



清華大學
Tsinghua University

Recent Advances in Transfer Learning

Mingsheng Long
Tsinghua University
Jan 25, 2018

<https://github.com/thuml>

Machine Learning

Learner: $f : \mathbf{x} \rightarrow y$

Distribution: $(\mathbf{x}, y) \sim P(\mathbf{x}, y)$



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

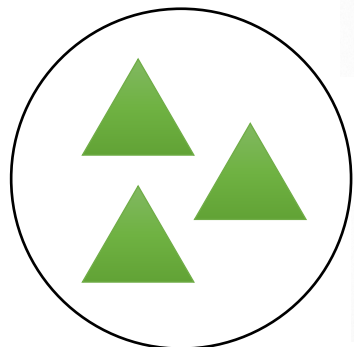
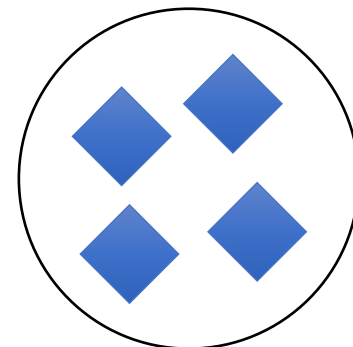


2D Renderings



Real Images

Target Domain



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



Model

$$f : \mathbf{x} \rightarrow y$$



Representation

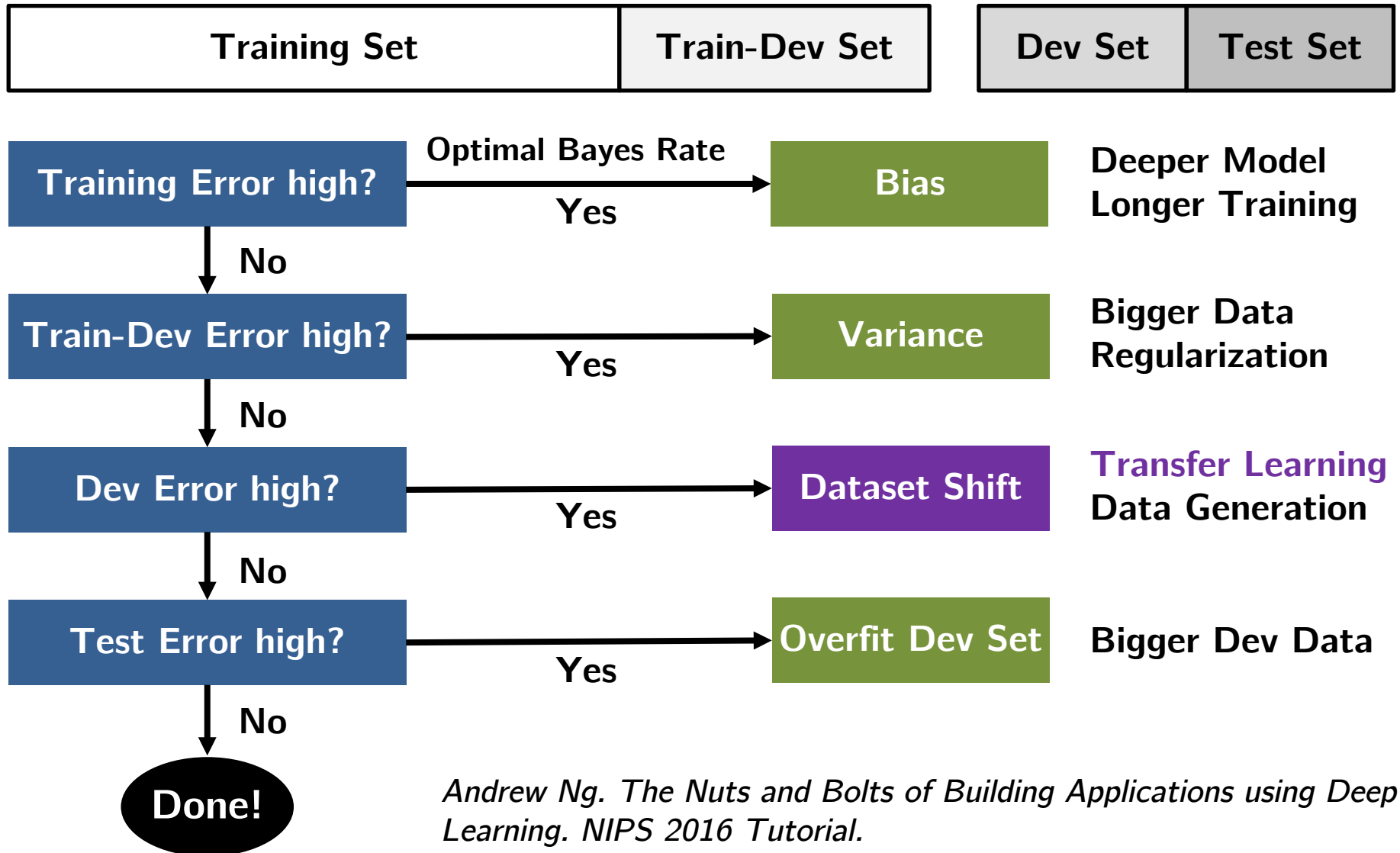
<http://ai.bu.edu/visda-2017/>



Model

$$f : \mathbf{x} \rightarrow y$$

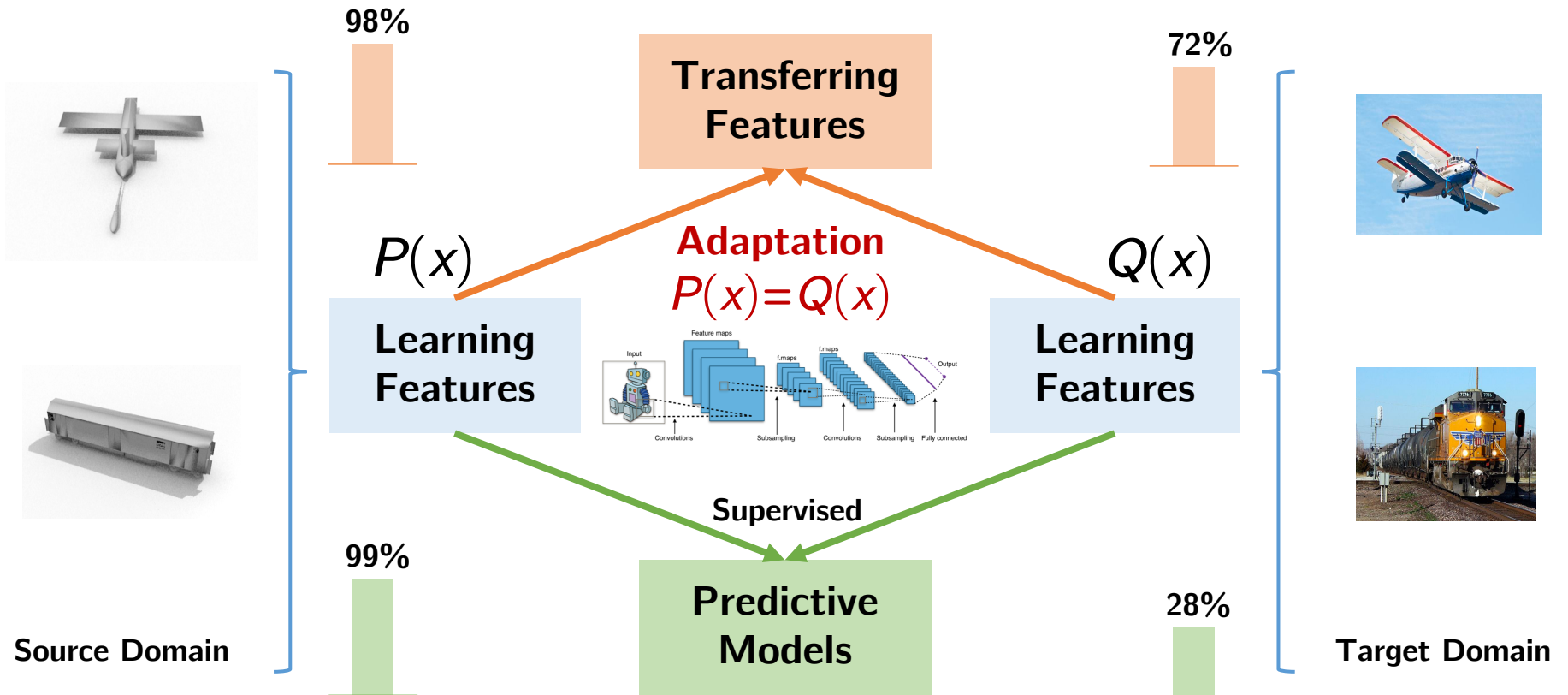
Transfer Learning: Why?



