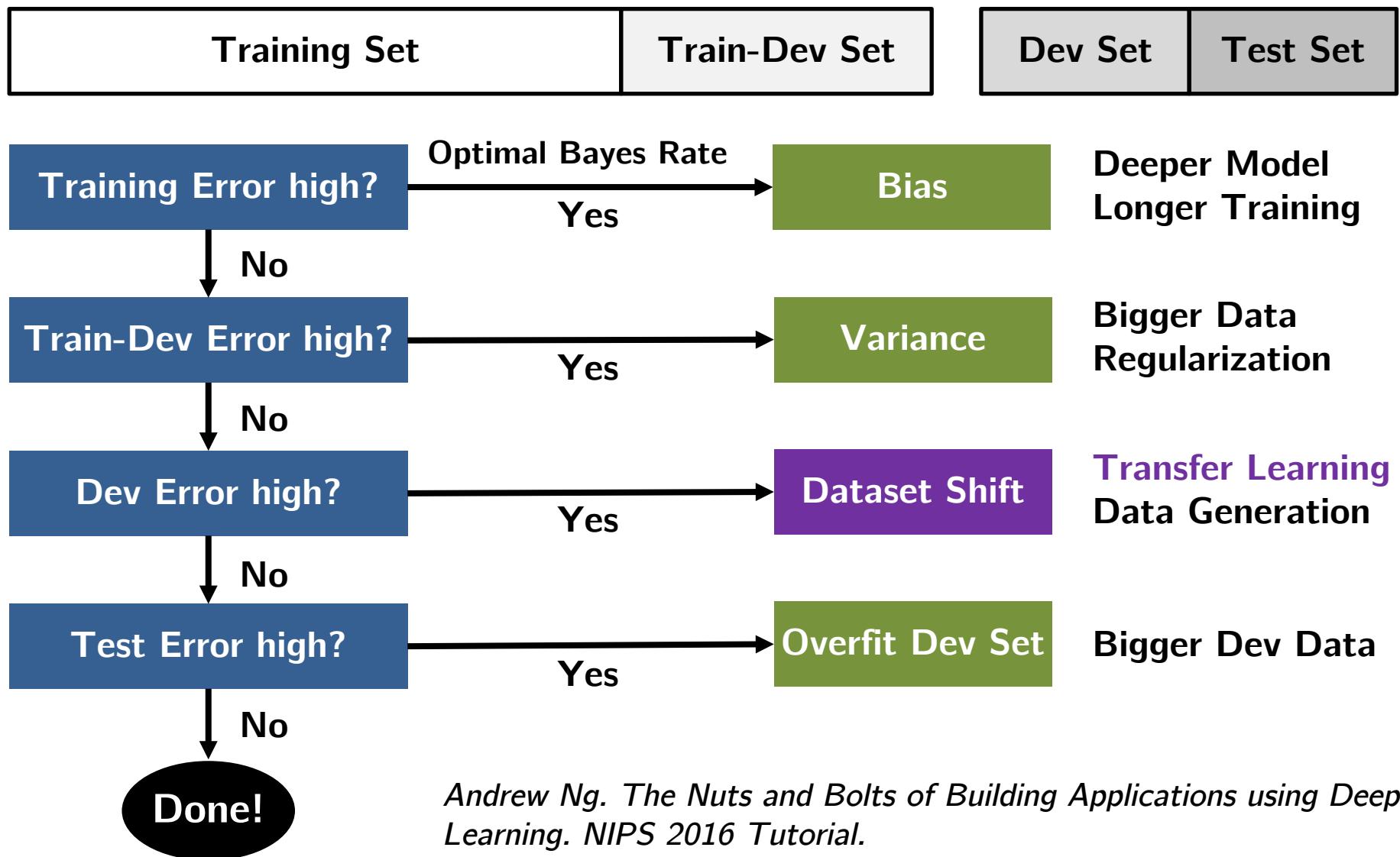
Model

$$f : x \rightarrow y$$

Model

$$f : x \rightarrow y$$

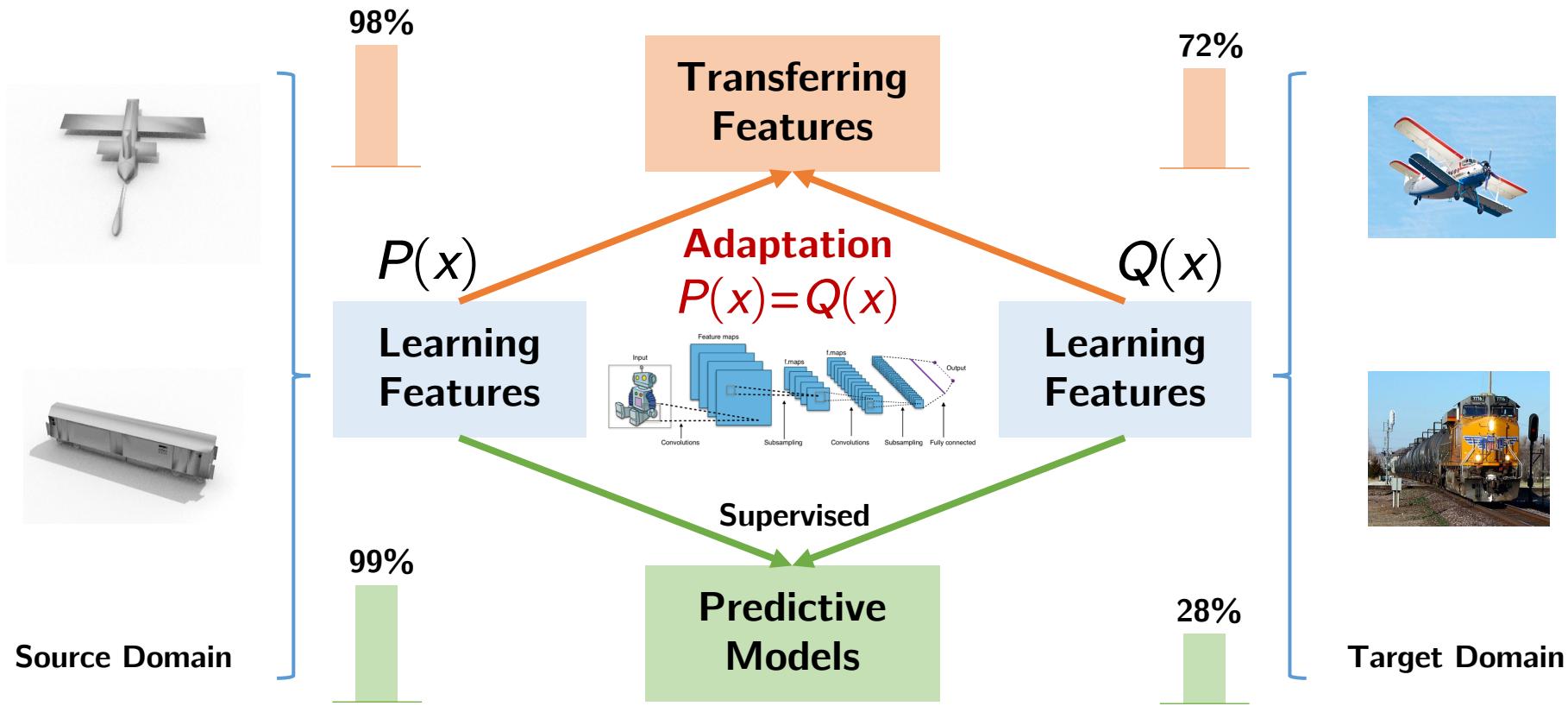
# Transfer Learning: Why?



