# Learning Transferable Features with Deep Adaptation Networks

Mingsheng Long[12], Yue Cao[1], Jianmin Wang[1], and Michael I. Jordan[2]

International Conference on Machine Learning, 2015

## Introduction

- Goal: Enhance the transferability of features from task-specific layers
- Proposed a Deep Adaptation Network DAN architecture
  - General features can generalize well to a novel task; however, for specific features they cannot bridge the domain discrepancy
- Some ways to enhance feature transferability:
  - By mean-embedding matching, feature transferability can be enhanced substantially
  - Utilizing multi-layer representations across domains in a reproducing kernel Hilbert space

# Main Breakthrough

- *Generalizes deep CNN to the domain adaptation*
- Deep adaptation of multiple task-specific layers, including output
- Optimal adaptation using multiple kernel two-sample matching

# Deep Learning For Domain Adaptation

- None or very weak supervision in the *target* task (new domain)
  - Target classifier cannot be reliably trained due to over-fitting
  - Fine-tuning is impossible as it requires substantial supervision
- Generalize related supervised source task to the target task
  - Deep networks can learn transferable features for adaptation
- Hard to find big source task for learning deep features from scratch
  - Transfer from deep networks pre-trained on unrelated big dataset
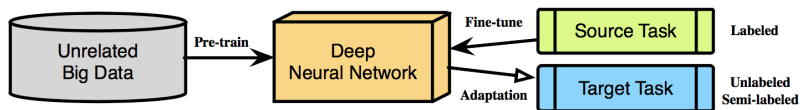  - Transferring features from distant tasks better than random features



Figure: *Deep Learning for Domain Adaptation Workflow*

# How Transferable Are Deep Features?

- Transferability is restricted by (Yosinski et al. 2014; Glorot et al. 2011)
- Specialization of higher layer neurons to original task (new task )
- Disentangling of variations in higher layers enlarges task discrepancy
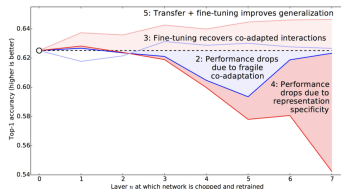- Transferability of features decreases while task discrepancy increases



Figure: Transferability of features decreases while task discrepancy increases

# Deep Adaptation Network (DAN)

## Key Observations (AlexNet) (Krizhevsky et al. 2012)

- Comprised of five convolutional layers $conv1 - conv5$ and three fully connected layers $fc6 - fc8$
- Convolutional layers learn general features: safely transferable
  - Safely freeze $conv1 - conv3$ & fine tuned $conv4 - conv5$
- Fully-connected layers fit task specificicy: $NOT$ safely transferable
  - Deeply adapt $fc6 - fc8$ using statistically optimal two-sample matching
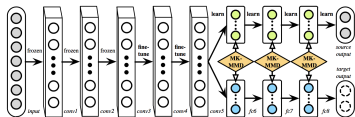


Figure: The DAN architecture for learning transferable features

# Objective Function

## Main Problems

- Feature transferability decreases with increasing task discrepancy
- Higher layers are tailored to specific tasks, NOT safely transferable
- Adaptation effect may vanish in back-propagation of deep networks

## Deep Adaptation with Optimal Matching

- Deep adaptation: match distributions in multiple layers includingoutput
- Optimal matching: maximize two-sample test of multiple kernels

$$\min_{\Theta} \max_{\kappa} \frac{1}{n_a} \sum_{i=1}^{n_a} J(\theta(x_i^a), y_i^a) + \lambda \sum_{\ell=l_1}^{l_2} d_k^2(D_s^\ell, D_t^\ell) \tag{1}$$

$\lambda > 0$ is a penality, $D_*^\ell = \{h_i^{*\ell}\}$ is the $\ell$-th layer hidden representaion

# MK-MMD

**Theorem (Two-Sample Test (Gretton et al. 2012))**

- $p = q$ iff $d_k^2(p,q) = 0$ (In practice, $d_k^2(p,q) < \epsilon$)
- $\max\limits_{k \in \kappa} d_k^2(D_s^\ell, D_t^\ell)\sigma_k^{-2} \Leftrightarrow \min$Type II Error ($d_k^2(p,q) < \epsilon$ when $p \neq q$)

**Multiple Kernel Maximum Mean Discrepancy (MK-MMD)**

$\triangleq$ RKHS distance between kernel embeddings of distributions p and q

$$d_k^2(p,q) \triangleq \|E_p[\phi(x^s)] - E_q[\phi(x^t)]\|_{\mathcal{H}_k}^2 \tag{2}$$