Andrew Ng. *The Nuts and Bolts of Building Applications using Deep Learning*. NIPS 2016 Tutorial.

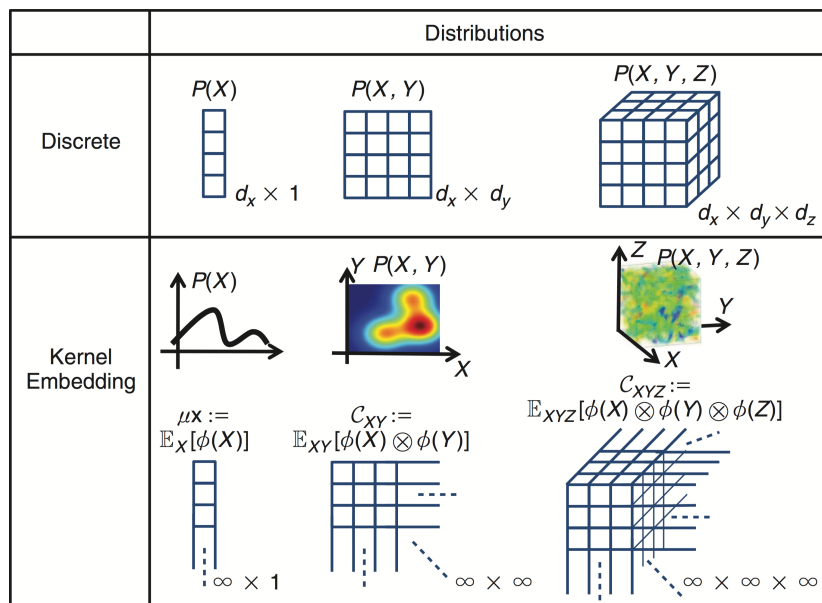
Transfer Learning: How?

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

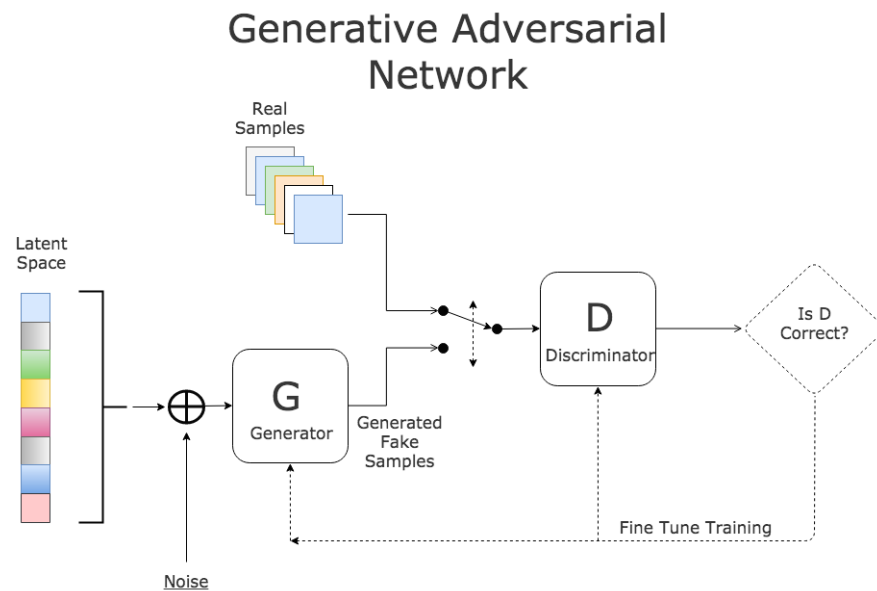


Distribution Matching

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



Kernel Embedding



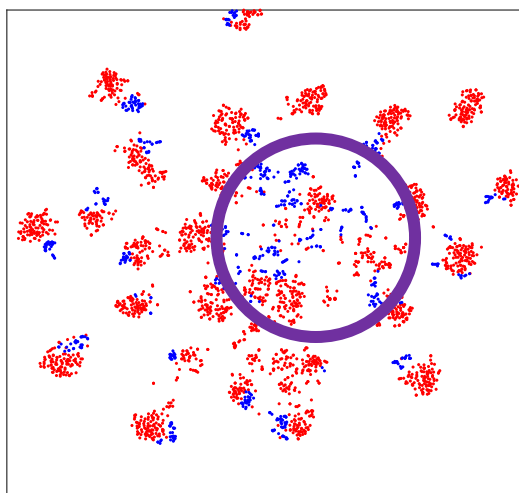
Adversarial Learning

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

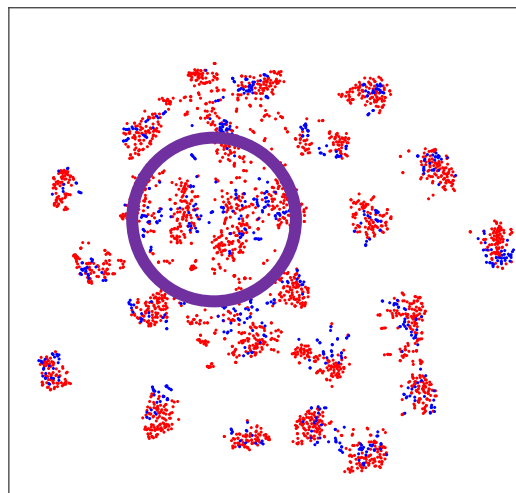
Goodfellow et al. Generative Adversarial Networks. *NIPS* 2014.

Distribution Matching

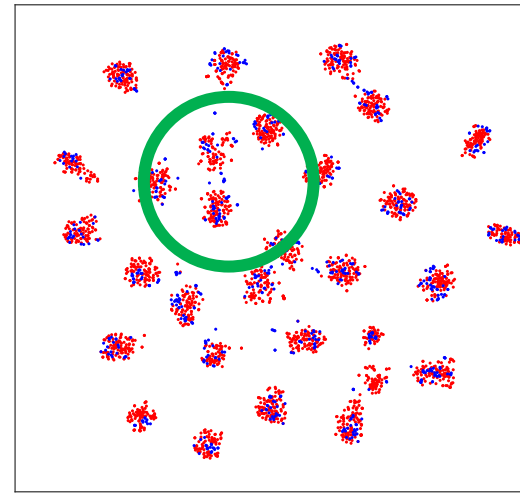
- Marginal distribution mismatch: $P(\mathbf{x}) \neq Q(\mathbf{x})$
- Conditional distribution mismatch: $P(\mathbf{y}|\mathbf{x}) \neq Q(\mathbf{y}|\mathbf{x})$



$$P(\mathbf{x}) \neq Q(\mathbf{x})$$
$$P(\mathbf{y}|\mathbf{x}) \neq Q(\mathbf{y}|\mathbf{x})$$



$$P(\mathbf{x}) \approx Q(\mathbf{x})$$
$$P(\mathbf{y}|\mathbf{x}) \neq Q(\mathbf{y}|\mathbf{x})$$