# Transfer Learning: How?

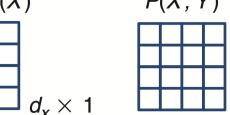
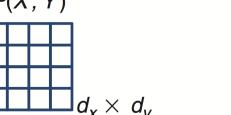
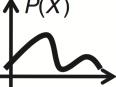
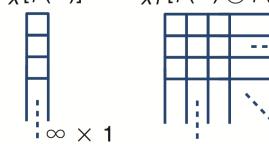
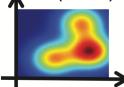
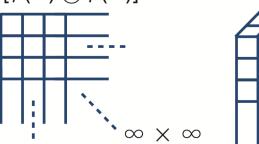
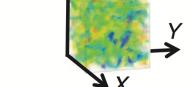
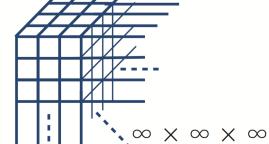


- Learning predictive models on transferable features s.t.  $P(x)=Q(x)$
- Distribution matching: MMD (ICML'15), GAN (ICML'15, JMLR'16)



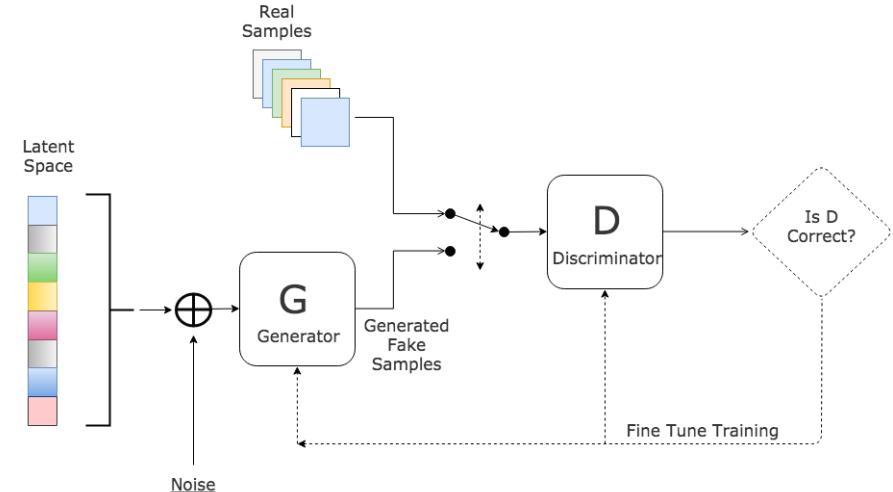
# Distribution Matching

- Marginal distribution mismatch:  $P(\mathbf{x}) \neq Q(\mathbf{x})$
- Conditional distribution mismatch:  $P(\mathbf{y}|\mathbf{x}) \neq Q(\mathbf{y}|\mathbf{x})$

	Distributions		
Discrete	$P(X)$  $d_x \times 1$	$P(X, Y)$  $d_x \times d_y$	$P(X, Y, Z)$  $d_x \times d_y \times d_z$
Kernel Embedding	$P(X)$  $\mu_X := \mathbb{E}_X[\phi(X)]$  $\infty \times 1$	$P(X, Y)$  $\mathbb{E}_{XY}[\phi(X) \otimes \phi(Y)]$  $\infty \times \infty$	$P(X, Y, Z)$  $\mathbb{E}_{XYZ}[\phi(X) \otimes \phi(Y) \otimes \phi(Z)]$  $\infty \times \infty \times \infty$

Kernel Embedding

Generative Adversarial Network

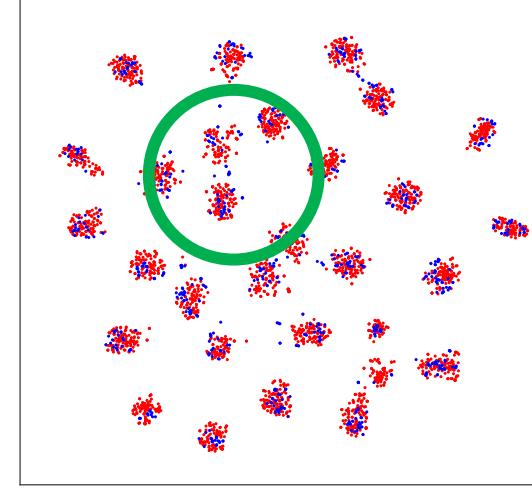
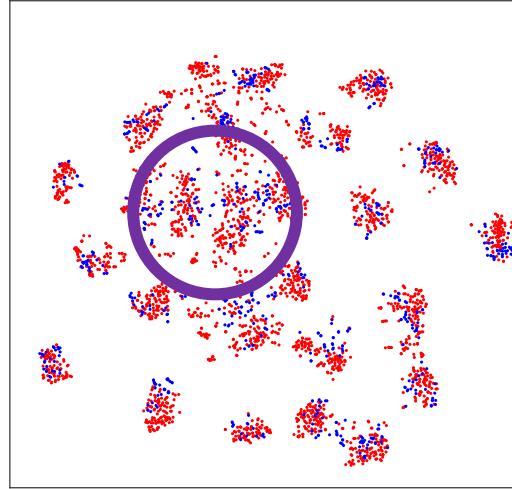
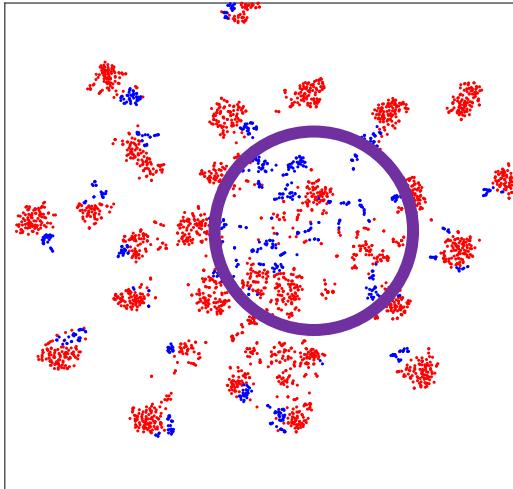


Adversarial Learning

Song et al. Kernel Embeddings of Conditional Distributions. *IEEE*, 2013.  
 Goodfellow et al. Generative Adversarial Networks. *NIPS* 2014.

# Distribution Matching

- Marginal distribution mismatch:  $P(\mathbf{x}) \neq Q(\mathbf{x})$
- Conditional distribution mismatch:  $P(\mathbf{y}|\mathbf{x}) \neq Q(\mathbf{y}|\mathbf{x})$



$$\begin{aligned} P(\mathbf{x}) &\neq Q(\mathbf{x}) \\ P(\mathbf{y}|\mathbf{x}) &\neq Q(\mathbf{y}|\mathbf{x}) \end{aligned}$$

$$\begin{aligned} P(\mathbf{x}) &\approx Q(\mathbf{x}) \\ P(\mathbf{y}|\mathbf{x}) &\neq Q(\mathbf{y}|\mathbf{x}) \end{aligned}$$

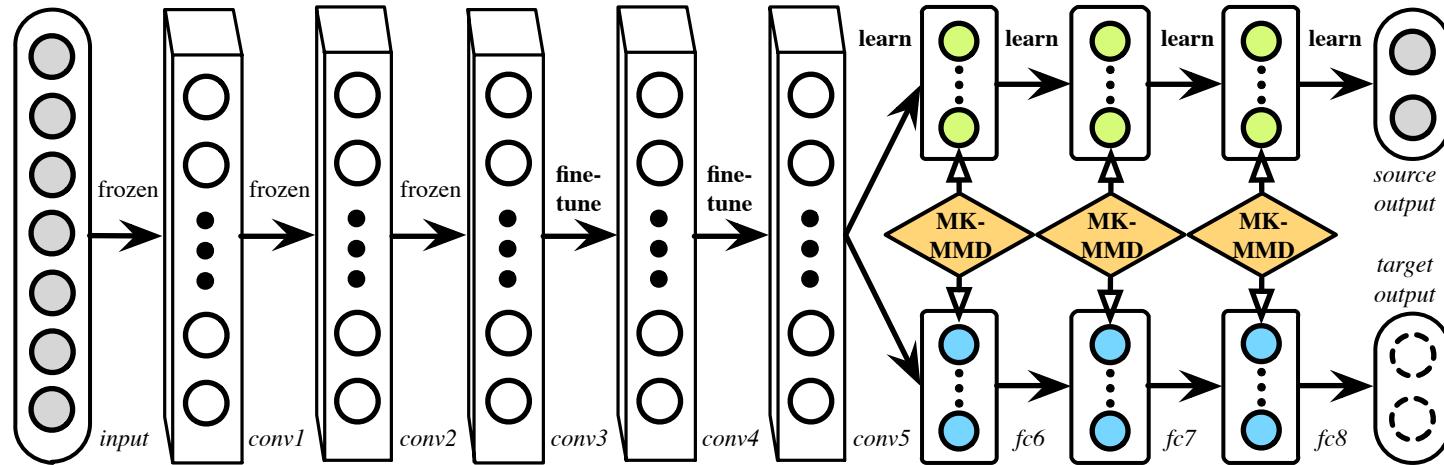
$$\begin{aligned} P(\mathbf{x}, \mathbf{y}) &\approx Q(\mathbf{x}, \mathbf{y}) \\ P(\mathbf{y}|\mathbf{x}) &\approx Q(\mathbf{y}|\mathbf{x}) \end{aligned}$$