$k(\mathbf{x}^s, \mathbf{x}^t) = \langle \phi(\mathbf{x}^s), \phi(\mathbf{x}^t) \rangle$ is a convex combination of m PSD kernels

$$\kappa \triangleq \left\{ k = \sum_{u=1}^m \beta_u k_u : \sum_{u=1}^m \beta_u = 1, \beta_u \geq 0, \forall u \right\} \tag{3}$$

# Learning CNN

## Linear-Time Algorithm of MK-MMD (Streaming Algorithm)

$O(n^2) : d_k^2(p, q) = \mathbf{E}_{\mathbf{x}^s \mathbf{x}'^s} k(\mathbf{x}^s, \mathbf{x}'^s) + \mathbf{E}_{\mathbf{x}^t \mathbf{x}'^t} k(\mathbf{x}^t, \mathbf{x}'^t) - 2\mathbf{E}_{\mathbf{x}^s \mathbf{x}^t} k(\mathbf{x}^s, \mathbf{x}^t)$

$d_k^2(p, q) = \frac{2}{n_s} \sum_{i=1}^{\frac{n_s}{2}} g_k(\mathbf{z}_i) \rightarrow$ linear-time unbiased estimate

- Quad-tuple: $\mathbf{z}_i \triangleq (\mathbf{x}_{2i-1}^s, \mathbf{x}_{2i}^s, \mathbf{x}_{2i-1}^t, \mathbf{x}_{2i}^t)$
- $g_k(\mathbf{z}_i) \triangleq k(\mathbf{x}_{2i-1}^s, \mathbf{x}_{2i}^s) + k(\mathbf{x}_{2i-1}^t, \mathbf{x}_{2i}^t) - k(\mathbf{x}_{2i-1}^s, \mathbf{x}_{2i}^t) - k(\mathbf{x}_{2i}^s, \mathbf{x}_{2i-1}^t)$

## Stochastic Gradient Descent(SGD)

For each layer $\ell$ and for each quad-tuple $\mathbf{z}_i \triangleq (\mathbf{x}_{2i-1}^s, \mathbf{x}_{2i}^s, \mathbf{x}_{2i-1}^t, \mathbf{x}_{2i}^t)$

$$\nabla_{\Theta^\ell} = \frac{\partial J(z_i)}{\partial \Theta^\ell} + \lambda \frac{\partial g_k(z_i^\ell)}{\partial \Theta^\ell} \tag{4}$$

# Learning Kernel

Learning optimal kernel $k = \sum_{u=1}^{m} \beta_u k_u$

Maximizing test power $\triangleq$ minimizing Type II error (Gretton et al. 2012)

$$\max_{k \in \kappa} d_k^2(D_s^\ell, D_t^\ell)\sigma_k^{-2} \qquad (5)$$

where $\sigma_k^2 = \mathbf{E_z}g_k^2(\mathbf{z}) - [\mathbf{E}_z g_k(\mathbf{z})]^2$ is the estimation variance.

Quadratic Program (QP), scaling linearly to sample size: $\mathrm{O}(m^{2n} + m^3)$

$$\min_{d^T\beta=1, \beta \geq 0} \beta^T(Q + \epsilon I)\beta \qquad (6)$$

where $\mathbf{d} = (d_1, d_2, ..., d_m)^T$, and each $d_u$ is MMD using base kernel $k_u$.

# Analysis

## Theorem (Adaptation Bound)

(Ben-David et al. 2010) Let $\theta \in H$ be a hypothesis, $\epsilon_s(\theta)$ and $\epsilon_t(\theta)$ be the expected risks of source and target respectively, then

$$\epsilon_t(\theta) \leq \epsilon_s(\theta) + 2d_k(p,q) + C \tag{7}$$

where $C$ is a constant for the complexity of hypothesis space, the empirical estimate of **H**-divergence, and the risk of an ideal hypothesis for both tasks.