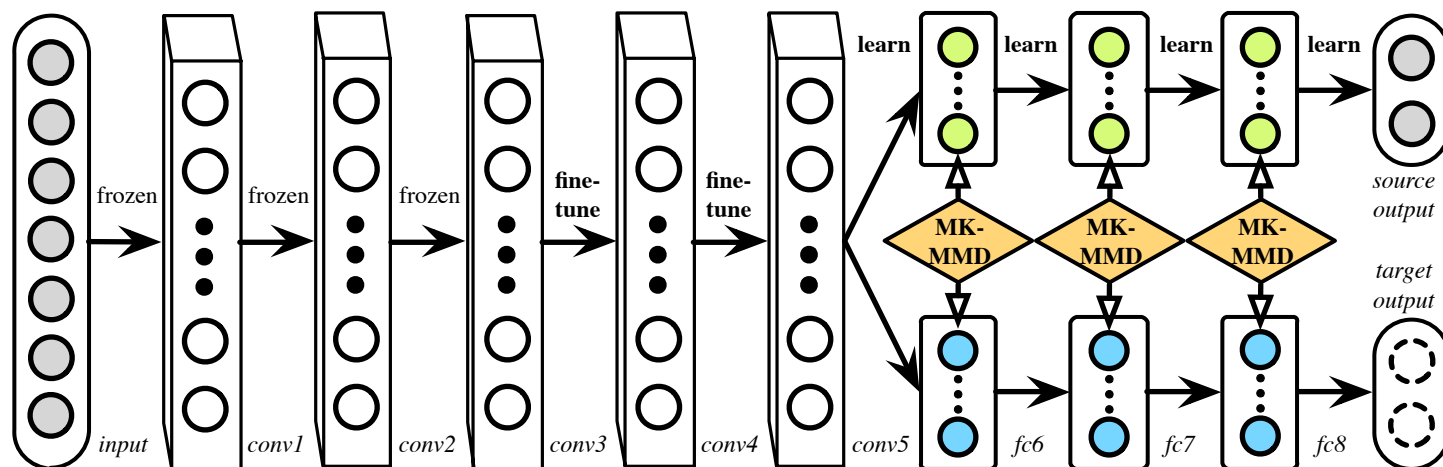
$$P(\mathbf{x}, \mathbf{y}) \approx Q(\mathbf{x}, \mathbf{y})$$
$$P(\mathbf{y}|\mathbf{x}) \neq Q(\mathbf{y}|\mathbf{x})$$

Problem 1



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

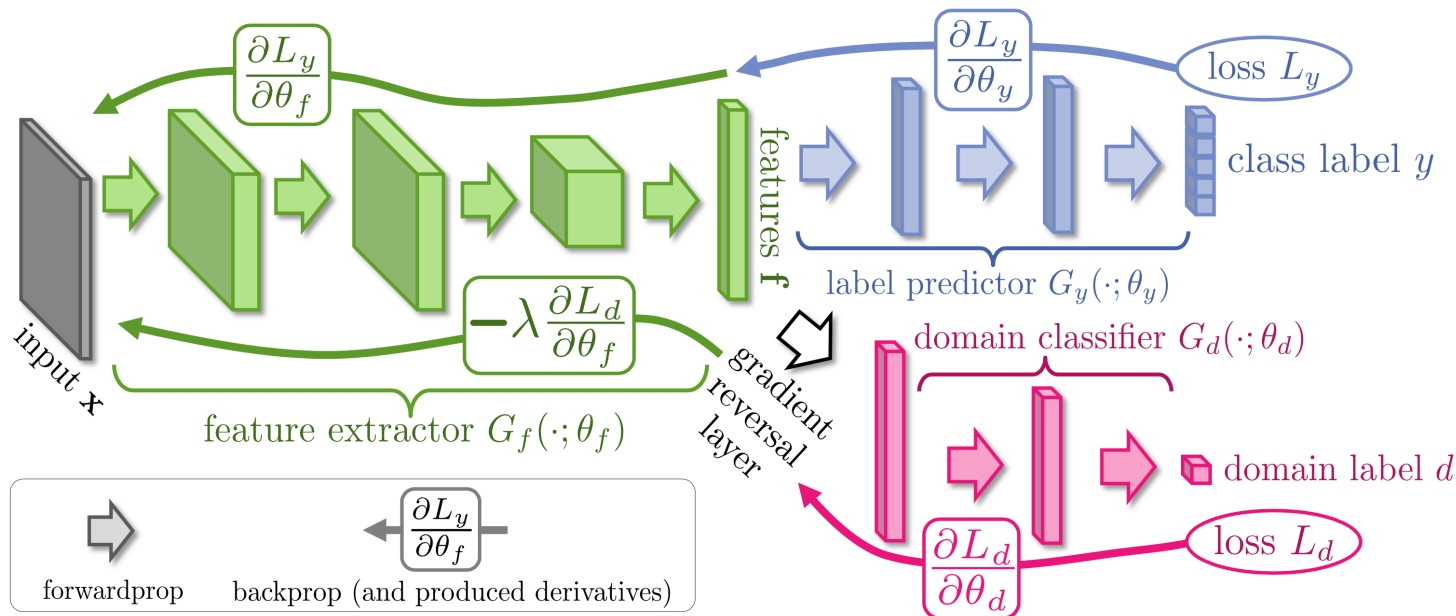
Deep Adaptation Network (DAN)



- 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)

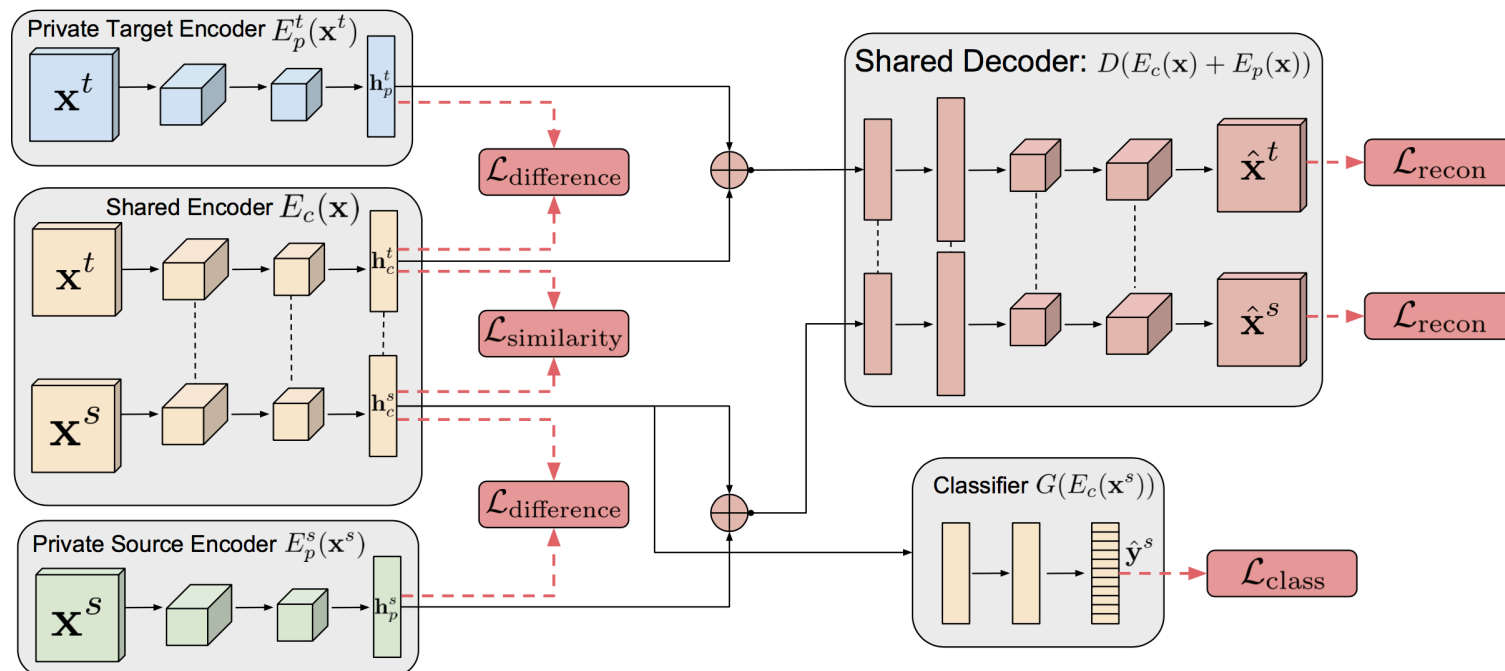


- **Adversarial** adaptation: learning features indistinguishable across domains

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

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

Domain Separation Network (DSN)