# Problem 1



$$P(x) \neq Q(x)$$

# Deep Adaptation Network (DAN)

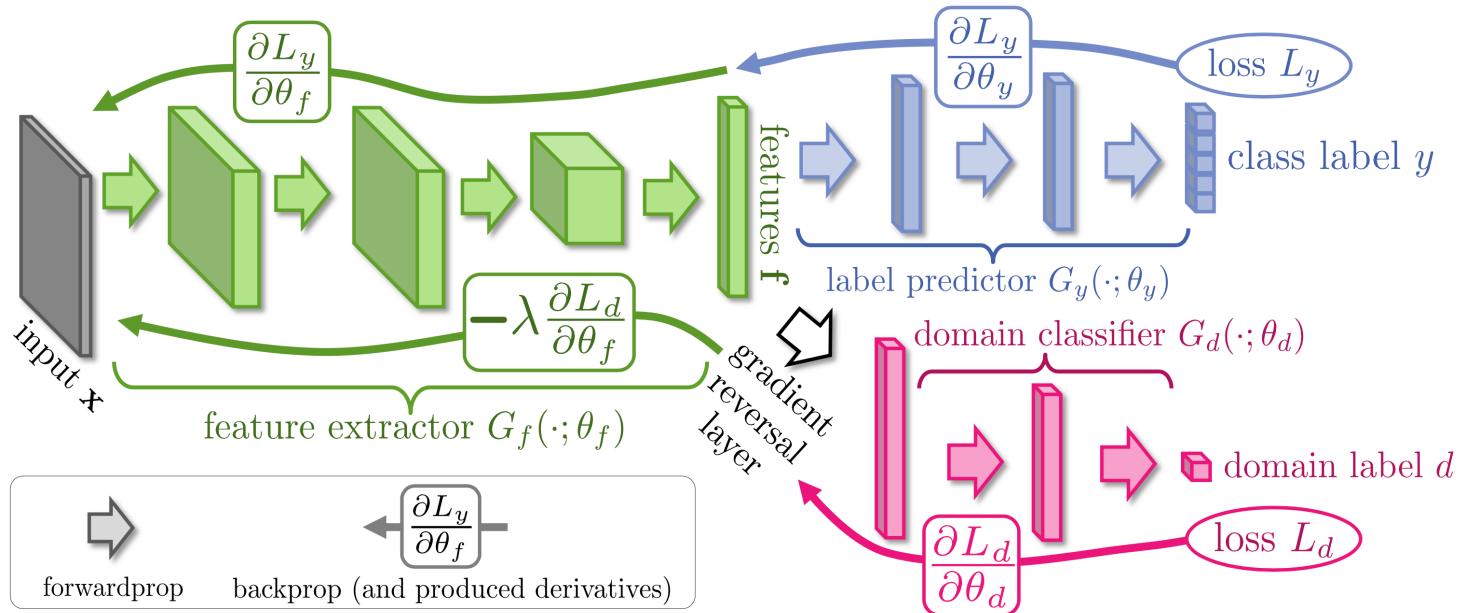


- Deep adaptation: match distributions in multiple domain-specific layers
- Optimal matching: maximize two-sample test power by multiple kernels

$$d_k^2(P, Q) \triangleq \left\| \mathbf{E}_P \left[ \phi(\mathbf{x}^s) \right] - \mathbf{E}_Q \left[ \phi(\mathbf{x}^t) \right] \right\|_{\mathcal{H}_k}^2$$

$$\min_{f \in \mathcal{F}} \max_{k \in \mathcal{K}} \frac{1}{n_a} \sum_{i=1}^{n_a} J(f(\mathbf{x}_i^a), y_i^a) + \lambda \sum_{\ell=l_1}^{l_2} d_k^2(\mathcal{D}_s^\ell, \mathcal{D}_t^\ell)$$

# Domain Adversarial Training (DANN)

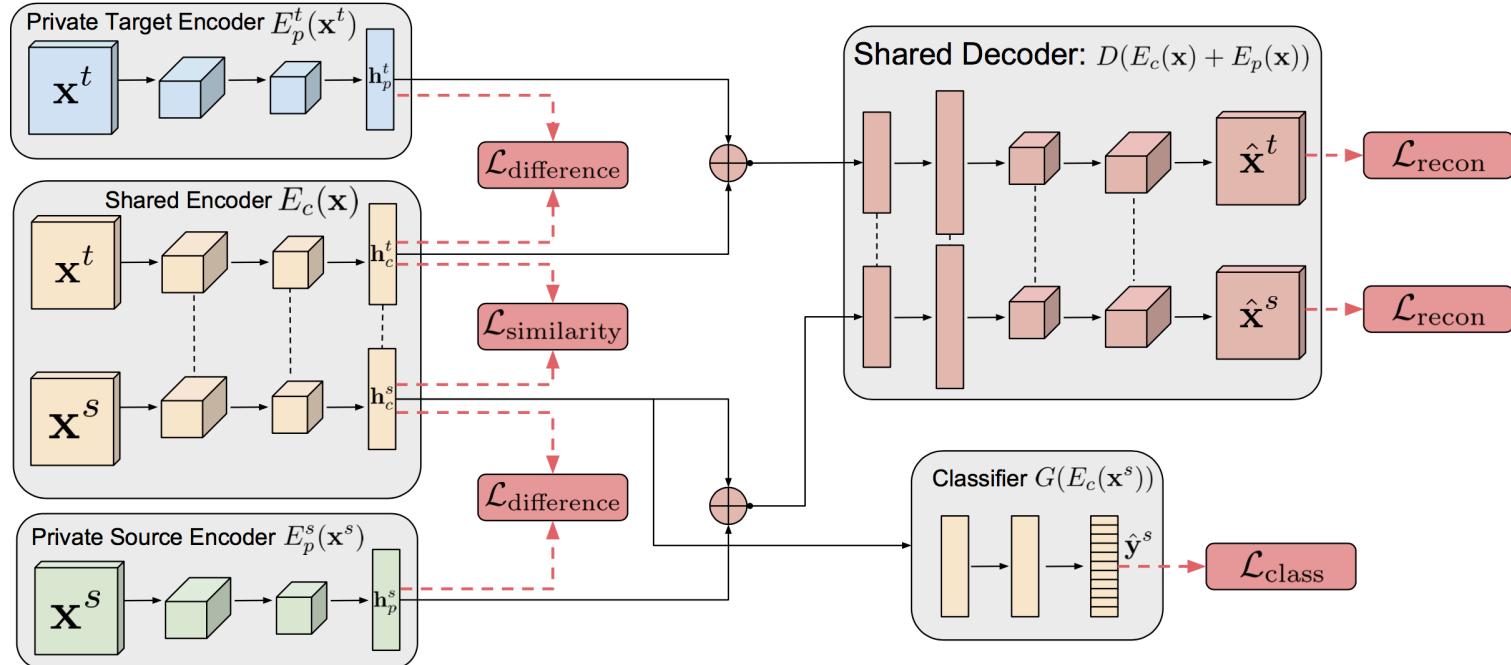


- **Adversarial adaptation:** learning features indistinguishable across domains

$$E(\theta_f, \theta_y, \theta_d) = \sum_{x_i \in \mathcal{D}_s} L_y(G_y(G_f(x_i)), y_i) - \lambda \sum_{x_i \in \mathcal{D}_s \cup \mathcal{D}_t} L_d(G_d(G_f(x_i)), d_i)$$

$$(\hat{\theta}_f, \hat{\theta}_y) = \operatorname{argmin}_{\theta_f, \theta_y} E(\theta_f, \theta_y, \theta_d) \quad (\hat{\theta}_d) = \operatorname{argmax}_{\theta_d} E(\theta_f, \theta_y, \theta_d)$$

# Domain Separation Network (DSN)

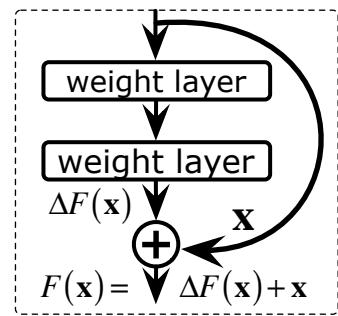
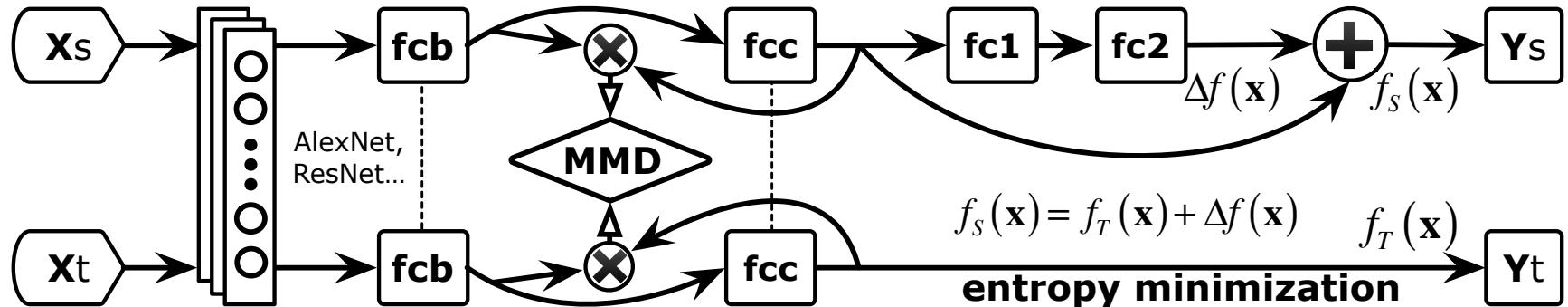


$$\hat{\mathbf{x}} = D(E_c(\mathbf{x}) + E_p(\mathbf{x})) \quad \hat{\mathbf{y}} = G(E_c(\mathbf{x}))$$

$$L = L_{\text{task}} + \alpha L_{\text{recon}} + \beta L_{\text{diff}} + \gamma L_{\text{sim}}$$

$$L_{\text{diff}} = \left\| \mathbf{H}_c^{s\top} \mathbf{H}_p^s \right\|_F^2 + \left\| \mathbf{H}_c^{t\top} \mathbf{H}_p^t \right\|_F^2$$