## Two-Sample Classifier: Nonparametric vs. Parametric

- Nonparametric MMD directly approximates $d_{\mathcal{H}}(p,q)$
- Parametric classifier: adversarial training to approximate $d_{\mathcal{H}}(p,q)$

# Experiment Setup

- Datasets: pre-trained on ImageNet, fined-tuned on Office&Caltech
- Tasks: 12 adaptation tasks $\Rightarrow$ An unbiased look at dataset bias
- Variants: DAN; single-layer: $DAN_7$, $DAN_8$; single-kernel: $DAN_{SK}$
- Protocols: unsupervised adaptation vs semi-supervised adaptation
- Parameter selection: cross-validation by jointly assessing
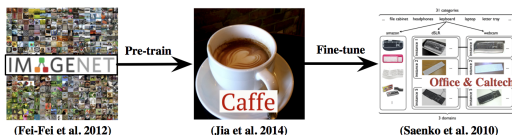  - test errors of source classifier and two-sample classifier (MK-MMD)



Figure: *the proposed DAN model is trained by fine-tuning from the AlexNet model (Krizhevsky et al., 2012) pre-trained on ImageNet, implemented in Caffe.*

# Results & Discussion

**Learning transferable features by deep adaptation and optimal matching**

- Deep adaptation of multiple domain-specific layers (DAN) vs. shallow adaptation of one hard-to-tweak layer (DDC)
- Two samples can be matched better by MK-MMD vs. SK-MMD

*Table 1.* Accuracy on *Office-31* dataset with standard unsupervised adaptation protocol (Gong et al., 2013).

| Method | A → W | D → W | W → D | A → D | D → A | W → A | Average |
|---|---|---|---|---|---|---|---|
| TCA | $21.5 \pm 0.0$ | $50.1 \pm 0.0$ | $58.4 \pm 0.0$ | $11.4 \pm 0.0$ | $8.0 \pm 0.0$ | $14.6 \pm 0.0$ | 27.3 |
| GFK | $19.7 \pm 0.0$ | $49.7 \pm 0.0$ | $63.1 \pm 0.0$ | $10.6 \pm 0.0$ | $7.9 \pm 0.0$ | $15.8 \pm 0.0$ | 27.8 |
| CNN | $61.6 \pm 0.5$ | $95.4 \pm 0.3$ | $\underline{99.0} \pm 0.2$ | $63.8 \pm 0.5$ | $51.1 \pm 0.6$ | $49.8 \pm 0.4$ | 70.1 |
| LapCNN | $60.4 \pm 0.3$ | $94.7 \pm 0.5$ | $\mathbf{99.1} \pm 0.2$ | $63.1 \pm 0.6$ | $51.6 \pm 0.4$ | $48.2 \pm 0.5$ | 69.5 |
| DDC | $61.8 \pm 0.4$ | $95.0 \pm 0.5$ | $98.5 \pm 0.4$ | $64.4 \pm 0.3$ | $52.1 \pm 0.8$ | $\underline{52.2} \pm 0.4$ | 70.6 |
| $DAN_7$ | $63.2 \pm 0.2$ | $94.8 \pm 0.4$ | $98.9 \pm 0.3$ | $65.2 \pm 0.4$ | $52.3 \pm 0.4$ | $52.1 \pm 0.4$ | 71.1 |
| $DAN_8$ | $\underline{63.8} \pm 0.4$ | $94.6 \pm 0.5$ | $98.8 \pm 0.6$ | $65.8 \pm 0.4$ | $52.8 \pm 0.4$ | $51.9 \pm 0.5$ | 71.3 |
| $DAN_{SK}$ | $63.3 \pm 0.3$ | $\underline{95.6} \pm 0.2$ | $99.0 \pm 0.4$ | $\underline{65.9} \pm 0.7$ | $\underline{53.2} \pm 0.5$ | $52.1 \pm 0.4$ | $\underline{71.5}$ |
| DAN | $\mathbf{68.5} \pm 0.4$ | $\mathbf{96.0} \pm 0.3$ | $\underline{99.0} \pm 0.2$ | $\mathbf{67.0} \pm 0.4$ | $\mathbf{54.0} \pm 0.4$ | $\mathbf{53.1} \pm 0.3$ | **72.9** |

Figure: Table 1. Accuracy on Office-31 dataset with standard unsupervised adaptation protocol

# Results & Discussion

Semi-supervised adaptation: source supervision vs. target supervision?