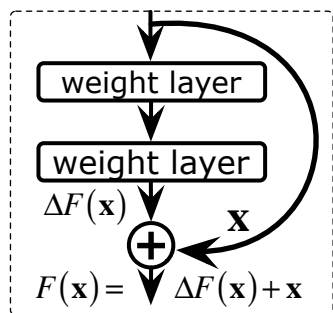
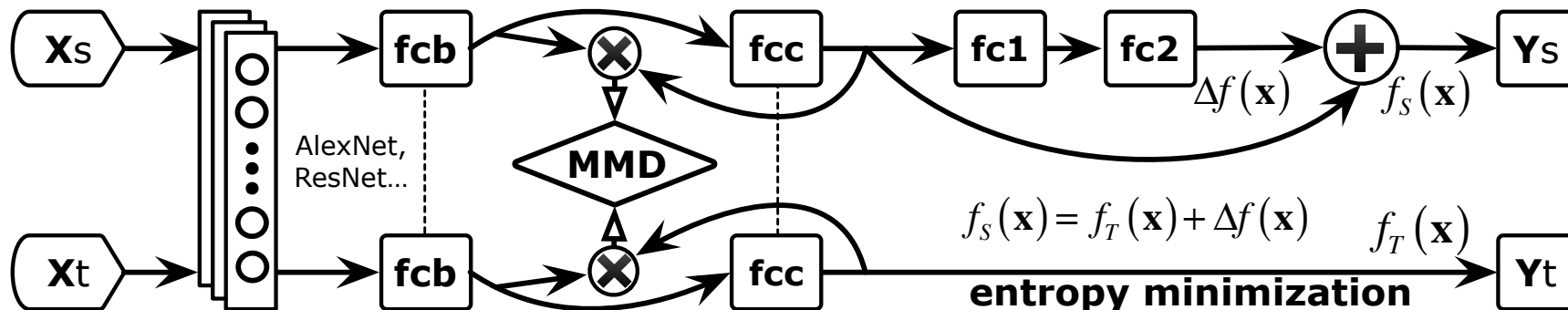


$$\hat{\mathbf{x}} = D(E_c(\mathbf{x}) + E_p(\mathbf{x})) \quad \hat{\mathbf{y}} = G(E_c(\mathbf{x}))$$

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

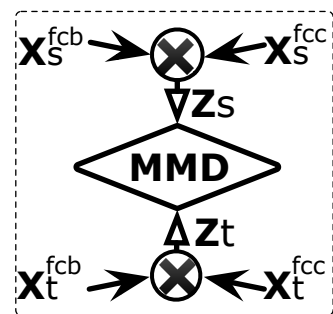
$$L_{diff} = \left\| \mathbf{H}_c^{sT} \mathbf{H}_p^s \right\|_F^2 + \left\| \mathbf{H}_c^{tT} \mathbf{H}_p^t \right\|_F^2$$

Residual Transfer (RTN)



Classifier
Adaptation

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

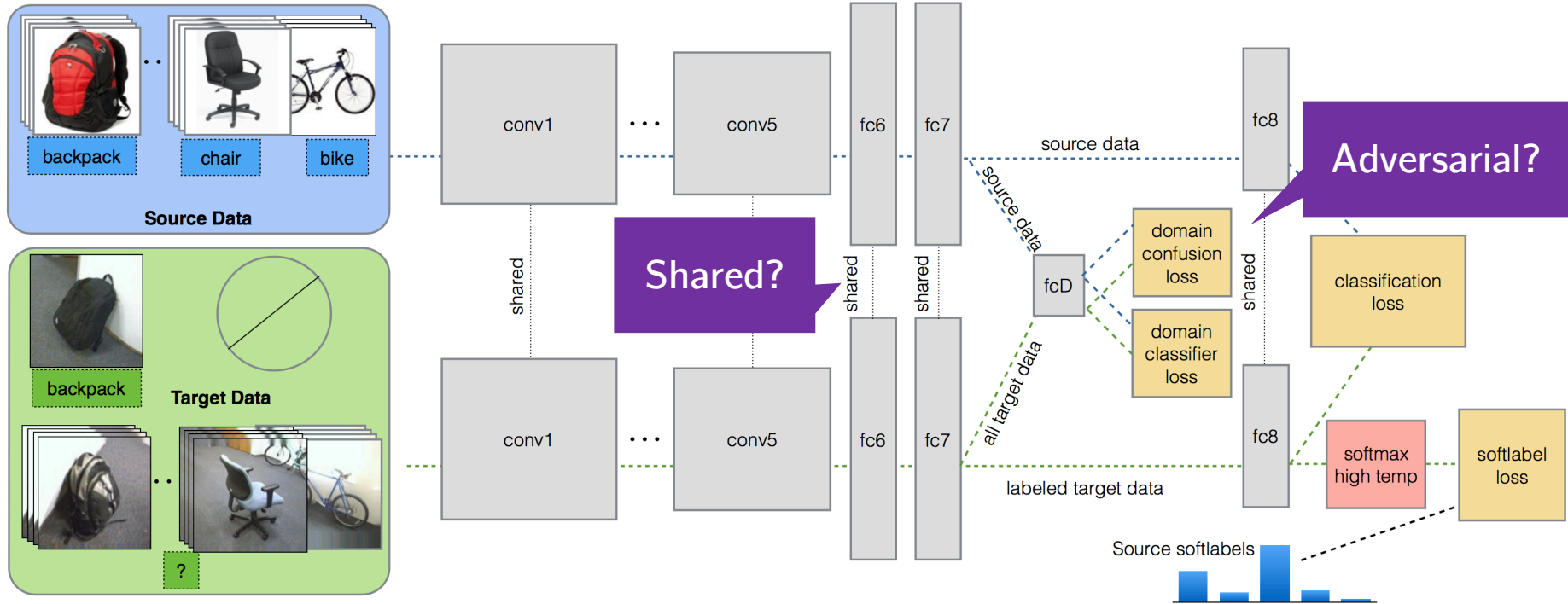


Feature
Adaptation

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

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

Asymmetric Transfer (ADDA)