# Residual Transfer (RTN)

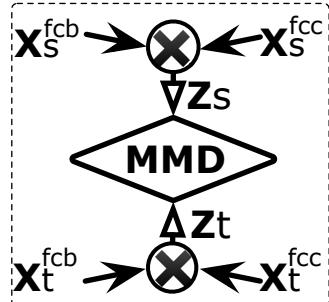


Classifier  
Adaptation

$$f_s = f_T + \Delta f \quad \min \frac{1}{n_s} \sum_{i=1}^{n_s} L(f_s(\mathbf{x}_i^s), y_i^s)$$

$$+ \frac{\gamma}{n_t} \sum_{i=1}^{n_t} H(f_t(\mathbf{x}_i^t))$$

$$+ \lambda D_{\mathcal{L}}(\mathcal{D}_s, \mathcal{D}_t),$$

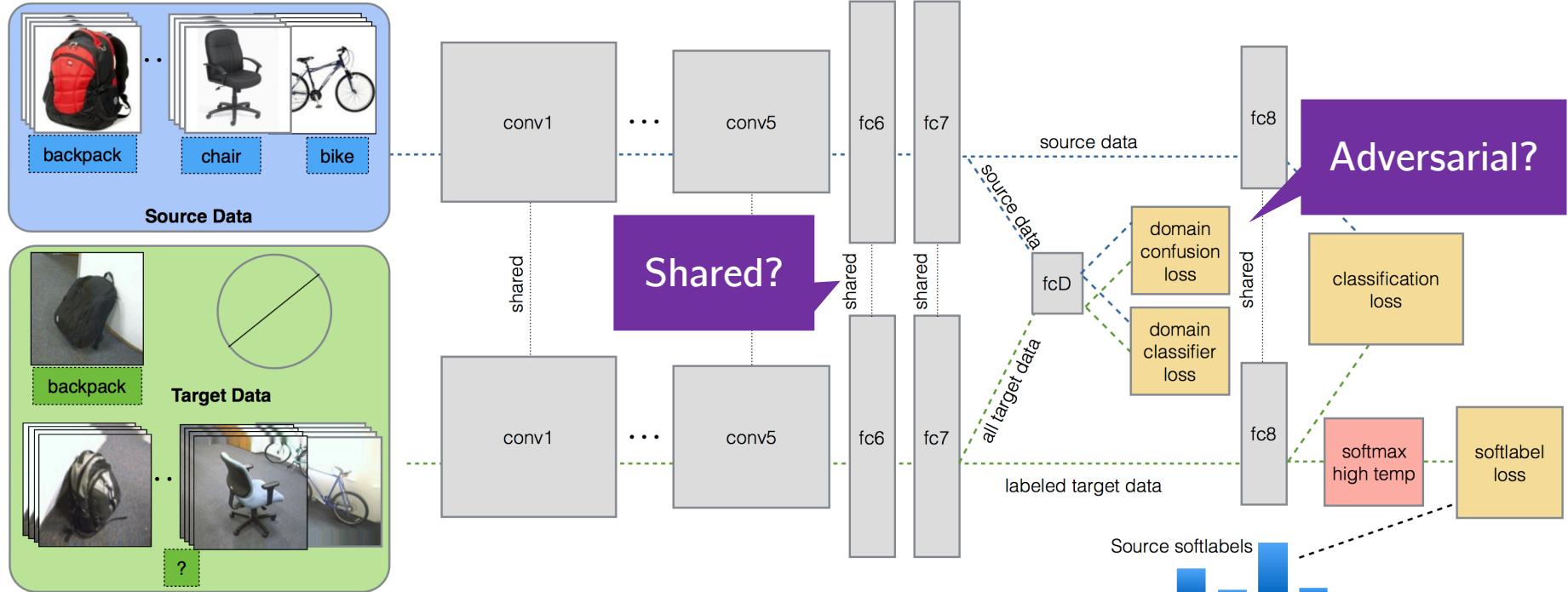


Feature  
Adaptation

# Asymmetric Transfer (ADDA)



清华大学  
Tsinghua University



$$\begin{aligned} & \min_D L_{adv_D}(\mathbf{X}_s, \mathbf{X}_t, M_s, M_t) \\ &= -\mathbf{E}_{\mathbf{x}_s} [\log D(M_s(\mathbf{x}_s))] \\ &\quad -\mathbf{E}_{\mathbf{x}_t} [\log(1 - D(M_t(\mathbf{x}_t)))] \end{aligned}$$

$$\begin{aligned} & \min_{M_s, M_t} L_{adv_M}(\mathbf{X}_s, \mathbf{X}_t, D) \\ &= -\mathbf{E}_{\mathbf{x}_t} [\log D(M_t(\mathbf{x}_t))] \end{aligned}$$