- Limited target supervision is prone to over-fitting the target task
- Source supervision can provide strong but inaccurate inductive bias
- Two-sample matching is more effective for bridging dissimilar tasks

Table 2. Accuracy on *Office-10* + *Caltech-10* dataset with standard unsupervised adaptation protocol (Gong et al., 2013).

| Method | A → C | W → C | D → C | C → A | C → W | C → D | Average |
|--------|-------|-------|-------|-------|-------|-------|---------|
| TCA | $42.7 \pm 0.0$ | $34.1 \pm 0.0$ | $35.4 \pm 0.0$ | $54.7 \pm 0.0$ | $50.5 \pm 0.0$ | $50.3 \pm 0.0$ | 44.6 |
| GFK | $41.4 \pm 0.0$ | $26.4 \pm 0.0$ | $36.4 \pm 0.0$ | $56.2 \pm 0.0$ | $43.7 \pm 0.0$ | $42.0 \pm 0.0$ | 41.0 |
| CNN | $83.8 \pm 0.3$ | $76.1 \pm 0.5$ | $80.8 \pm 0.4$ | $91.1 \pm 0.2$ | $83.1 \pm 0.3$ | $89.0 \pm 0.3$ | 84.0 |
| LapCNN | $83.6 \pm 0.6$ | $77.8 \pm 0.5$ | $80.6 \pm 0.4$ | $\mathbf{92.1} \pm 0.3$ | $81.6 \pm 0.4$ | $87.8 \pm 0.4$ | 83.9 |
| DDC | $84.3 \pm 0.5$ | $76.9 \pm 0.4$ | $80.5 \pm 0.2$ | $91.3 \pm 0.3$ | $85.5 \pm 0.3$ | $89.1 \pm 0.3$ | 84.6 |
| DAN$_7$ | $\underline{84.7} \pm 0.3$ | $78.2 \pm 0.5$ | $\underline{81.8} \pm 0.3$ | $91.6 \pm 0.4$ | $87.4 \pm 0.3$ | $88.9 \pm 0.5$ | 85.4 |
| DAN$_8$ | $84.4 \pm 0.3$ | $\underline{80.8} \pm 0.4$ | $81.7 \pm 0.2$ | $91.7 \pm 0.3$ | $\underline{90.5} \pm 0.4$ | $89.1 \pm 0.4$ | $\underline{86.4}$ |
| DAN$_{SK}$ | $84.1 \pm 0.4$ | $79.9 \pm 0.4$ | $81.1 \pm 0.5$ | $91.4 \pm 0.3$ | $86.9 \pm 0.5$ | $\underline{89.5} \pm 0.3$ | 85.5 |
| DAN | $\mathbf{86.0} \pm 0.5$ | $\mathbf{81.5} \pm 0.3$ | $\mathbf{82.0} \pm 0.4$ | $\underline{92.0} \pm 0.3$ | $\mathbf{92.0} \pm 0.4$ | $\mathbf{90.5} \pm 0.2$ | $\mathbf{87.3}$ |

Figure: Table 2. Accuracy on Office-10 + Caltech-10 dataset with standard unsupervised adaptation protocol

# Data Visualization

**How transferable are DAN features? t-SNE embedding for visualization**

- target points form clearer class boundaries
- target points can be classified more accurately
- Source and target categories are aligned better



(a) DDC Features on Source    (b) DDC Features on Target    (c) DAN Features on Source    (d) DAN Features on Target
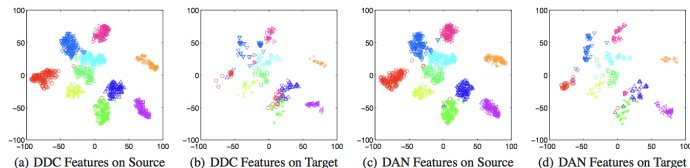
Figure: t-SNE of DDC features on source (a) and target (b) versus. t-SNE of DAN features on source (c) and target (d)

# Empirical Analysis

How is generalization performance related to two-sample discrepancy?

- $\hat{d}_A$ on CNN & DAN features > $\hat{d}_A$ on Raw features
- $\Rightarrow \hat{d}_A$ on DAN feature is much smaller than $\hat{d}_A$ on CNN feature
- $\hat{d}_A$ on DAN feature < $\hat{d}_A$ on CNN feature
- $\Rightarrow$ Domain adaptation can be boosted by reducing domain discrepancy
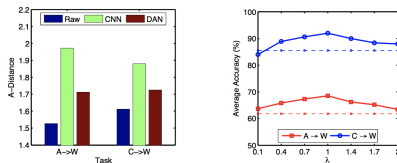


Figure: (e) A-Distance of CNN & DAN features; (f) sensitivity of $\lambda$

# Summary

- DAN: A deep adaptation network for learning transferable features
- Deep adaptation of multiple task-specific layers (including output)
- Optimal adaptation using multiple kernel two-sample matching
- Optimal adaptation using multiple kernel two-sample matching