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

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

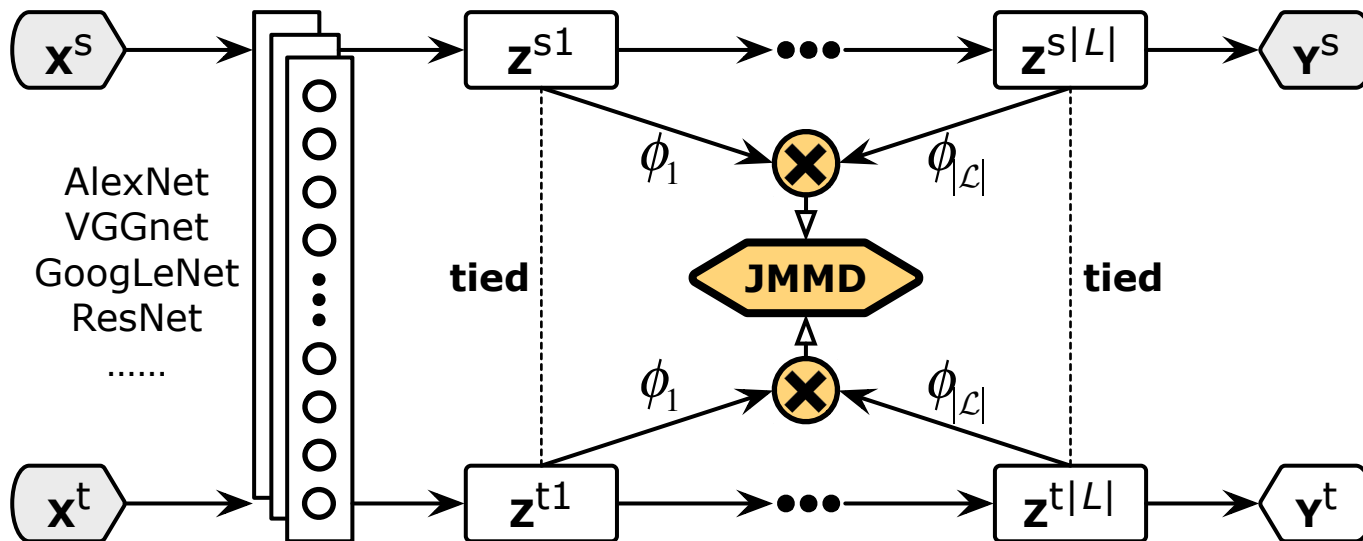
Asymmetric

Problem 2



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

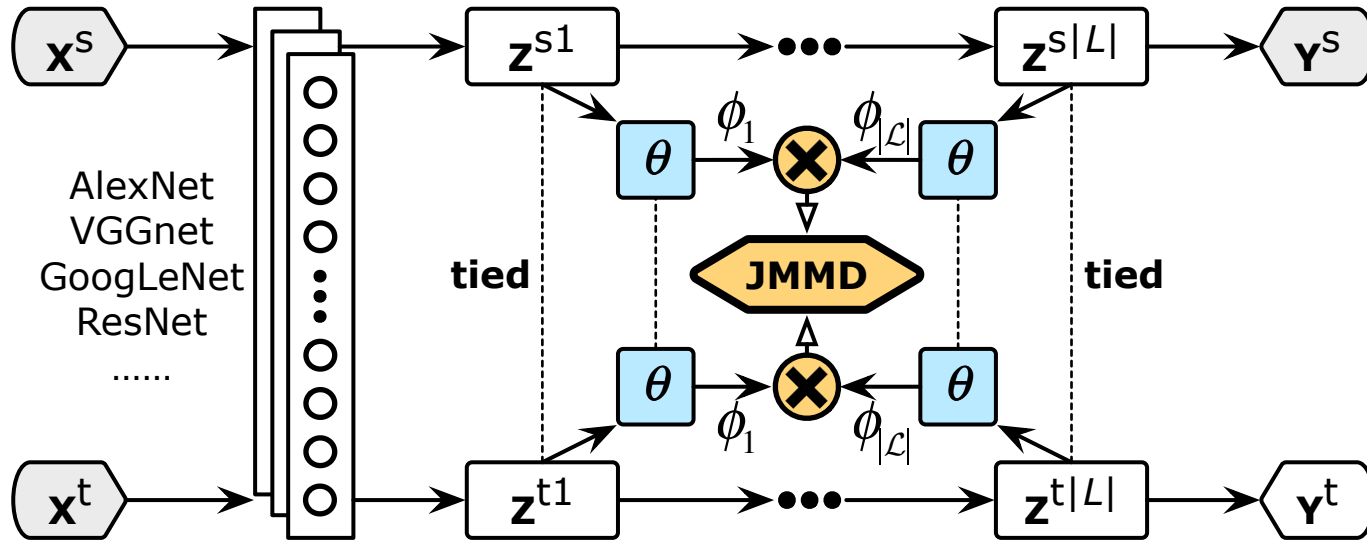
Joint Adaptation Network (JAN)



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

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

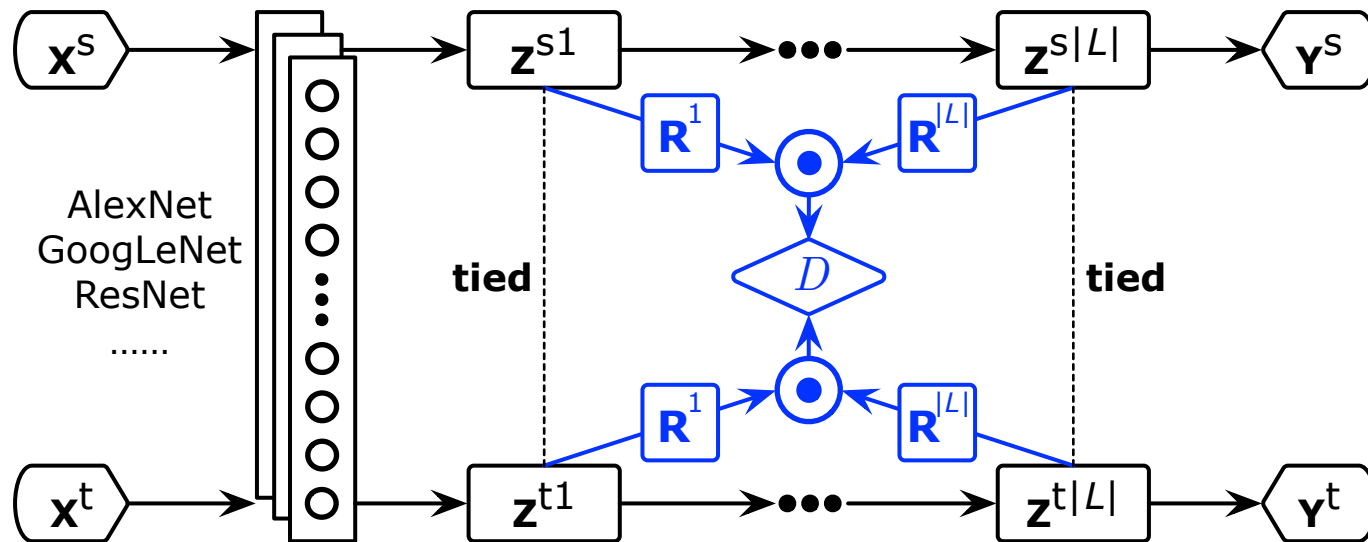
Adversarial JAN



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

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

Multilinear Adversarial Network (MAN)