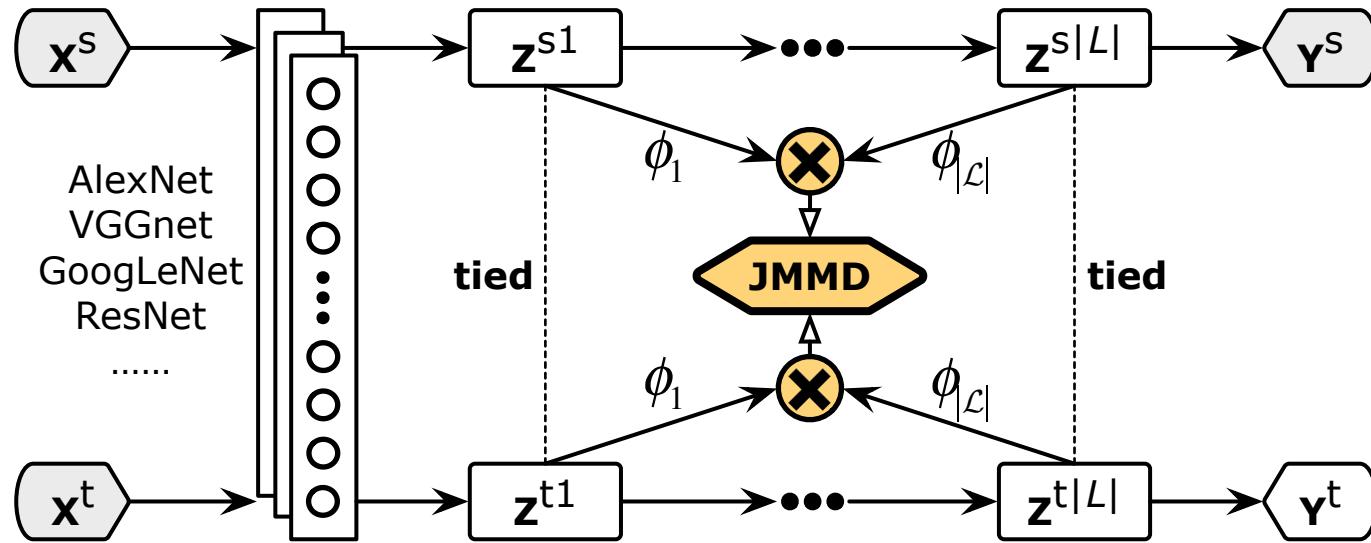
**Asymmetric**

# Problem 2



$$P(x,y) \neq Q(x,y)$$

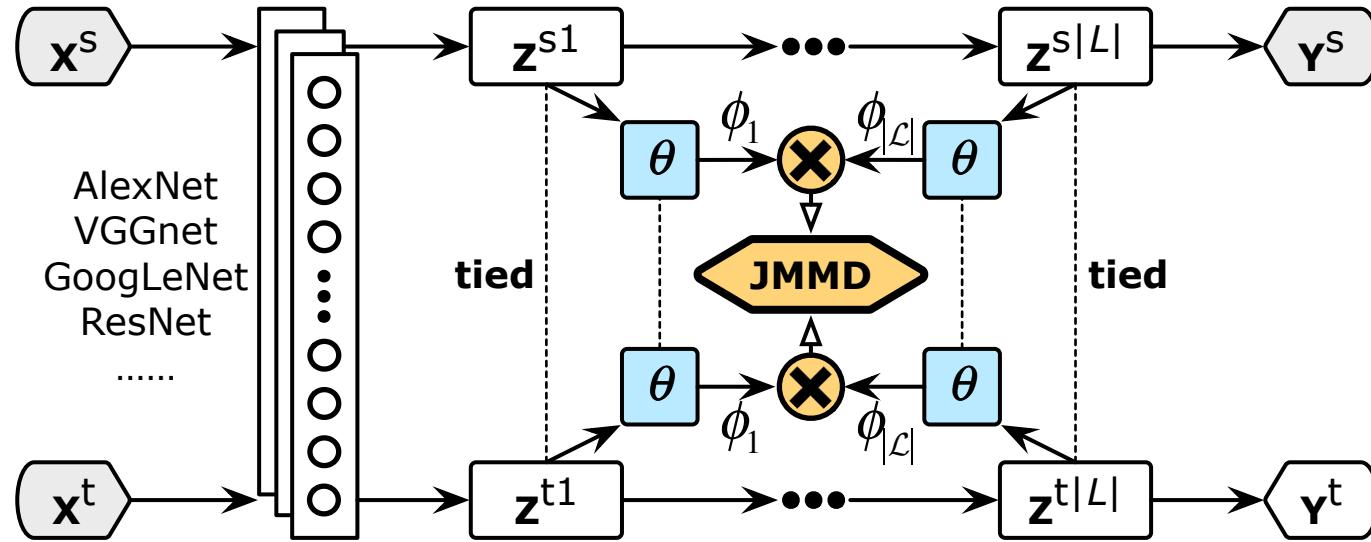
# Joint Adaptation Network (JAN)



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

$$D_{\mathcal{L}} \triangleq \left\| \mathbf{E}_P \left[ \bigotimes_{\ell=1}^{|\mathcal{L}|} \phi^\ell \left( \mathbf{z}^{s\ell} \right) \right] - \mathbf{E}_Q \left[ \bigotimes_{\ell=1}^{|\mathcal{L}|} \phi^\ell \left( \mathbf{z}^{t\ell} \right) \right] \right\|_{\bigotimes_{\ell=1}^{|\mathcal{L}|} \mathcal{H}^\ell}^2$$

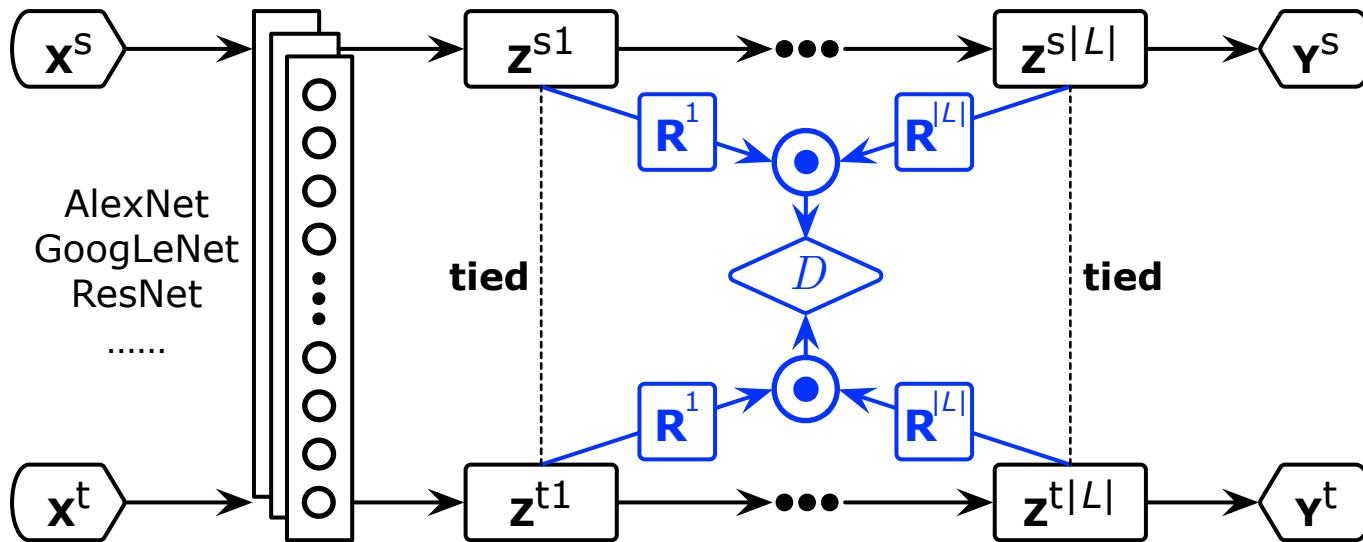
# Adversarial JAN



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

$$D_{\mathcal{L}} \triangleq \left\| \mathbf{E}_P \left[ \bigotimes_{\ell=1}^{|\mathcal{L}|} \phi^\ell \left( \theta^\ell \left( \mathbf{z}^{s\ell} \right) \right) \right] - \mathbf{E}_Q \left[ \bigotimes_{\ell=1}^{|\mathcal{L}|} \phi^\ell \left( \theta^\ell \left( \mathbf{z}^{t\ell} \right) \right) \right] \right\|_{\bigotimes_{\ell=1}^{|\mathcal{L}|} \mathcal{H}^\ell}^2$$

# Multilinear Adversarial Network (MAN)



$$\phi_{\mathcal{L}}(\mathbf{z}_i^s) = \frac{1}{\sqrt{d}} \left( \odot_{\ell=1}^{|L|} \mathbf{R}^\ell \mathbf{z}_i^{s\ell} \right), \phi_{\mathcal{L}}(\mathbf{z}_j^t) = \frac{1}{\sqrt{d}} \left( \odot_{\ell=1}^{|L|} \mathbf{R}^\ell \mathbf{z}_j^{t\ell} \right)$$

$$\min_F \frac{1}{n_s} \sum_{i=1}^{n_s} J(F(\mathbf{x}_i^s), \mathbf{y}_i^s) + \frac{\lambda}{n_s} \sum_{i=1}^{n_s} \log D(\phi_{\mathcal{L}}(\mathbf{z}_i^s)) + \frac{\lambda}{n_t} \sum_{j=1}^{n_t} \log (1 - D(\phi_{\mathcal{L}}(\mathbf{z}_j^t)))$$

$$\min_D -\frac{1}{n_s} \sum_{i=1}^{n_s} \log D(\phi_{\mathcal{L}}(\mathbf{z}_i^s)) - \frac{1}{n_t} \sum_{j=1}^{n_t} \log (1 - D(\phi_{\mathcal{L}}(\mathbf{z}_j^t)))$$