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

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

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

Empirical Benchmark

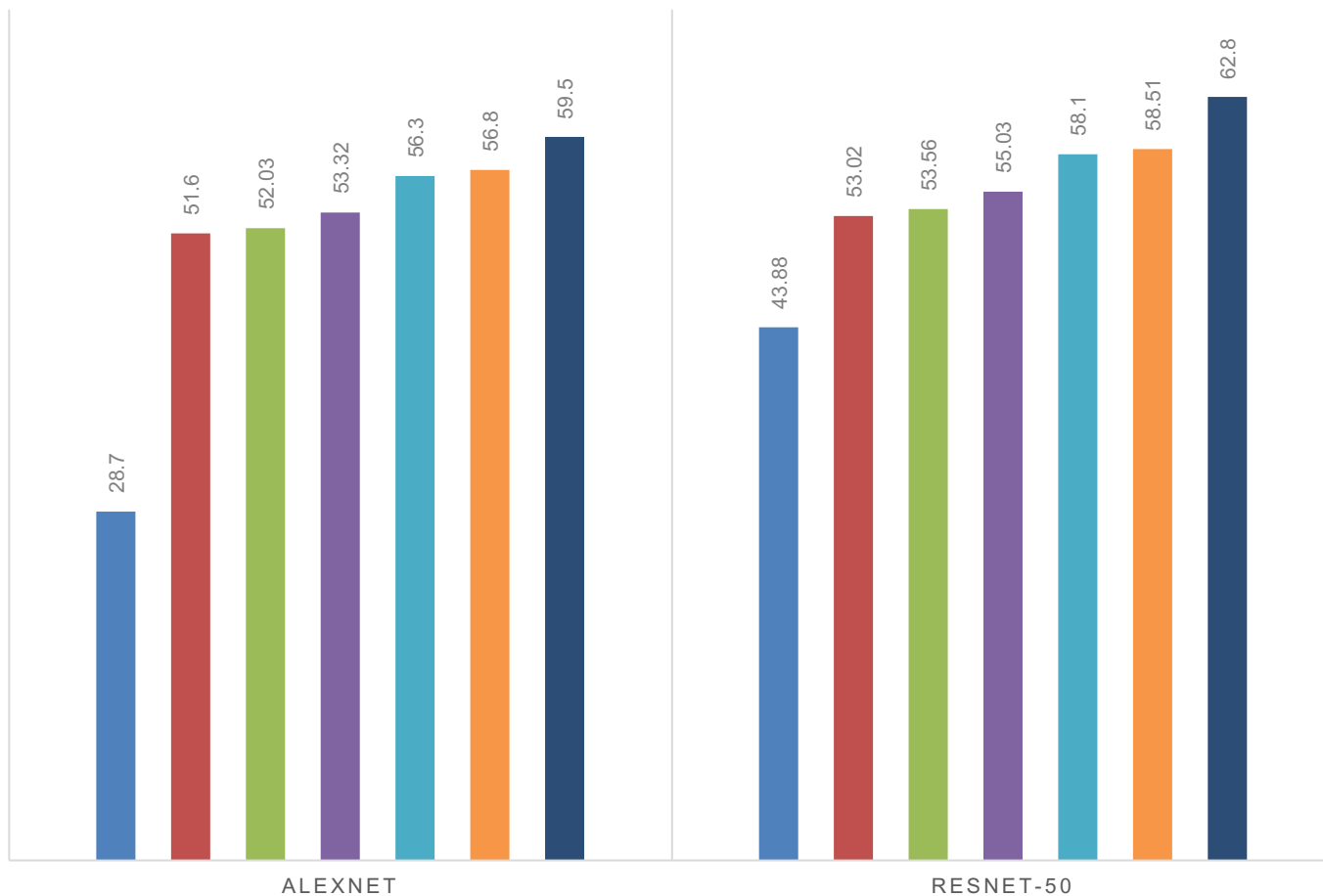
VISDA CHALLENGE 2017

Source Domain



Target Domain

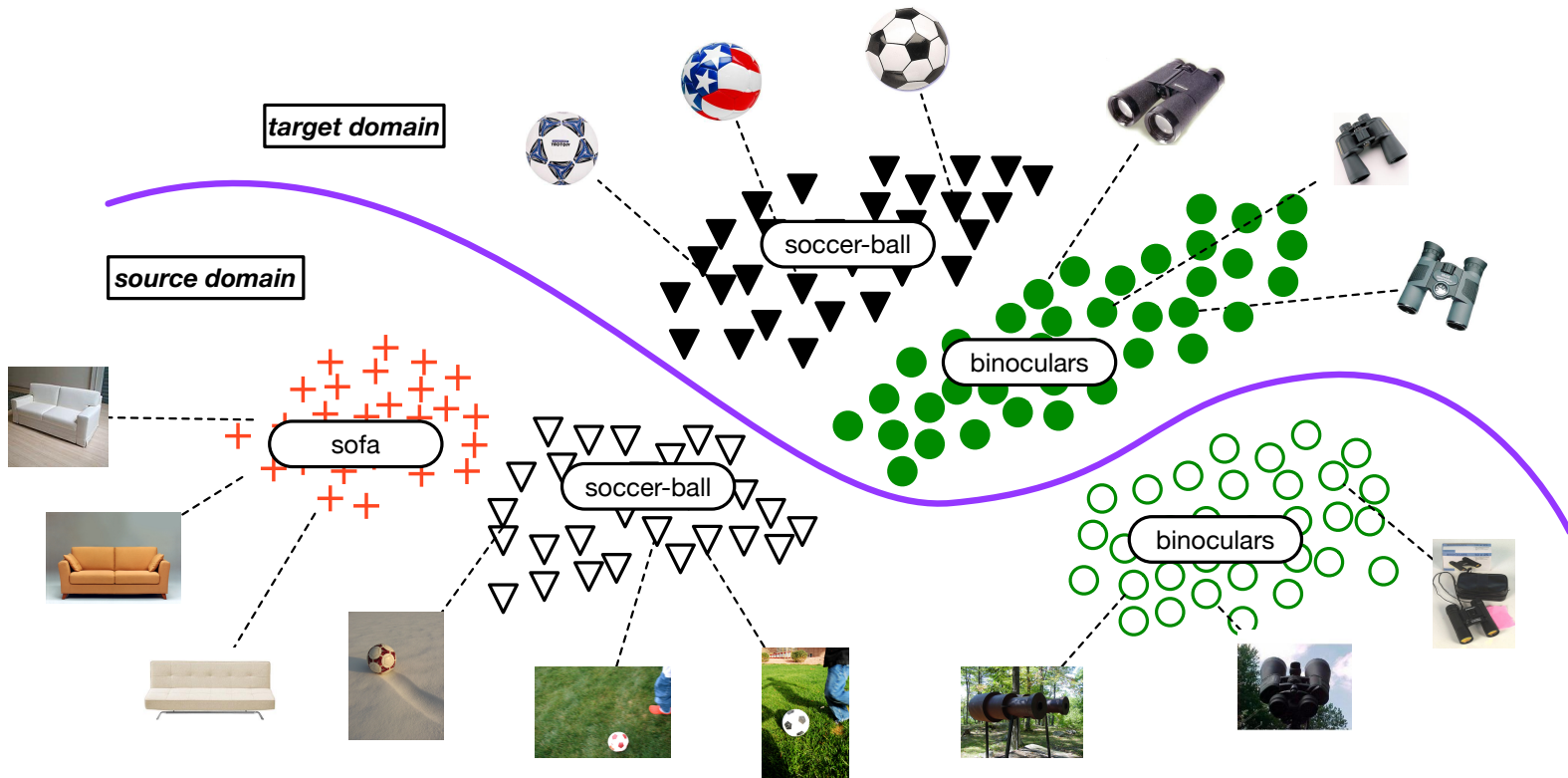
■ CNN ■ DAN ■ RTN ■ RevGrad ■ JAN ■ JAN-A ■ MAN



Problem 3

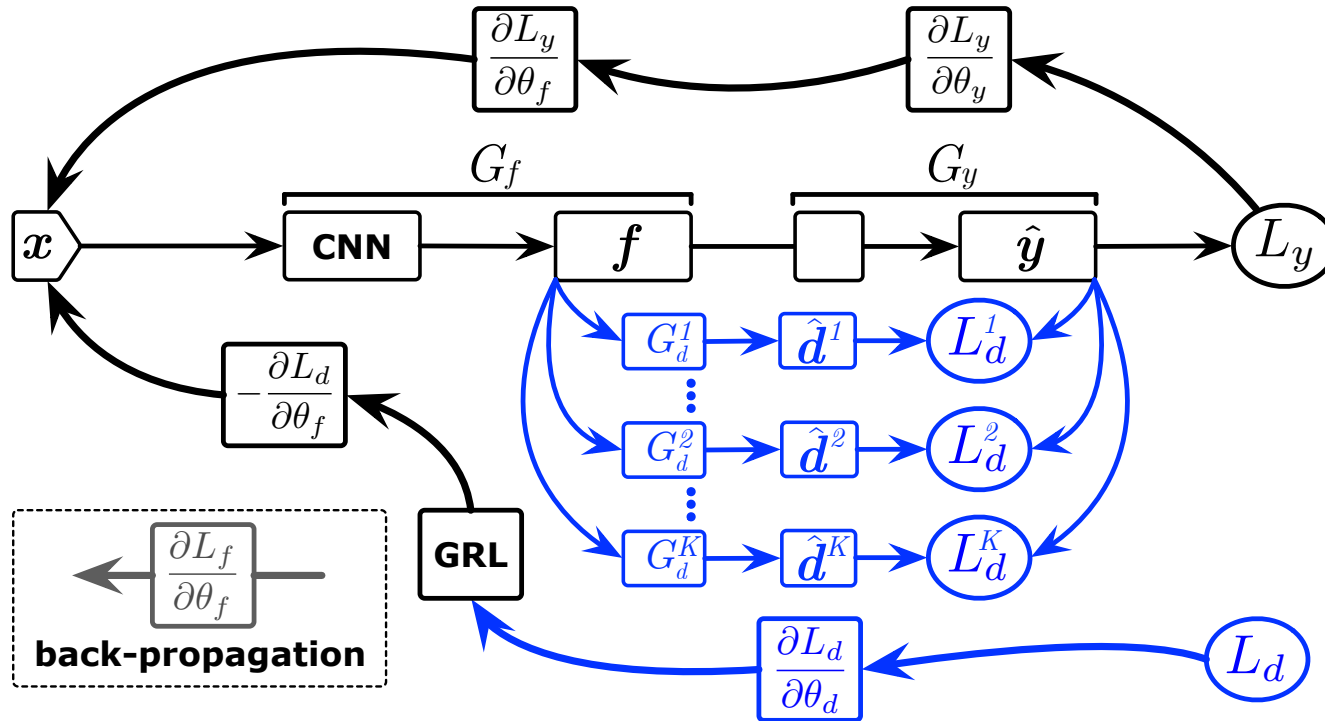
$$Y_s \neq Y_t$$

Partial Transfer Learning



$$Y_s \supset Y_t$$

Selective Adversarial Network (SAN)