# Empirical Benchmark



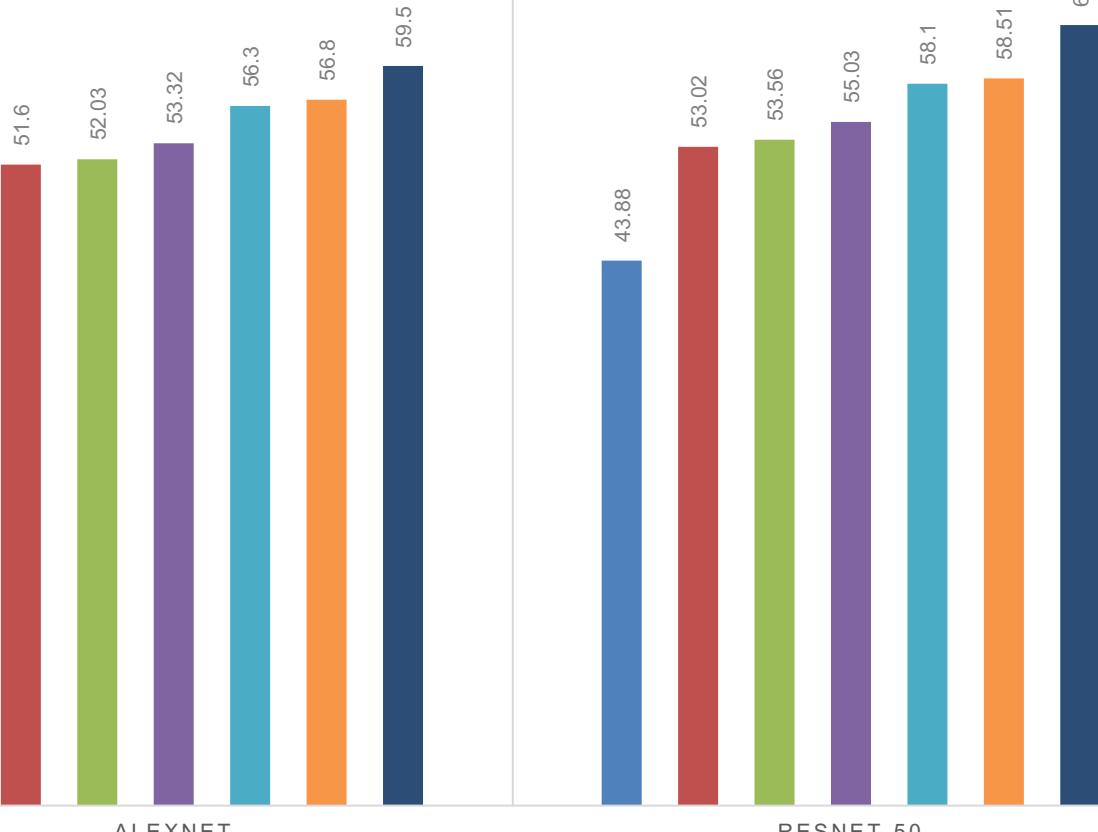
## VISDA CHALLENGE 2017

### Source Domain



### Target Domain

CNN DAN RTN RevGrad JAN JAN-A MAN



# Problem 3

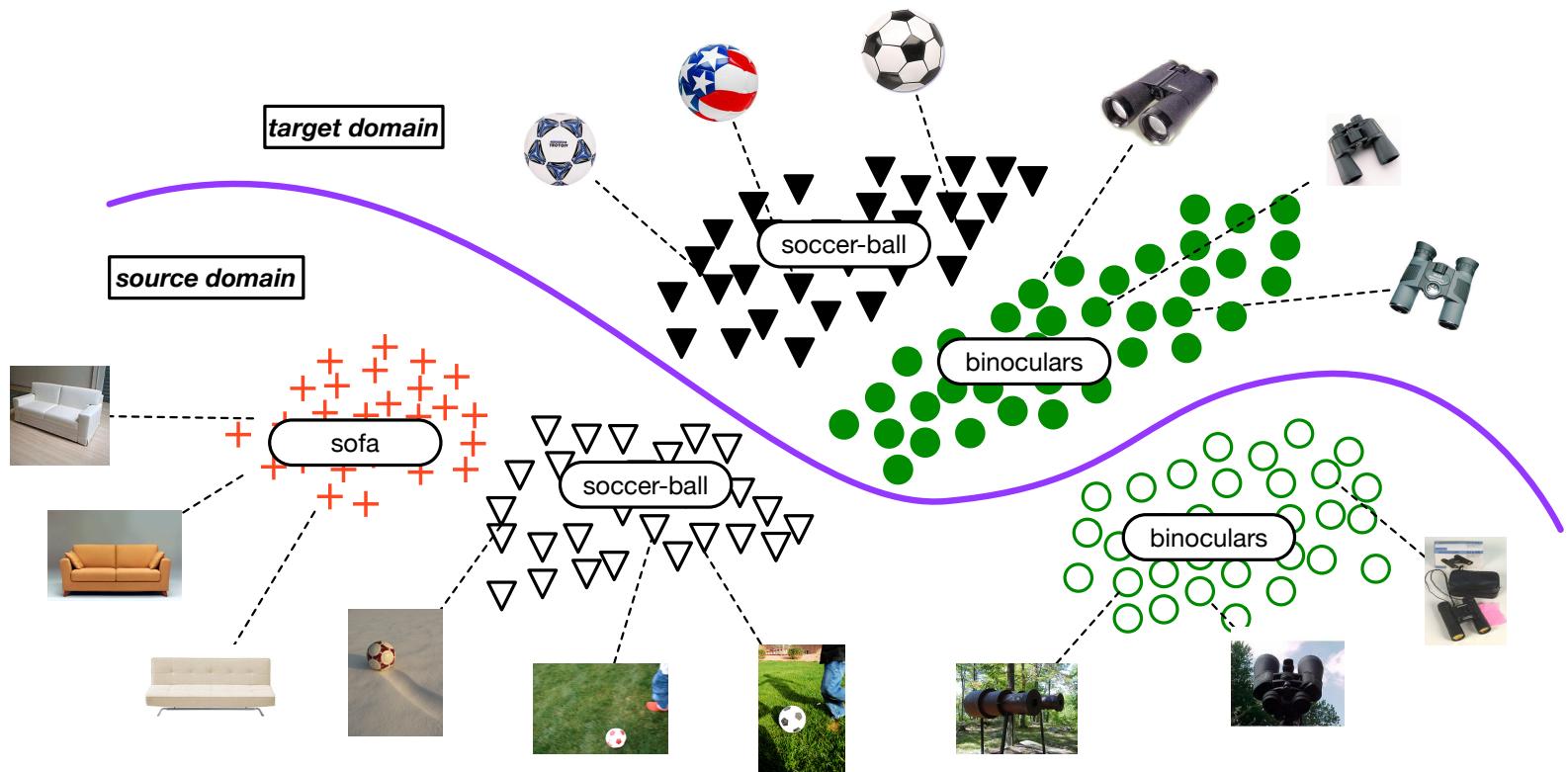


$$Y_s \neq Y_t$$

# Partial Transfer Learning

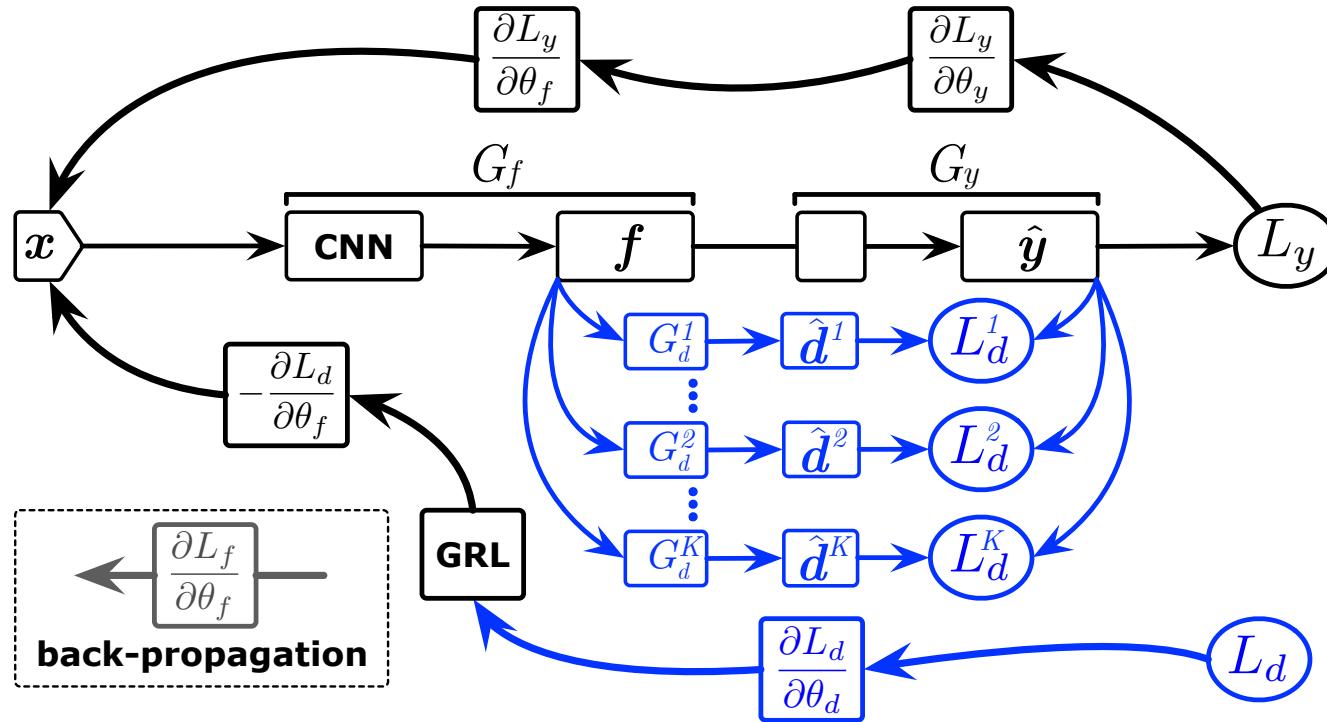


清华大学  
Tsinghua University



$$Y_s \supset Y_t$$

# Selective Adversarial Network (SAN)

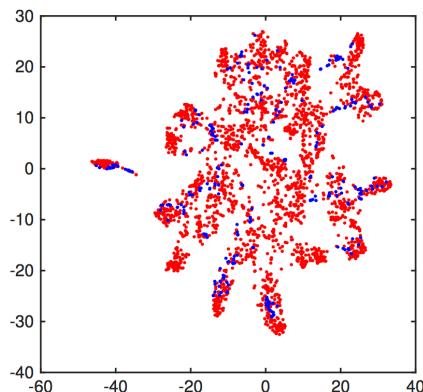


$$C(\theta_f, \theta_y, \theta_d^k |_{k=1}^{|\mathcal{C}_s|}) = \frac{1}{n_s} \sum_{\mathbf{x}_i \in \mathcal{D}_s} L_y(G_y(G_f(\mathbf{x}_i)), y_i) + \frac{1}{n_t} \sum_{\mathbf{x}_i \in \mathcal{D}_t} H(G_y(G_f(\mathbf{x}_i))) \\ - \frac{\lambda}{n_s + n_t} \sum_{k=1}^{|\mathcal{C}_s|} \left[ \frac{1}{n_t} \sum_{\mathbf{x}_i \in \mathcal{D}_t} \hat{y}_i^k \right] \sum_{\mathbf{x}_i \in \mathcal{D}_s \cup \mathcal{D}_t} \hat{y}_i^k L_d^k(G_d^k(G_f(\mathbf{x}_i)), d_i)$$