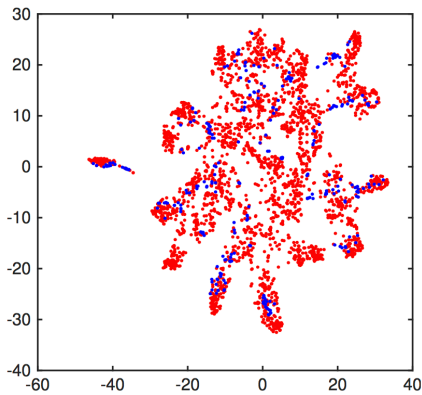
$$C(\theta_f, \theta_y, \theta_d^k |_{k=1}^{|\mathcal{C}_s|}) = \frac{1}{n_s} \sum_{\mathbf{x}_i \in \mathcal{D}_s} L_y(G_y(G_f(\mathbf{x}_i)), y_i) + \frac{1}{n_t} \sum_{\mathbf{x}_i \in \mathcal{D}_t} H(G_y(G_f(\mathbf{x}_i)))$$

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

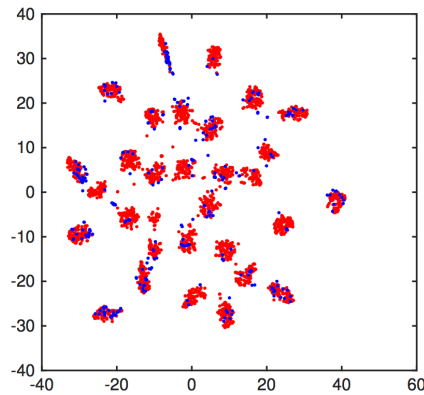
Selective Adversarial Network (SAN)



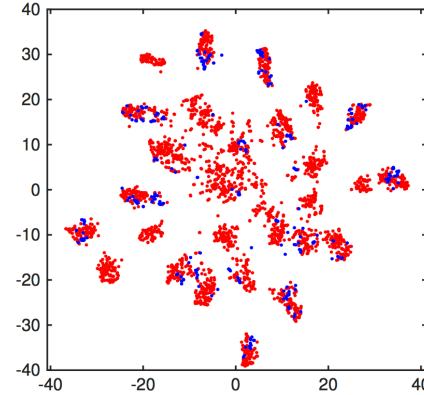
Method	Caltech-Office				ImageNet-Caltech		
	C 256 \rightarrow W 10	C 256 \rightarrow A 10	C 256 \rightarrow D 10	Avg	I 1000 \rightarrow C 84	C 256 \rightarrow I 84	Avg
AlexNet [14]	58.44	76.64	65.86	66.98	52.37	47.35	49.86
DAN [15]	42.37	70.75	47.04	53.39	54.21	52.03	53.12
RevGrad [6]	54.57	72.86	57.96	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	88.33	83.82	85.35	85.83	68.45	55.61	62.03