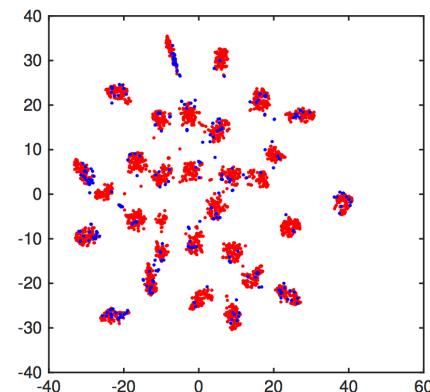
# Selective Adversarial Network (SAN)



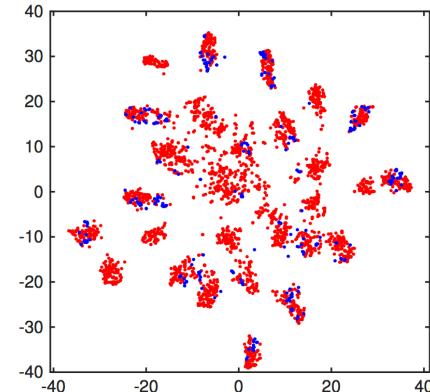
Method	Caltech-Office				ImageNet-Caltech		
	C 256 → W 10	C 256 → A 10	C 256 → D 10	Avg	I 1000 → C 84	C 256 → I 84	Avg
AlexNet [14]	58.44	76.64	65.86	66.98	52.37	47.35	49.86
DAN [15]	42.37	70.75	47.04	53.39	54.21	52.03	53.12
RevGrad [6]	54.57	72.86	57.96	61.80	51.34	47.02	49.18
RTN [17]	71.02	81.32	62.35	71.56	63.69	50.45	57.07
ADDA [26]	73.66	78.35	74.80	75.60	64.20	51.55	57.88
SAN-selective	76.44	81.63	80.25	79.44	66.78	51.25	59.02
SAN-entropy	72.54	78.95	76.43	75.97	55.27	52.31	53.79
SAN	<b>88.33</b>	<b>83.82</b>	<b>85.35</b>	<b>85.83</b>	<b>68.45</b>	<b>55.61</b>	<b>62.03</b>



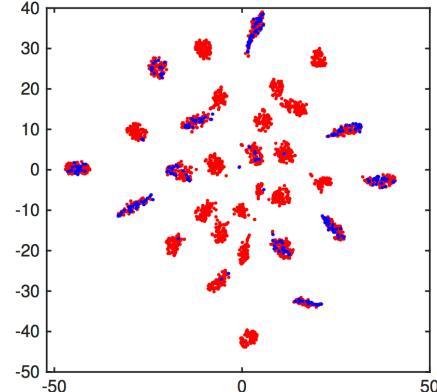
(a) DAN



(b) RevGrad



(c) RTN

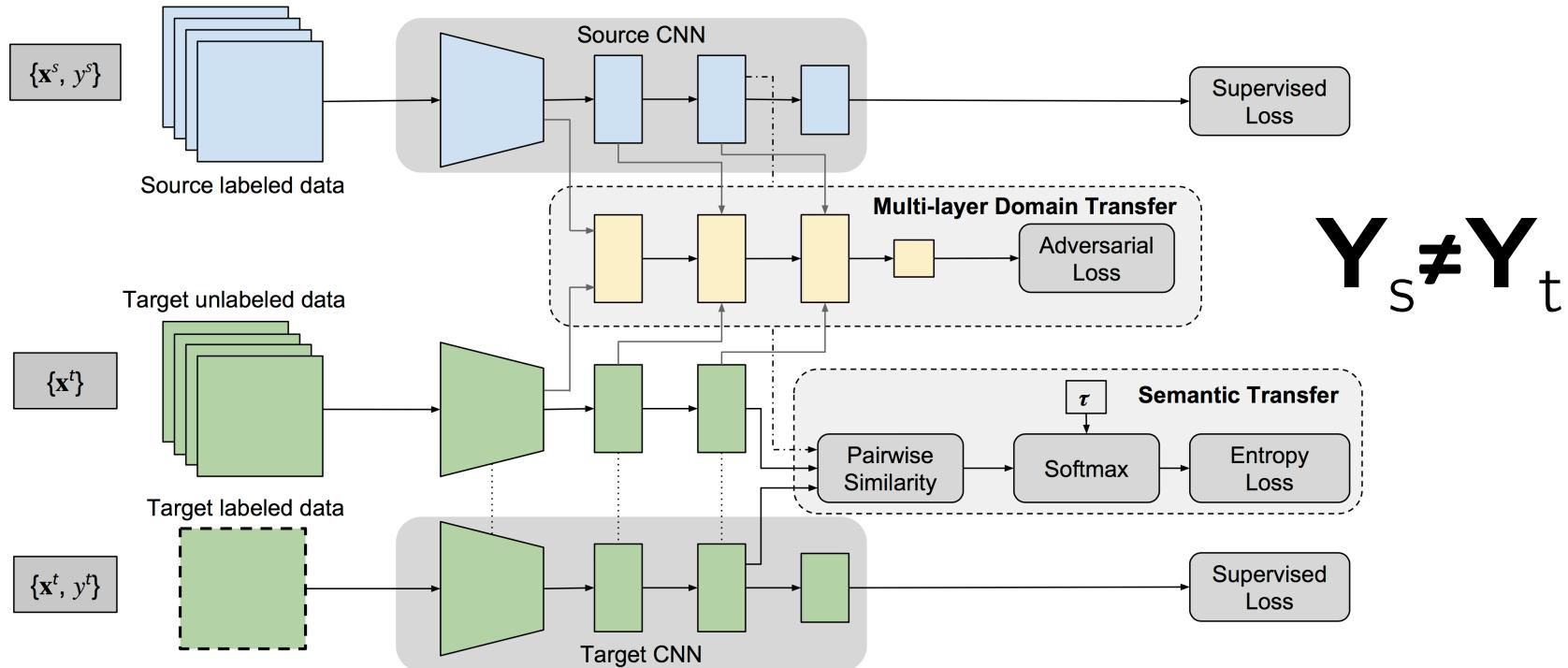


(d) SAN

# Joint Domain and Semantic Transfer



清华大学  
Tsinghua University



$$L_{ST}(\tilde{\mathbf{x}}_t, \mathbf{x}_s) = \sum_{\tilde{\mathbf{x}}_t \in \tilde{\mathbf{X}}_t} H(\sigma(v_s(\tilde{\mathbf{x}}_t) / \tau))$$

$$L_{ST,upsup}(\tilde{\mathbf{x}}_t, \mathbf{x}_t) = \sum_{\tilde{\mathbf{x}}_t \in \tilde{\mathbf{X}}_t} H(\sigma(v_t(\tilde{\mathbf{x}}_t) / \tau))$$