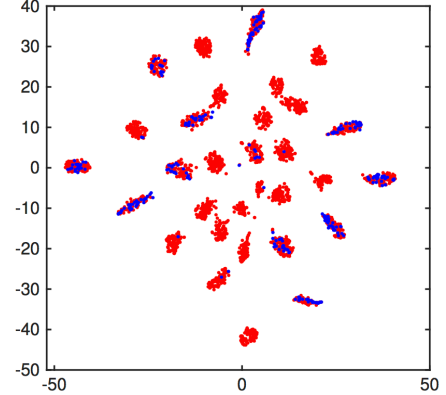
(a) DAN



(b) RevGrad

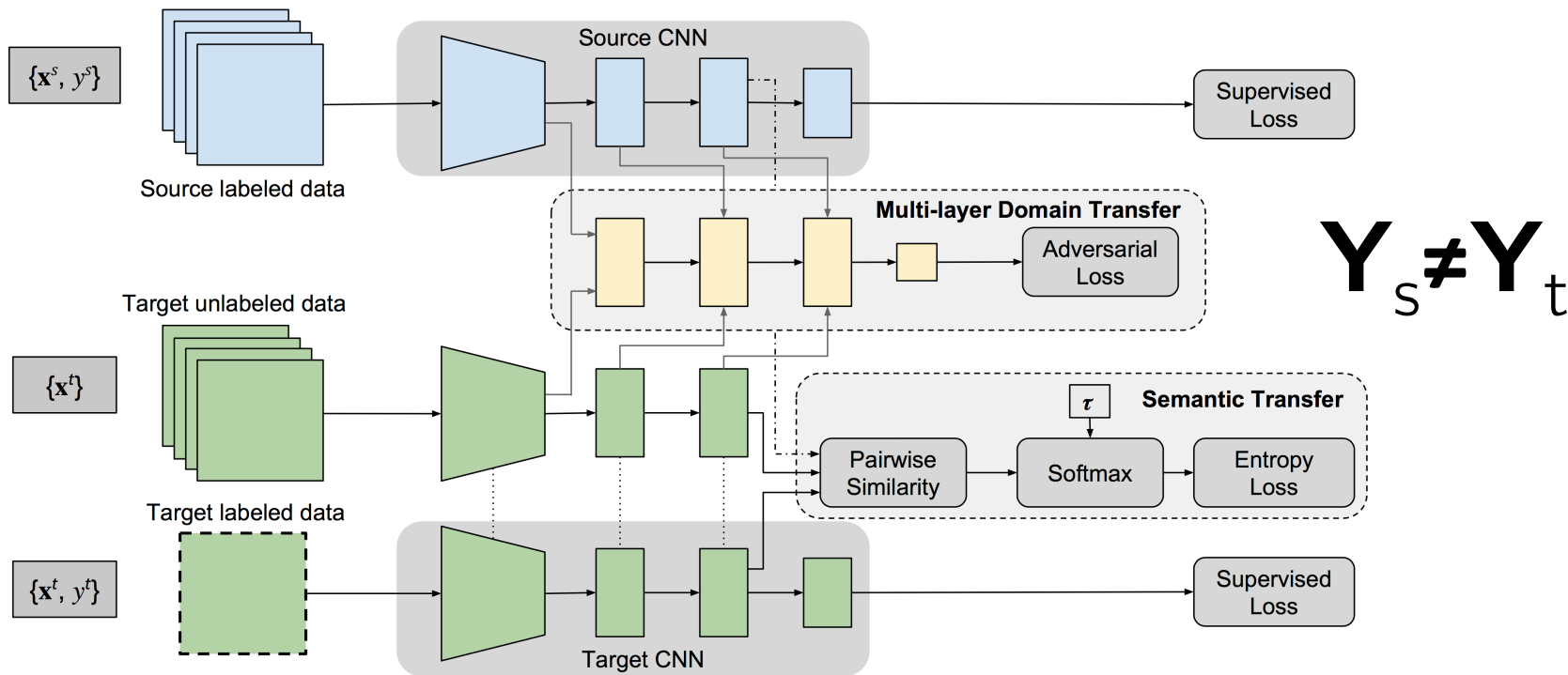


(c) RTN



(d) SAN

Joint Domain and Semantic Transfer

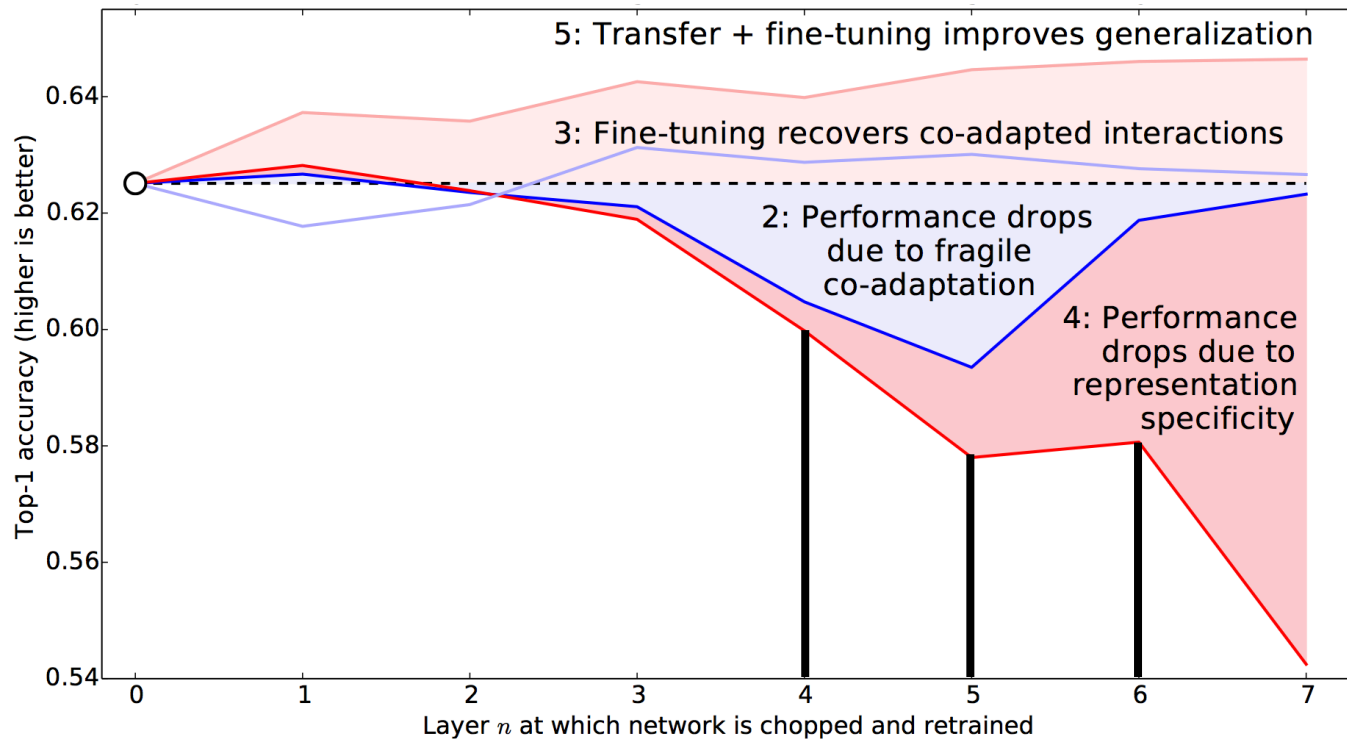
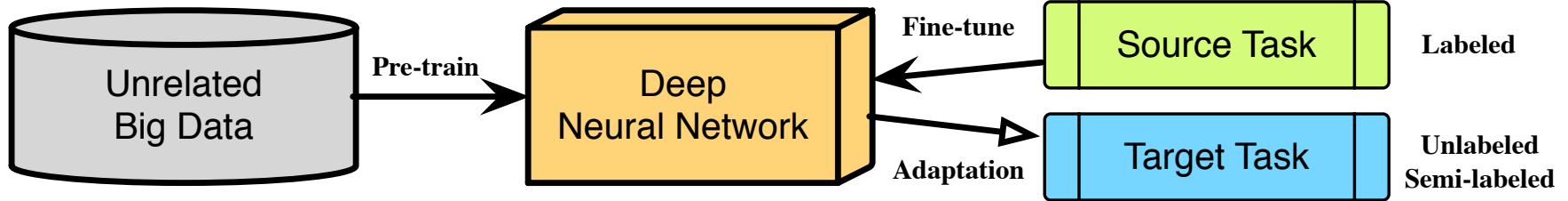


$$L_{ST}(\tilde{\mathbf{X}}_t, \mathbf{X}_s) = \sum_{\tilde{\mathbf{x}}_t \in \tilde{\mathbf{X}}_t} H(\sigma(v_s(\tilde{\mathbf{x}}_t) / \tau)) \quad L_{ST, upsup}(\tilde{\mathbf{X}}_t, \mathbf{X}_t) = \sum_{\tilde{\mathbf{x}}_t \in \tilde{\mathbf{X}}_t} H(\sigma(v_t(\tilde{\mathbf{x}}_t) / \tau))$$

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

Transferable Architecture

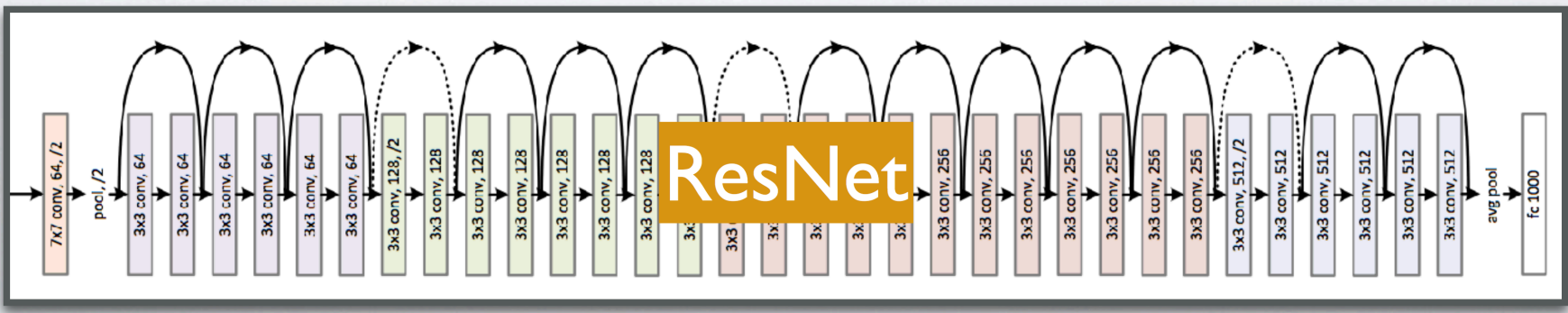
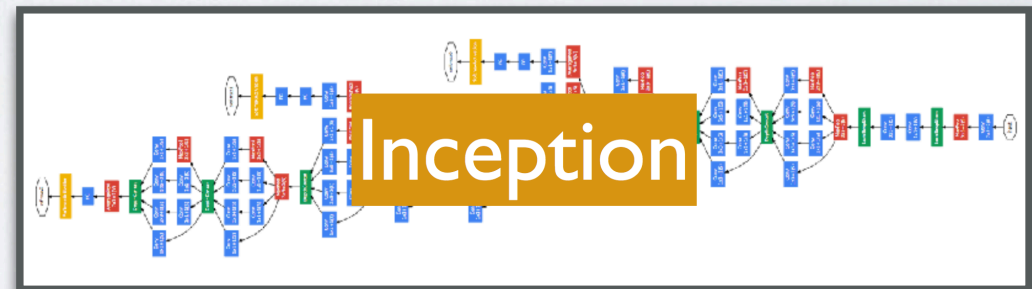
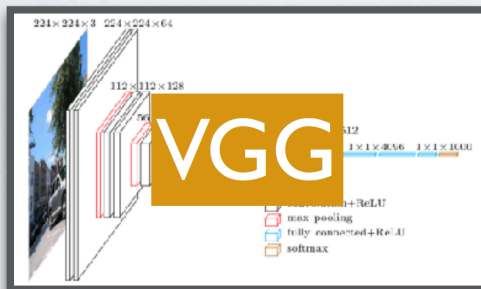
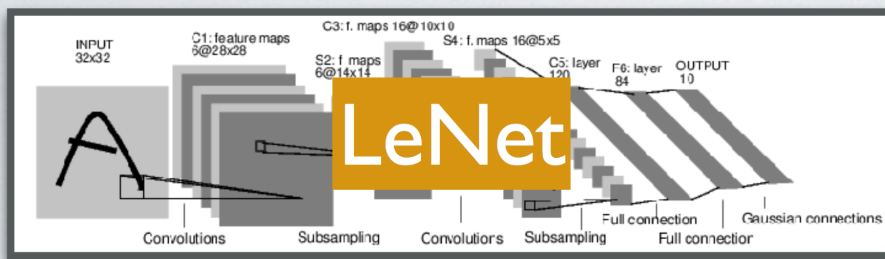
Transferability



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

$$\mathbf{X}_s \neq \mathbf{X}_t \wedge \mathbf{Y}_s \neq \mathbf{Y}_t$$

- Pixel-Level Transfer Learning

$$P(\mathbf{x}) \neq Q(\mathbf{x}) \wedge P(\mathbf{z}) \neq Q(\mathbf{z})$$

- Learning Transferable Architectures

QA



清華大學
Tsinghua University

Thank You!