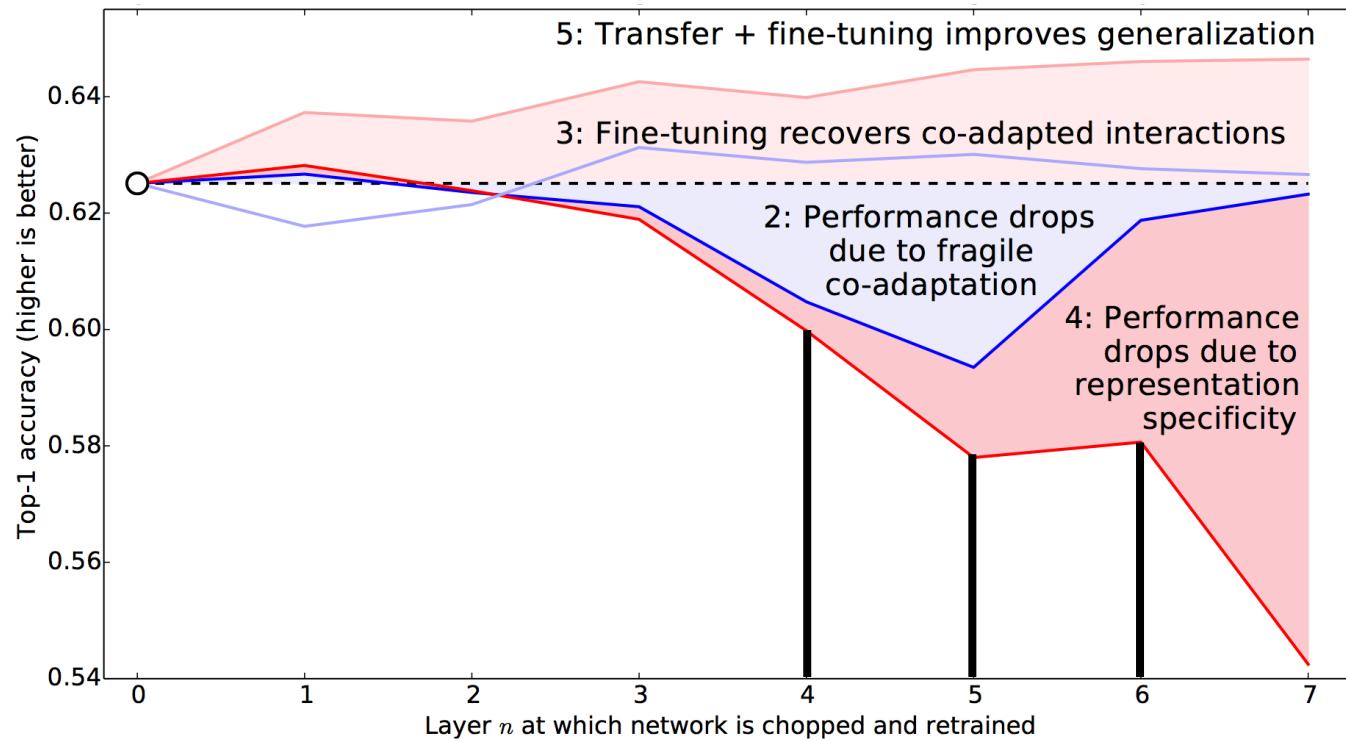
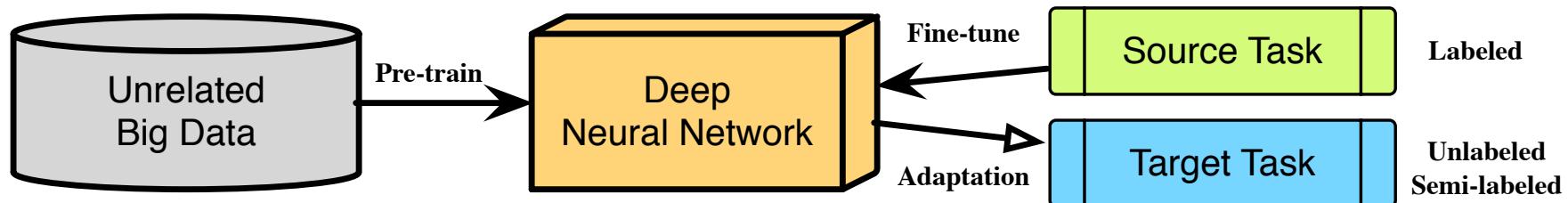
$$L_{ST,sup}(\mathbf{x}_t) = - \sum_{\{\mathbf{x}_s, \mathbf{x}_t\} \in \mathbf{X}_t} \log \frac{\exp([v_t(\mathbf{x}_t)]_{\mathbf{y}_t})}{\sum_{i=1}^n \exp([v_t(\mathbf{x}_t)]_i)}$$

# Problem 4



# Transferable Architecture

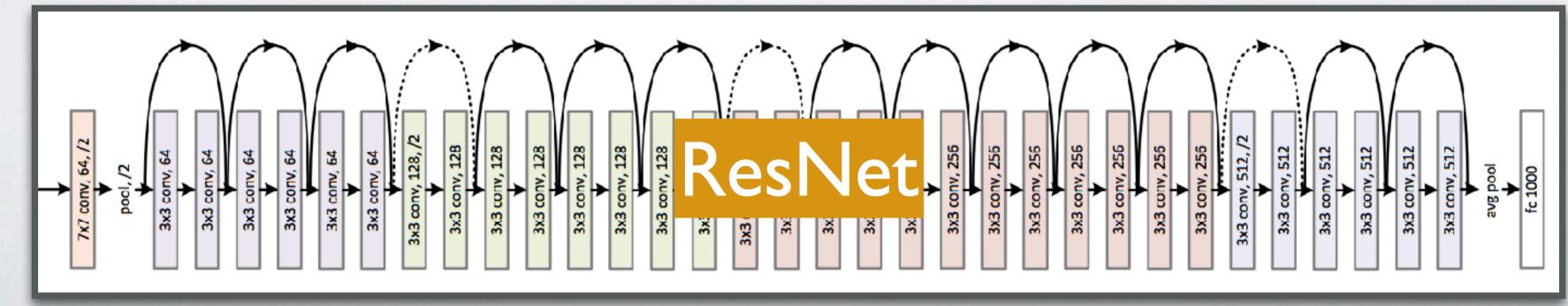
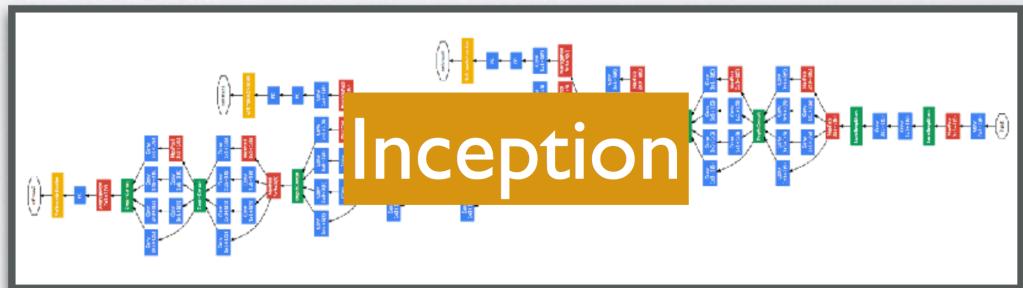
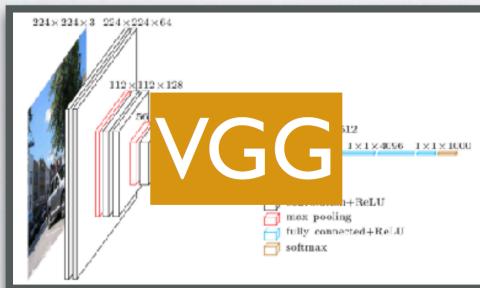
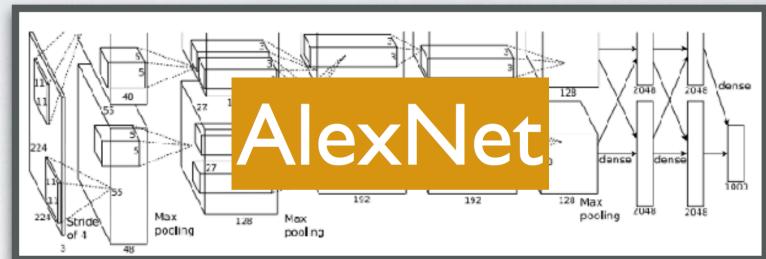
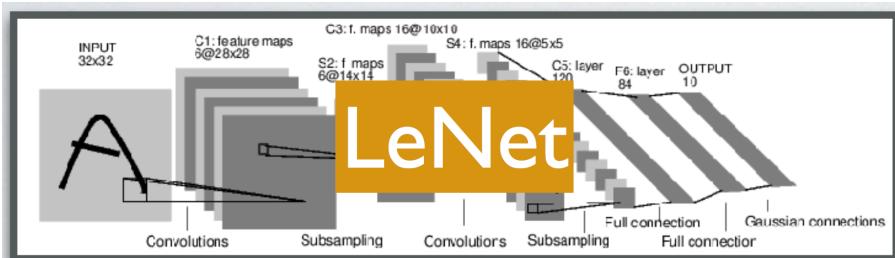
# Transferability



# Transferable Architecture



清华大学  
Tsinghua University



Some modules may not influence in-domain accuracy but influence the transferability

# Open Problems



- Heterogeneous Transfer Learning

$$X_s \neq X_t \wedge Y_s \neq Y_t$$

- Pixel-Level Transfer Learning

$$P(x) \neq Q(x) \wedge P(z) \neq Q(z)$$

- Learning Transferable Architectures

# Thank You!